

JMIR Public Health and Surveillance

Impact Factor (2020): 4.11
Volume 7 (2021), Issue 8 ISSN: 2369-2960 Editor in Chief: Travis Sanchez, PhD, MPH

Contents

Viewpoints

- Substandard and Falsified Medicines: Proposed Methods for Case Finding and Sentinel Surveillance ([e29309](#))
Elizabeth Pisani, Amalia Hasnida, Mawaddati Rahmi, Maarten Kok, Steven Harsono, Yusi Anggriani. 3
- With Great Hopes Come Great Expectations: Access and Adoption Issues Associated With COVID-19 Vaccines ([e26111](#))
Zhaohui Su, Dean McDonnell, Ali Cheshmehzangi, Xiaoshan Li, Daniel Maestro, Sabina Šegalo, Junaid Ahmad, Xiaoning Hao. 97

Original Papers

- Identifying Communities at Risk for COVID-19–Related Burden Across 500 US Cities and Within New York City: Unsupervised Learning of the Coprevalence of Health Indicators ([e26604](#))
Andrew Deonarine, Genevieve Lyons, Chirag Lakhani, Walter De Brouwer. 18
- Online News Coverage of the Sugar-Sweetened Beverages Tax in Malaysia: Content Analysis ([e24523](#))
Muhammad Mohd Hanim, Budi Md Sabri, Norashikin Yusof. 34
- Forecasting COVID-19 Hospital Census: A Multivariate Time-Series Model Based on Local Infection Incidence ([e28195](#))
Hieu Nguyen, Philip Turk, Andrew McWilliams. 45
- Census Tract Patterns and Contextual Social Determinants of Health Associated With COVID-19 in a Hispanic Population From South Texas: A Spatiotemporal Perspective ([e29205](#))
Cici Bauer, Kehe Zhang, Miryoung Lee, Susan Fisher-Hoch, Esmeralda Guajardo, Joseph McCormick, Isela de la Cerda, Maria Fernandez, Belinda Reininger. 58
- Natural Language Processing Insight into LGBTQ+ Youth Mental Health During the COVID-19 Pandemic: Longitudinal Content Analysis of Anxiety-Provoking Topics and Trends in Emotion in LGBTQ+ Microcommunity Subreddit ([e29029](#))
Hannah Stevens, Irena Acic, Sofia Rhea. 69
- The Roles of General Health and COVID-19 Proximity in Contact Tracing App Usage: Cross-sectional Survey Study ([e27892](#))
Dirk Witteveen, Pablo de Pedraza. 84

Corrigenda and Addenda

Correction: Census Tract Patterns and Contextual Social Determinants of Health Associated With COVID-19 in a Hispanic Population From South Texas: A Spatiotemporal Perspective ([e32870](#))

Cici Bauer, Kehe Zhang, Miryoung Lee, Susan Fisher-Hoch, Esmeralda Guajardo, Joseph McCormick, Isela de la Cerda, Maria Fernandez, Belinda Reininger.

43

Viewpoint

Substandard and Falsified Medicines: Proposed Methods for Case Finding and Sentinel Surveillance

Elizabeth Pisani^{1,2,3}, MA, MSc, PhD; Amalia Hasnida¹, BSc, MSc; Mawaddati Rahmi³, MSc; Maarten Olivier Kok^{1,4}, BA, MSc, PhD; Steven Harsono⁵, BA; Yusi Anggriani³, BSc, MPH, PhD

¹Erasmus School of Health Policy and Management, Erasmus University, Rotterdam, Netherlands

²School of Public Health, Imperial College, London, United Kingdom

³Faculty of Pharmacy, Universitas Pancasila, Jakarta, Indonesia

⁴Department of Health Sciences, Vrije Universiteit Amsterdam, Amsterdam, Netherlands

⁵IQVIA Public Health, Singapore, Singapore

Corresponding Author:

Elizabeth Pisani, MA, MSc, PhD

Erasmus School of Health Policy and Management

Erasmus University

3000 DR

Rotterdam

Netherlands

Phone: 31 2072541654

Email: pisani@ternyata.org

Abstract

The World Health Organization and others warn that substandard and falsified medicines harm health and waste money, especially in low- and middle-income countries. However, no country has measured the market-wide extent of the problem, and no standardized methods exist to estimate the prevalence of either substandard or falsified medicines. This is, in part, because the task seems overwhelming; medicine markets are huge and diverse, and testing medicines is expensive. Many countries do operate some form of postmarket surveillance of medicine, but their methods and goals differ. There is currently no clear guidance on which surveillance method is most appropriate to meet specific public health goals. In this viewpoint, we aimed to discuss the utility of both case finding and risk-based sentinel surveillance for substandard and falsified medicines, linking each to specific public health goals. We posit that choosing the system most appropriate to the goal, as well as implementing it with a clear understanding of the factors driving the production and sale of substandard and falsified medicines, will allow for surveillance resources to be concentrated most efficiently. We adapted principles used for disease outbreak responses to suggest a case-finding system that uses secondary data to flag poor-quality medicines, proposing risk-based indicators that differ for substandard and falsified medicines. This system potentially offers a cost-effective way of identifying “cases” for market withdrawal, enhanced oversight, or another immediate response. We further proposed a risk-based sentinel surveillance system that concentrates resources on measuring the prevalence of substandard and falsified medicines in the risk clusters where they are most likely to be found. The sentinel surveillance system provides base data for a transparent, spreadsheet-based model for estimating the national prevalence of substandard and falsified medicines. The methods we proposed are based on ongoing work in Indonesia, a large and diverse middle-income country currently aiming to achieve universal health coverage. Both the case finding and the sentinel surveillance system are designed to be adaptable to other resource-constrained settings.

(*JMIR Public Health Surveill* 2021;7(8):e29309) doi:[10.2196/29309](https://doi.org/10.2196/29309)

KEYWORDS

substandard drugs; falsified medicine; counterfeit medicine; medicine quality; sentinel surveillance; public health surveillance; substandard; pharmaceuticals; surveillance; public health

Introduction

Background on Substandard and Falsified Medicines

In late 2017, World Health Organization's (WHO's) press department issued a press release with the bold headline: "1 in 10 medical products in developing countries is substandard or falsified" [1]. More recently, with governments scrambling to secure supplies of diagnostic tests, medicines, and vaccines to cope with the COVID-19 pandemic, WHO and others have issued new warnings stating that the world may face an increased threat of poor-quality medicines [2-5]. These include substandard medicines, which are made by registered pharmaceutical companies in regulated factories but do not meet the quality standards set out in their market authorization paperwork, either because they were poorly made or because they have degraded since manufacture. This increased threat of poor-quality medicines also includes falsified medicines, which are made, repackaged, or sold by criminals who seek deliberately to misrepresent the identity, composition, or source of the product [6].

Poor-quality medicines can use up family and national budgets without curing patients; indeed, they sometimes poison or kill people instead of curing them. Underdosing infectious pathogens also allows drug-resistant infections to spread [7,8]. Thus, if these "medicines" are indeed common, they may substantially undermine physical and financial health. Estimates based on available data for particularly well-studied molecules provide an order of magnitude: poor-quality antimalarials were estimated to cost US \$130 million per year in a single region of the Democratic Republic of Congo, US \$141.5 million in Zambia, and US \$830 million across Nigeria. In the latter country, substandard antimalarials are estimated to contribute to 12,300 deaths per year [9-11]. A 2015 study reported that, in 39 sub-Saharan African countries, there were 122,350 deaths attributable to poor-quality antimalarials among children under 5 years of age. However, the authors noted that "there is considerable uncertainty surrounding our results because of gaps in data on case fatality rates and prevalence of poor-quality antimalarials" [12]. Similarly, the meta-analysis that gave rise to WHO's press release, which said that 10% of medical products in developing countries are substandard or falsified, is careful to note the many limitations of that estimate. This meta-analysis was based on studies of uneven sizes and methods, conducted largely in low-income countries with limited domestic pharmaceutical industries, and heavily skewed toward antimalarials and a few other medicines that most interest global health agencies. Even within that constrained pool and looking only at studies that included sample sizes of 50 or more, reported prevalence of substandard or falsified medicines ranged from 0% to 91% [13]. Reviews have reported similar data constraints and findings [14-18]. For example, Ozawa and colleagues [18] found that studies reported a prevalence of substandard or falsified medicines between 0.8% and 89% in Africa and a prevalence of between 0.7% and 50% in Asia.

WHO actively maintains a case-reporting system for substandard and falsified medical products, including medicines, contraceptives, vaccines, and point-of-sale diagnostics. For

brevity, we use the term medicines throughout this paper to cover all these medical products. Regular training provided to individuals designated as in-country focal points increases the use of the system, but, similar to all case-reporting systems, it provides no information on denominators (the number of products inspected or tested), so interpretation of trends and comparisons between countries is difficult.

It may be that the problem of medicine quality is understated because of a vicious cycle of limited systematic measurement leading to limited visibility and limited awareness of the problem that in turn restricts resources available for systematic measurement. Alternatively, the problem may be that WHO and researchers are cherry-picking data to overstate the problem, perhaps for reasons of self-interest, as Hodges and Garnett [19] suggest.

We do not know which of these dynamics holds true. There is, to our knowledge, no clear understanding of the prevalence of substandard or falsified medicines in any single country, let alone across all the "developing countries," as suggested by the press release's headline. No country has yet made systematic estimates of the prevalence of substandard or falsified medicines across all therapeutic categories in its medicine market, and no standardized methods for calculating such an estimation yet exist.

In this viewpoint, we aimed to briefly review different approaches to surveillance and estimation in public health, discuss their relevance in the context of medicine quality, and lay out ideas for 2 potentially cost-minimizing methods that may improve our ability to measure or reduce the prevalence of poor-quality medicines, especially, in low- and middle-income settings.

Approaches to Surveillance

Overview of Surveillance Systems

We follow WHO, United States Centers for Disease Control and Prevention, the World Bank, and others in defining public health surveillance as the ongoing and systematic collection and use of data to inform policy, plan and evaluate interventions, and improve health outcomes [20,21]. Surveillance systems monitor the prevalence of infectious and noncommunicable diseases; of disability; and, increasingly, of the behavioral, social, corporate, and environmental causes of ill-health. In addition, surveillance systems have, in recent decades, expanded to include the systematic monitoring of health system factors such as service use, prescription practices, or access to medicine.

Surveillance can take many forms, each serving a slightly different purpose within the catch-all definition of "improving health outcomes." However, most can be categorized into either "passive" or "active" surveillance. Passive surveillance involves reporting events such as disease diagnoses as they arise. An early example of passive sentinel surveillance in the United States was the weekly reporting of diseases by designated physicians, which began in Massachusetts in 1874 [22]. The informatics era has greatly expanded the potential for secondary data to be used to inform public health decision-making. Examples include the use of both retail data of over-the-counter medicine sales and data from internet searches to flag potential

disease outbreaks, and the use of medical-claims data to track trends in noncommunicable diseases [23,24]. Active surveillance tends to be more resource intensive, usually involving purposive data collection—often the collection and screening of blood or other biological samples or, more recently, medical imaging.

Active surveillance systems systematically collect and test samples for the purpose of tracking ill-health or health-related risks. Some active surveillance, such as active sentinel surveillance, test a cross-section of a defined population to establish disease prevalence. Others, such as case finding, specifically, target individuals at highest risk of needing

services. These 2 types of active surveillance have different purposes. Sentinel surveillance, similar to many other surveillance systems, such as those that track noncommunicable diseases, determinants of health, and health system factors, provides data intended to guide medium- or longer-term health program planning. Case-finding systems, frequently used in infectious disease outbreaks but also used for early detection of treatable noncommunicable conditions, provide data intended to inform immediate therapeutic or preventative action. These different goals, upon which this paper focuses in the context of medicine quality, affect the design and use of surveillance systems and data, as show in [Table 1](#).

Table 1. Major types of surveillance systems in public health.

Purpose	Outbreak response	Health program planning
System design:	Case finding: identify infected individuals	Sentinel surveillance: track prevalence over time
Resulting action:	Isolate and treat	Adjust policies and programs
Key characteristic:	Specific: pinpoint individuals for rapid follow-up	Comparable: standardized methods allowing comparison over time
Cannot be used to:	Estimate prevalence; track trends over time	Respond at individual level

Case Finding

The control of outbreaks and epidemics of infectious diseases requires that chains of transmission be broken. In these circumstances, surveillance systems try to identify infected individuals, isolating and, if possible, treating them to interrupt transmission. We term these systems “case finding.” They are relatively rare but have seen a resurgence during the COVID-19 pandemic.

Sentinel Surveillance

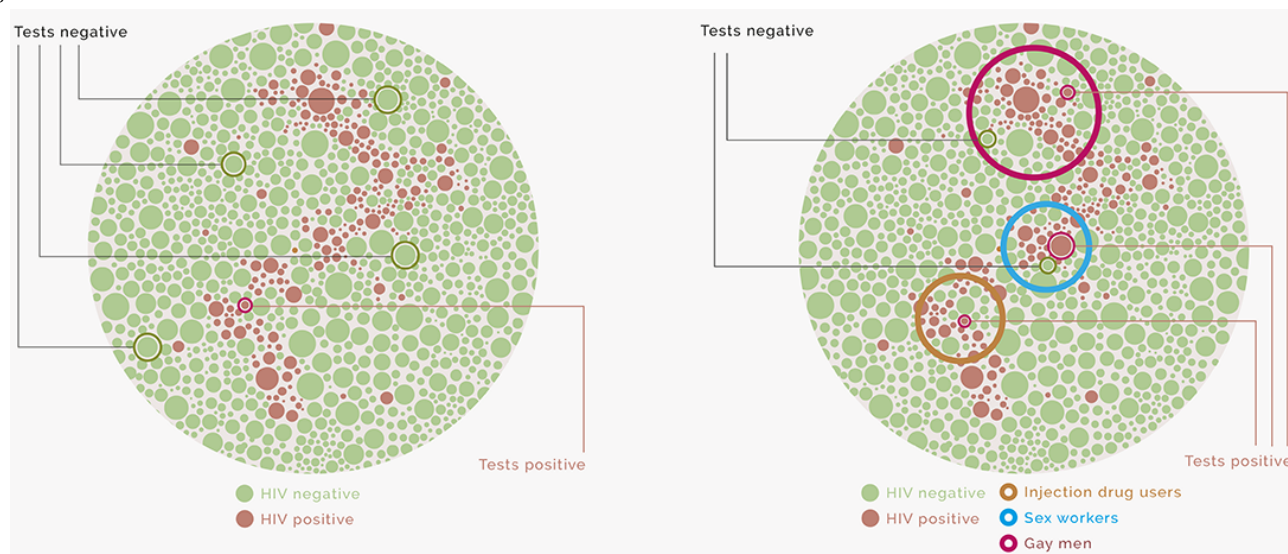
Sentinel surveillance systems are more common. Designed to track trends in infection over time, they use standardized methods to measure the prevalence of a disease within a defined population, comparing the result with prevalence measured in the same way in earlier years or in different locations. Sentinel surveillance is used to estimate the burden of disease, to target prevention and treatment interventions, and to monitor the impact of these interventions.

While passive surveillance can achieve these goals, it is of limited use for tracking rare diseases, which are easily missed by these systems [25]. While epidemiological orthodoxy holds that active surveillance involving regular screening of randomly selected samples provides the best approximation of disease trends across a population as a whole, this is also impractical for rare conditions that would require very large samples.

The HIV pandemic entrenched the idea of active sentinel surveillance in populations defined not by geography but by risk of exposure to the virus [26,27]. This allowed health authorities to focus surveillance resources in subpopulations where the majority of cases of the largely invisible disease were to be found while still producing comparable data and tracking trends over time. In many countries, those groups included people who inject drugs, gay men, sex workers of all genders, and sex workers’ most frequent clients.

[Figure 1](#) illustrates randomized and risk-based approaches to HIV sentinel surveillance. For the same limited resources (in this simplification, 5 tests), random sampling, on the left, yields just 1 positive test, while sentinel surveillance, on the right, yields 3. Combined with robust estimates of the size of those subpopulations, this approach will provide a more accurate estimate of the prevalence of infection nationally, the data produced will more accurately reflect the effect of targeted risk-reduction interventions, and these benefits will be achieved at a lower cost compared with random sampling.

The HIV example is of considerable relevance when thinking about surveillance of substandard and falsified medicines, because it is largely invisible until actively tested and clusters around known risk factors. A similar model for medicine quality is presented in this paper’s section “Proposed Method for Sentinel Surveillance.”

Figure 1. Illustrative difference between random and risk-based surveillance for HIV.

Existing Surveillance of Medicine Quality

At the international level, surveillance of medicine quality takes the form of passive reporting of detected cases. As with disease case reporting, this provides information by demographic, geographic, environmental, or other factors, which is invaluable in helping to identify clusters of risk. However, case reporting does not provide any idea of the number of products tested. No case reports may mean there is no problem in a particular country, but it may also mean there is no capacity or willingness to detect or report cases.

At the national level, medicine regulators in many countries conduct some form of postmarket surveillance. In some countries, this is largely passive, limited to collating reports of adverse events submitted by health care providers through pharmacovigilance systems. Again, this means the denominator is unknown. However, in other countries, the regulator actively samples medicines from supply chains for inspection and testing. Where they report the number of products inspected or tested, as well as the number of out-of-specification products, this active surveillance allows for the calculation of prevalence in the segment of the market from which samples were drawn.

Sample selection in active surveillance varies widely, from random to convenience sampling, although regulators do not always state which method they use. Academic groups and WHO have published recommended methods for conducting surveys of medicine quality [28,29], as well as for sampling high-risk medicines from the internet [30]. While not focused specifically on sentinel surveillance, these methods have informed the guidance provided by technical and regulatory agencies on surveillance approaches that focus on selecting products at the highest risk for inspection and testing, including at the point of import [31-36]. To date, the criteria for determining risk have focused largely on the risk of impact to public health, factors intrinsic to the molecules (eg, stability and therapeutic index), and regulatory history. Not all agencies share information about risk profiling, for fear of helping those who produce poor-quality medicines to circumvent targets. However, as far as we know, market-related drivers of the risk

of falsification are rarely considered. Furthermore, none of the guidelines or tools currently in the public domain explicitly differentiate between the risks for falsification and the risks for substandard production or degradation. However, attention paid to risk-based, postmarket surveillance is growing; the WHO Member State Mechanism on falsified and substandard medicines chose the development of methods and tools for risk-based surveillance as a prioritized activity in its current workplan [6], and work is ongoing.

Sentinel surveillance of substandard or falsified medicines is a form of postmarket surveillance designed explicitly to select samples in reproducible ways over time, so that trends can reliably be measured. This is rare in the case of substandard or falsified medicines, although some repeat random surveys have been conducted [37]. Site-based sentinel sampling has also been attempted in some locations. This may suffer from bias if people, including falsifiers, become aware of the practice and change behavior to avoid supplying known sentinel sites [29].

In the context of medicine quality, case-finding efforts focus on trying to identify individual products that are most likely to be substandard or falsified, so that they can quickly be recalled or otherwise removed from the market. In some high-income countries with strong pharmacovigilance systems, these efforts coexist with active sampling from the supply chain.

Market-wide estimates of the prevalence of substandard or falsified medicines are virtually nonexistent. Reasons for this include the apparent complexity of the task and the expense of pharmacopeial testing. Even small countries will typically have many thousands of registered medicines and vaccines on the market. Meanwhile, well-staffed medicine-testing laboratories are scarce—there are fewer than 50 WHO-prequalified drug-testing laboratories across all low- and middle-income settings [38]. Local pharmaceutical reference standards and reagents that allow for testing of the content and quality of medicines are often unavailable, while international, gold-standard products can cost several hundred dollars for even the most common molecules [39]. In addition, medicine regulators may be wary of systematic approaches, seeing

transparent surveillance and robust estimation processes as an unwelcome evaluation of regulatory performance.

In short, while postmarket surveillance exists in different forms, there is currently no global guidance on the purpose or shape of national surveillance systems for substandard or falsified medicines and no standardized methods for translating the results of surveillance into market-wide estimates of prevalence.

The remainder of this paper proposes candidate methods, expanding on existing risk-based approaches. We aimed to review the risk factors that underlie (1) substandard and (2) falsified medicines; to propose a method for case finding based on the identified risk factors; to propose a sentinel surveillance method based on the identified risk factors; and to propose a method for developing nation-wide estimates for the prevalence of substandard and falsified medicines, based on sentinel surveillance. Our proposal is based on exploratory work undertaken in Indonesia. We believe the proposed methods are feasible in many resource-limited settings.

Table 2. Market risk factors for substandard medicines.

Risk factor	Risks to quality
High or rising pressure on profit margins	Incentivizes cost cutting
Stretched technical capacity	Increases risk of production errors or degradation during distribution
Limited oversight	Allows substandard products to flow through the supply chain
Low risk of damage to corporate reputation	Reduces incentive to invest in quality assurance

Some of these market factors operate at the level of a particular brand, others operate at the company level, others relate to the level of the supply chain, and others relate to a specific molecule. These market factors interact and further combine with other factors already considered in risk-based surveillance for medicine quality, such as the stability of a molecule or the complexity of the production or packaging process, to signal the likelihood that a medicine will be substandard or falsified.

Table 3. Market risk factors for falsified medicines.

Risk factor	Falsifier incentive
Shortage of (or restricted access to) affordable, desired product	Criminals prefer to make and sell products where there is a ready market (where demand exceeds accessible supply)
High-priced medicine or relatively high-priced brand	Profit opportunity influences choice of product and brand falsification
Limited risk of discovery or punishment	Risk of retribution shapes choice of distribution channel

We propose adding indicators of the market factors shown in [Tables 2](#) and [3](#) to increase the specificity of existing risk-based sampling and to more easily distinguish between products at risk for falsification and those more likely to be substandard.

Proposed Method for Case Finding: Sample Based on an Index of Risk

Effective case-finding systems may appeal to regulators, politicians, and the public, because they inform product recalls and other immediate actions to protect patients. On the downside, these systems are data-hungry, and sampling is relatively resource intensive. They do not systematically test a

Risk Factors for Poor-Quality Medicines

Overview of Risk Factors for Poor-Quality Medicines

Both falsified and substandard medicines exist because there is money to be made selling them. In the same way that a spike in opportunistic infections once signalled a potential cluster of undetected HIV infections, dynamics in medicine markets can act as crude predictors of clustered cases of substandard or falsified medicines. In earlier works, we reviewed available academic literature; examined reports of the case-reporting database, WHO Global Surveillance and Monitoring System for substandard and falsified medical products; and conducted detailed case studies in 4 middle-income countries, using an epidemiological approach to identify risk factors associated with substandard and falsified medicines [\[6,40\]](#). We identified a limited number of market-related factors that combine to increase the possibility that certain products in a market will be substandard, as shown in [Table 2](#).

A comprehensive review of academic literature describing interventions to control falsified medicines found few studies that addressed market drivers of falsification [\[41\]](#). However, we find market-related factors strongly shape incentives for falsifiers, leading to increased risk of falsification, as shown in [Table 3](#).

specific number of samples from a well-defined population, and, thus, cannot easily be used to measure trends over time or to estimate the magnitude of the problem.

These limitations notwithstanding, many sources of routinely collected data related to medicine markets do exist, including in middle- and some lower-income countries. These include data collected by medicine regulators, health authorities and insurers, customs and excise departments, and market research firms. We propose to use these data to guide case finding, as shown in [Textbox 1](#). Steps 2-5 should be undertaken separately for substandard and falsified medicines.

Textbox 1. Steps for systematic case finding for falsified and substandard medicines.

- Step 1: Define indicators of public health importance (eg, burden of disease, vulnerability of affected population, narrow therapeutic index, sales volume of brand, or dosage form).
- Step 2: Define 1 or more objective indicators for each of the risk factors for substandard medicines and for falsified medicines, specifying the level at which it operates. Identify potential collinearity, and eliminate duplication.
- Step 3: Create risk scores for each numeric indicator (eg, none, minimal, some, or high), and calculate indicators and scores for each product (examples in Supplementary Tables A and B in [42]).
- Step 4: For each product, add risk scores to create a total index of risk. Select products to be sampled, prioritizing those with the highest risk index.
- Step 5: For sampled products, weigh by risks related to geography and supply chain, and draw up a sample frame.
- Step 6: Sample selected products (from specified locations, if indicated in Step 5).
- Step 7: Test sampled products. For potentially falsified products, screen visually and using rapid or low-cost devices such as a hand-held spectrometer or field-based thin layer chromatography. For potentially substandard products, perform quantitative assay and dissolution tests.

Step 1 is carried out in consultation with health authorities, while steps 2 and 3 take into account available data sources and the opinion of experts from the many sectors involved in the production, procurement, sale, and use of medicines. In [Table 4](#), we provided a single example of a possible indicator for each area of risk for substandard production or degradation. These suggestions derive from ongoing exploratory work in Indonesia, a large middle-income country with substantial domestic medicine production and a single-payer health insurance system. Exact specifications of the indicators, as well as decisions about potential weighting, may differ by country and will be determined, in large part, by the data available. A more comprehensive list of alternative indicators, together with suggested data sources, is provided in the supplementary tables in [42].

Because the medicine market is extremely heterogeneous, several indicators use relative measures, such as ratios compared with the median. These indicators must then be turned into scores that can be added together to create a total index of risk as described in [Textbox 1](#). Supplementary Table A in [42] suggests methods for turning indicators into risk scores.

To provide a single example for the first indicator in [Table 4](#), examine the ratio of price to weighted market median price for the same product. If the ratio is above 1, the product is priced above the market median and, thus, not deemed irrationally cheap or at risk of cost cutting. Deciles of risk are calculated only for those products with a price-to-median-price ratio of less than one. Products closest to the median (deciles 7-10) may also be considered at no risk. Those in deciles 5 and 6 may score

at 1 risk point (at minimal risk for cost cutting), and those in the second to fourth deciles score 2 points (at some risk). Brands (or nonbranded products from a specific market authorization holder) that fall into the first decile—the products selling at the deepest discount—are awarded 3 risk points (at high risk for cost cutting). Narrower gradations would allow for greater specificity; expert committees may decide what is most appropriate in the local context.

A similar process can be undertaken for products at risk for falsification, but the indicators will be different. [Table 5](#) provides examples for each of the major risk-factor groupings. Again, a more comprehensive list of alternative indicators, together with suggested data sources, is provided in the supplementary tables in [42].

The success of the case-finding approach will depend, to a significant extent, on the willingness of data custodians to share these data with those conducting case finding. The sensitivity and specificity of case finding will additionally depend on the ways in which indicators are combined. While [Textbox 1](#) describes a simple index, weighting is possible. If weighting is used, it is likely that brand-specific indicators, which have greater specificity, will carry a greater weight than market-wide indicators relating to molecules. However, we propose working with regulators to use retrospective data to find the model that best predicts poor-quality products in specific markets. Regulators with higher capacity for analysis may wish to develop more complex algorithms, including “big data” approaches, that combine price and volume data in ways that more closely pinpoint risk in specific markets.

Table 4. Indicative components to flag potentially substandard medicines.

Indicates	Indicator	Level at which indicator applies	Rationale
Profit pressure: cost cutting	Ratio of price to weighted market median for same product (same molecule and dosage form)	Brand and dosage form	Although premium brands are usually available, products produced by a large number of companies will tend toward the lowest cost of quality-assured production plus a fair profit [43]. If a particular product sells significantly below the market median, it may signal insufficient investment in quality assurance or other cost-cutting measures.
Technical limitation: production errors	Number of years continuously producing this molecule	Manufacturer (per molecule)	As companies and their staff gain experience and streamline their standard operating procedures in the production of a new medicine, the risk of production errors falls. Mistakes in production are more common among newly registered manufacturers.
Limited oversight	Time since most recent GMP ^a inspection of any facility	Manufacturer (by production site)	Medicine regulators aim to inspect production facilities on a regular basis; some additionally include risk-based inspection. In practice, frequency of inspection depends on regulatory capacity, and intervals may vary. The risk of detectable deviations from GMP grows with time since last inspection.
Production history	Number of regulatory warnings or sanctions given to manufacturer over reference period	Manufacturer (all products)	Investment in quality assurance is embedded in corporate culture. Manufacturers who repeatedly receive warnings for GMP violations may systematically underinvest in quality assurance, meaning all their products are at higher risk.
Reputational risk	Number of years of MA ^b holder in market	MA holder (all products)	Most companies are incentivized to invest in QA ^c , in part, because they wish to maintain their reputation as a provider of quality goods. New companies may be established opportunistically, especially, in rapidly growing markets. With less investment in building a reputation than older firms, new companies may have less to lose if found to be marketing substandard products.
Intrinsic risk: degradation	Stability of molecule	Molecule (all products)	Some molecules are less stable than others and more sensitive to variations in humidity, temperature, light, or other factors. Less stable molecules are more likely to degrade, becoming substandard before consumption.
Ecological risk: degradation	Classification of district accessibility	All products (or less stable products), by point of sale	Long supply chains and poor infrastructure pose challenges for maintaining temperature and humidity and may also reduce frequency of distribution. These factors increase the risk of degradation, especially, for less stable products.

^aGMP: good manufacturing practice.^bMA: market authorization.^cQA: quality assurance.

Table 5. Indicative components to flag potentially falsified medicines.

Indicates	Indicator	Level at which indicator applies	Rationale
Market opportunity: limited affordability	Product is on patent but not listed in current national formulary	Brand	On-patent products usually have premium prices. When they are not listed in current national formularies, they are usually not covered in the national insurance scheme, indicating limited affordability for patients. Patients or health care providers may seek these products at cut prices outside of the regulated supply chain.
Market opportunity: desirability	Molecule is used recreationally or off-label	Molecule	Some narcotics and psychotropic medicines are used recreationally or otherwise abused, including use for purposes for which they are not licensed. Additionally, access to some medicines is tightly restricted for political reasons, such as their potential use as abortifacients. Since the sale of these products is regulated, users without prescriptions commonly seek them outside of the regulated supply chain or from vendors who do not observe due diligence.
Profitability	Ratio of (price × retail channel sales volume) to market median, for the same dosage form	Brand	Falsifiers want to sell products for which there is a lucrative market, for which a large number of patients are prepared to pay a high price. For any given medicine for which there is a choice of brands, those brands with a combination of a relatively high retail price and a relatively large sales volume will be attractive targets.
Low risk of detection	Number of listings for product on 2 largest internet marketplaces	Brand	General internet marketplaces provide an unregulated but commonly used space for trading medicines without official licenses. The vast number of online transactions creates difficulty for regulatory monitoring, and anonymity limits the possibility of repercussions. More listings of products on the largest general online marketplaces also indicate high demand.

Proposed Method for Sentinel Surveillance: Tracking Trends in Risk Groups

Sentinel surveillance is a form of postmarket surveillance that is less data-intensive than case finding and has a different purpose. Systematic testing of comparable samples over time allows health authorities to: establish the likely prevalence of substandard medicines and, separately, of falsified medicines; inform estimates of the health and economic impacts of these medicines; make a case for additional investment in quality assurance in production or procurement, if necessary, including more investment in regulatory enforcement; plan and implement policies and programs to reduce prevalence of poor-quality medicines; and track progress over time toward achieving that goal.

The principal challenge in developing a robust, risk-based sentinel surveillance system for substandard and falsified medicines is in identifying “risk groups” of medicines, within which most poor-quality medicines cluster. These risk groups

are the functional equivalent of the risk behaviors that circumscribe sentinel populations for another invisible threat to health, HIV infection.

Similar to HIV sentinel surveillance, the specific sentinel groups may vary from country to country, depending on market dynamics and the risks and opportunities they create. The critical point is that groups are defined based on the feasibility of drawing samples and in ways that are replicable over time.

Drawing from existing risk-based approaches and the additional market risk factors identified in Table 2 and Table 3, and, again, with reference to ongoing exploratory work in the Indonesian market, we propose sentinel groups for both substandard medicines (Table 6) and falsified medicines (Table 7). We underline that these groupings are not intended to encompass all at-risk products, nor do we suggest that all of the products in these groupings are at risk. Rather, these factors act as proxies that may yield a higher concentration of at-risk products, compared with a random sample.

Table 6. Suggested sentinel groups for substandard medicines.

Sentinel group	Definition	Signals potential problem of	Rationale
Irrationally low-priced essential medicines	In public systems, medicine price <75% of international reference price; in retail pharmacies, cheapest available version of target medicine	Cost cutting	Irrationally low prices are a strong predictor for cost cutting. Selection of samples from the public system can be brand specific, and price data are available in advance, so a clear price threshold can be set. For private provision and retail sampling, brand-specific sampling is not feasible, so the cheapest medicine available should be sampled. Molecules should be selected for local public health importance.
Contract-manufactured medicines	Products randomly selected from those that are manufactured by a company other than market authorization holder	Reduced oversight	Contract manufacturing is a frequent means of lowering production costs (by outsourcing to companies that can achieve economies of scale). The market authorization holder does not always have clear oversight of quality assurance practices at contract manufacturers.
Poor regulatory history	Medicines randomly selected from those made by companies with history of regulatory violations and involuntary recalls within a time reference period	Inadequate quality assurance	Though regulators work with past violators to improve manufacturing and distribution practice, corporate culture and incentives appear to act as enablers of cost cutting and other practices that increase the risk of substandard production or degradation.
Technically vulnerable	Medicines randomly selected from those that are technically challenging to manufacture or distribute, including those with unstable molecules, limited therapeutic index, or sterile forms	Higher potential for production errors	These product-specific characteristics require particular investment in quality assurance. It would be possible to restrict this sentinel group to newer, less experienced manufacturers or market authorization holders. For unstable molecules, samples may be drawn from outlets in geographically remote areas.

Table 7. Suggested sentinel groups for falsified medicines.

Sentinel group	Definition	Signals potential problem of	Rationale
High irrational demand	Medicines that are used for recreation or other off-label purposes, for which alternatives are restricted or expensive; randomly sampled from retail outlets	Market opportunity	Demand planning is based on authorized uses only. Off-label use creates shortages, which provide market opportunities for falsifiers. Sample frame may be weighted toward independent pharmacies and medicine shops.
Life-saving but unaffordable	Medicines that are known or reputed to be life-saving, that retail at >10% of average per capita household spending, but that are not covered by national insurers; brand-specific sample from retail outlets	Market opportunity	Patients with life-threatening conditions are highly motivated to acquire these medicines. High profit margins incentivize their sale, which may diminish due diligence even in the regulated supply chain. This sample may include products not locally authorized; these should also be screened for falsification.
Sold on unregulated internet platforms	Random sample of “prescription-only” medicines sold through unlicensed internet sellers (sentinel group may be combined with signals of profit potential though purposive sampling of brands retailing at >200% of market median for the dosage form)	Evasion of regulation	Falsifiers favor internet sales because the potential for detection and successful prosecution is low. While the sample may be weighted toward medicines with high irrational demand, it should include other medicines of public health importance, such as antibiotics.

Most sentinel groups should be sampled in the public, private, and (if applicable) nonprofit sectors; the obvious exceptions are samples specific to unregulated channels, which should not

exist in the public sector. A sentinel surveillance system may be established following a process similar to that described in [Textbox 2](#).

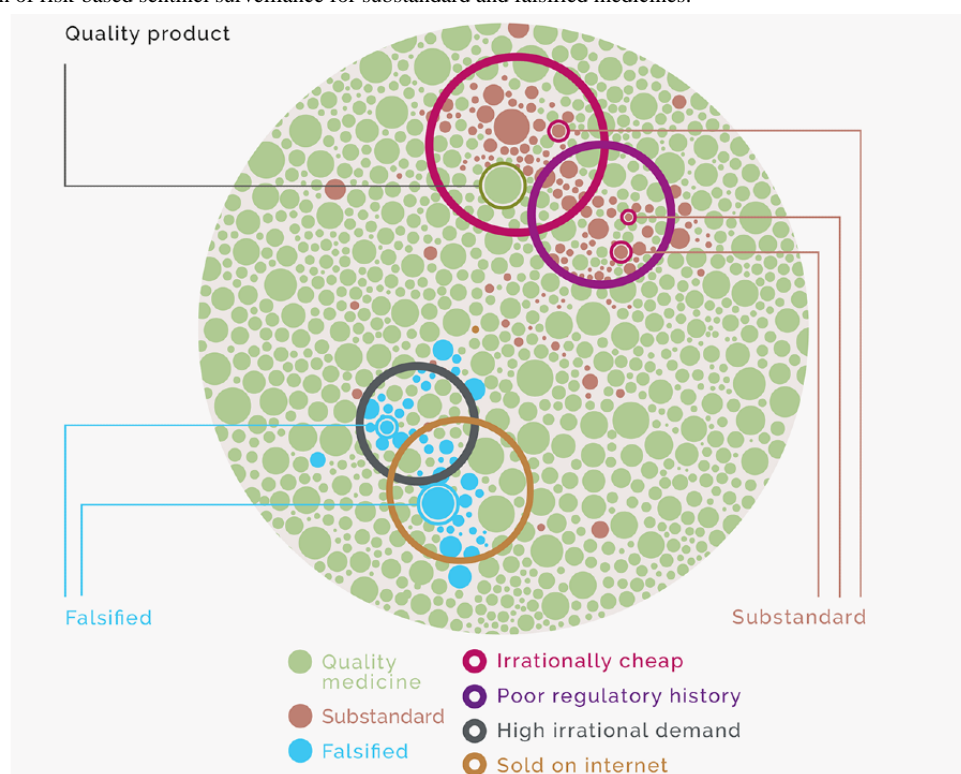
Textbox 2. Steps to establish sentinel surveillance for falsified and substandard medicines.

- Step 1: Define indicators of public health importance (eg, burden of disease, vulnerability of population, sales volume of brand or dosage form, or publicly procured); list medicines by public health importance.
- Step 2: Review local data sources and market conditions to define proxy “sentinel groups,” in which the highest concentrations of (A) falsified and (B) substandard medicines of public health importance are likely to be found.
- Step 3: Define indicators (or combinations of indicators) that best circumscribe those sentinel groups.
- Step 4: Draw up a sample frame for each sentinel group, including in it the public and private sectors, as appropriate, and sample a predetermined number of products for testing.
- Step 5: Test samples. For sentinel groups containing potentially falsified products, screen visually and using rapid or low-cost devices such as hand-held spectrometers or field-based thin layer chromatography. For sentinel groups containing potentially substandard products, perform quantitative assays and dissolution tests.

Pharmacopeial testing is expensive. The risk-based sentinel approach aims to reduce costs of routine surveillance in 2 ways. First, it increases the “yield” of testing by focusing it on the clusters of medicines most likely to be at risk (Figure 2). Second, it provides an initial triage for testing technologies. Products

selected in sentinel groupings for falsification risk can be screened visually and using lower-cost field-based devices [44]; only those at high risk of substandard production need to undergo assay and dissolution testing.

Figure 2. Illustration of risk-based sentinel surveillance for substandard and falsified medicines.



Developing National Estimates for the Prevalence of Substandard and Falsified Medicines

Steps to Estimate the Prevalence of Poor-Quality Medicines

In the same way that HIV prevalence among sex workers or drug injectors does not represent the prevalence of the virus in a whole population, the prevalence of poor-quality medicines

in risk-based sentinel groups does not represent medicine quality across a whole national market. However, if the size of each of those sentinel groups can be calculated and if assumptions can be made about their relationships with the wider population, then sentinel surveillance can provide a starting point for making robust estimates of national prevalence [45]. The same is true for medicine surveillance groups and their relationships with the wider market. Textbox 3 suggests a process for developing such estimates.

Textbox 3. Steps to estimate the national prevalence of falsified and substandard medicines.

- Step 1: Conduct risk-based sentinel surveillance for falsified and substandard medicines to obtain prevalence estimates for each sentinel group.
- Step 2: Calculate the interaction between sentinel groups.
- Step 3: Use market, procurement, and regulatory data to estimate the total volumes of medicines in each sentinel group.
- Step 4: Apply the prevalence estimates (step 1) to the volume data (step 3) to estimate the total number of substandard medicines across all sentinel groups for substandard medicines, correcting for interactions where necessary. Repeat for falsified medicines.
- Step 5: Use market, procurement, and regulatory data to estimate the total volumes of medicines in (lower-risk) market subsectors that do not fall in the sentinel groups.
- Step 6: Use all available data sources (eg, regulatory data, academic studies, case-reporting data, and data from other countries) to make assumptions about the residual prevalence in these nonsentinel sectors, comparing with prevalence in sentinel groups. Document each assumption, and, then, estimate the prevalence of substandard and falsified medicines in each sector.
- Step 7: Apply the prevalence estimates (step 6) to volumes (step 5) to estimate the total number of poor-quality medicines outside of sentinel groups.
- Step 8: Calculate the national prevalence estimate, where national prevalence estimate = (step 4 + step 6)/(step 3 + step 5).

The approach in [Textbox 3](#) has the great advantage of transparency. In addition, the process will highlight important data gaps that can be rectified over time. Assumptions can easily be corrected as more complete data are collected or shared. For example, step 6 might initially include the assumption that imported, on-patent oncology drugs that do not fall into any of the groups in [Table 6](#) are substandard 0% of the time. The assumption may change if regulators in producing countries issue product recalls for products that are also exported.

Estimation Process

Health-related estimates tend to improve in accuracy and local relevance (and, thus, potential utility) if potential end users are involved in their production [46]. This is, in part, because these actors can identify and contribute data to the process, and their expertise provides critical insights for informing necessary assumptions [47]. Estimation of the prevalence of substandard and falsified medicine should be led by medicine regulators and ministries of health. We would strongly encourage active consultation with pharmaceutical manufacturers and distributors, insurers, procurement authorities, consumer and patient advocates, and professional associations (eg, doctors and pharmacists) for deciding on methods and assumptions. The policies, decisions, and behaviors of these actors shape markets and influence the quality of medicines that patients consume [48]. Besides enriching the process, their participation increases the likelihood that estimates will be broadly accepted and acted upon. However, the motivations and interests of these groups are rarely aligned; careful annotation and transparent publication of all assumptions and data sources used in the estimation process can guard against capture by any interest group, protecting the integrity of the estimates [47].

Next Steps

The national estimates that result from this process will not capture the full complexity of medicine markets, especially, in initial rounds of estimation. However, we think it is important to begin to work toward that goal with tools that are most likely to be adopted by regulators in resource-constrained settings. We believe that a feasible and important first step in better quantifying the threat posed by substandard and falsified medicines includes simple, spreadsheet-based national models

that clearly document all data sources and assumptions and that are based on clearly defined and repeatable sentinel surveillance. Later, more sophisticated models may embrace more complexity and be expanded to estimate the extent to which these products undermine health and well-being and the damage they do to family and national budgets.

Our suggested methods may seem complex, and the processes may seem institutionally daunting, but, again, we draw inspiration from the experience of HIV surveillance. The current state of information systems for medicine quality closely resembles HIV surveillance systems circa 1995. Many low- and middle-income countries had no system at all. In those that did, incomplete case reporting was the norm; sentinel surveillance focused mainly on pregnant women even in countries where most infections were in men; behavioral risk surveillance was in its infancy; and estimation of population size was unheard of. Meanwhile, most estimates of national prevalence were made by a handful of people working for international organizations, using assumptions that did not reflect the diversity of national epidemics [49].

Now, the picture is very different. Most countries have developed surveillance systems based around the specific risks that drive their national epidemic and gather data related to risk behaviors and treatment outcomes in ways that are comparable across time. Population size estimation allows for the development of prevalence estimates that are useful in informing programming and measuring risks. Many different sectors cooperate in the implementation of these systems, which are largely country-led [50].

This transition in HIV surveillance was made possible by the vast, disease-specific investments in HIV seen from the early 2000s, investments that were themselves triggered, in part, by findings in countries such as Thailand, an “early-adopter” of risk-based surveillance for HIV. While we do not imagine that similar investments will be forthcoming in the field of medicine quality, we believe that increased national investments in medicine procurement through expanded efforts to achieve universal health coverage will increase the urgency of ensuring that public money is invested in medicines that actually cure patients or prevent disease, rather than in their substandard or

falsified doppelgangers. It is, thus, a good time to start building the capacities and systems that will improve the ways in which health systems measure; understand; and, ultimately, curtail the extent and distribution of substandard and falsified medicines.

Prioritizing Surveillance Approaches for Medicine Quality

As with infectious disease surveillance, the 2 approaches we have suggested for surveillance of medicine quality—case finding and sentinel surveillance—have different purposes. Case finding aims to pinpoint problems for an immediate response, while sentinel surveillance allows for more reliable quantification of the problem and for the monitoring of the effectiveness of interventions. It is unclear in situations where resources are constrained, which one a regulator should prioritize.

In the short term, especially, from the point of view of the medicine regulator who will be held responsible if substandard or falsified medicines are shown to harm patients, case finding will probably be the more attractive option. This is true despite the fact that case finding is more data intensive and, likely, more technically challenging to implement, because it requires greater specificity to succeed than sentinel surveillance. Case finding is likely to be the more valuable approach in settings where the regulator is well resourced and where substandard and falsified medicines are comparatively rare.

Where it is suspected that falsified and, especially, substandard medicines may be rather more prevalent, however, the calculus changes, at least, from a broader public health point of view. Here, a more robust understanding of the extent of the problem, provided by estimates based on sentinel surveillance, may be more valuable. In such settings, which may include many low- and middle-income countries, the prevalence of substandard and, to a lesser extent, falsified medicines will likely have a system-wide effect on health outcomes, as well as on public and private finances. Only after estimating prevalence can one reliably estimate impact. Robust estimates of impact are politically persuasive and may encourage policy makers to change health financing, procurement, and industrial policies in ways intended to erode the factors that incentivize the production and sale of substandard medicines and to shrink the market for falsified products. In addition, a clear understanding of the magnitude of the problem may prove a powerful argument for adequate resourcing for medicine regulators, especially, in lower-income settings.

We are currently consulting closely with the Indonesian national medicine regulator to trial risk-based case finding, as well as to plan and implement sentinel surveillance and develop national estimates using the methods suggested in this paper. We welcome challenges to our thinking and suggestions to improve the proposed methods, and we look forward to continued debate on the subject.

Acknowledgments

AH's work was supported through a fellowship from the United States Pharmacopeia Quality Institute, which also provided support for supervision to EP and MOK. MOK's work was also supported by the Netherlands Research Excellence Initiative (REI). These funders have had no part in preparing this manuscript.

We thank Annie Wang for contributing to initial definitions for case finding. We thank Michael Deats, Sara de Valente, Katharina Hauk, Paul Newton, and Koray Parmaksiz for helpful comments on an earlier draft.

Authors' Contributions

The paper was conceived and drafted by EP. EP, SH, AH, YA, and MR, together, developed the case finding indicators. MOK provided supervision for AH. All authors commented on the manuscript and approved the final version.

Conflicts of Interest

AH's work on risk-based surveillance was supported by a research fellowship from the United States of Pharmacopeial Convention (USP) Quality Institute; USP Quality Institute also provided supervision support for EP and MK. In addition, EP, YA, AH, and MR currently receive research funding from UK National Institute for Health Research (GHPSR Project: NIHR131145) to develop and trial sentinel surveillance and estimation methods related to medicine quality. MOK receives research support from the Netherlands Research Excellence Initiative. SH is an employee of the public health division of health care data science firm IQVIA. IQVIA provides data and insights on pharmaceutical market dynamics in over 100 countries. The company contributed research and staff time for exploration of the ideas proposed in this paper at no cost.

References

1. 1 in 10 medical products in developing countries is substandard or falsified. World Health Organization. 2017 Nov 28. URL: <https://www.who.int/news/item/28-11-2017-1-in-10-medical-products-in-developing-countries-is-substandard-or-falsified> [accessed 2020-12-21]
2. Newton PN, Bond KC, Adeyeye M, Antignac M, Ashenef A, Awab GR, et al. COVID-19 and risks to the supply and quality of tests, drugs, and vaccines. *The Lancet Global Health* 2020 Jun;8(6):e754-e755. [doi: [10.1016/s2214-109x\(20\)30136-4](https://doi.org/10.1016/s2214-109x(20)30136-4)]
3. Willmer, G. New alliance seeks to fight plague of fake medicines. *SciDevNet*. 2021 Dec 14. URL: <https://www.scidev.net/global/news/new-alliance-seeks-to-fight-plague-of-fake-medicines/> [accessed 2020-12-21]

4. World Health Organization. Medical Product Alert N°3/2020. WHO Rapid Alerts. 2020. URL: <https://www.who.int/news-room/detail/31-03-2020-medical-product-alert-n-3-2020> [accessed 2020-04-09]
5. World Health Organization. Medical Product Alert N°4/2020. WHO Rapid Alerts. 2020. URL: <https://www.who.int/news-room/detail/09-04-2020-medical-product-alert-n4-2020> [accessed 2020-04-09]
6. Report of the fifth meeting of the Member State mechanism on substandard/spurious/falsely-labelled/falsified/counterfeit medical products. World Health Organization. 2017 Jan 17. URL: https://apps.who.int/gb/sf/pdf_files/MSM5/A_MS5_8-en.pdf [accessed 2017-02-14]
7. WHO Global Surveillance Monitoring System for Substandard and Falsified Medical Products. World Health Organization. 2017. URL: http://www.who.int/medicines/regulation/ssffc/publications/GSMS_Report.pdf [accessed 2018-02-09]
8. Buckley G, Gostin L. Countering the Problem of Falsified and Substandard Drugs: Committee on Understanding the Global Public Health Implications of Substandard, Falsified, and Counterfeit Medical Products. In: B, Lawrence O. Gostin. Washington DC: Institute of Medicine of the National Academies; 2013.
9. Jackson KD, Higgins CR, Laing SK, Mwila C, Kobayashi T, Ippolito MM, et al. Impact of substandard and falsified antimalarials in Zambia: application of the SAFARI model. BMC Public Health 2020 Jul 09;20(1):1083 [FREE Full text] [doi: [10.1186/s12889-020-08852-w](https://doi.org/10.1186/s12889-020-08852-w)] [Medline: [32646393](https://pubmed.ncbi.nlm.nih.gov/32646393/)]
10. Beagrie SM, Higgins CR, Evans DR, Laing SK, Erim D, Ozawa S. The economic impact of substandard and falsified antimalarial medications in Nigeria. PLoS One 2019 Aug 15;14(8):e0217910 [FREE Full text] [doi: [10.1371/journal.pone.0217910](https://doi.org/10.1371/journal.pone.0217910)] [Medline: [31415560](https://pubmed.ncbi.nlm.nih.gov/31415560/)]
11. Ozawa S, Haynie D, Bessias S. Modeling the Economic Impact of Substandard and Falsified Antimalarials in the Democratic Republic of the Congo. Am J Trop Med Hyg ? 2019;100:57. [doi: [10.4269/ajtmh.18-0334](https://doi.org/10.4269/ajtmh.18-0334)]
12. Renschler R, Walters K, Newton P, Laxminarayan R. Estimated under-five deaths associated with poor-quality antimalarials in sub-Saharan Africa. Am J Trop Med Hyg 2015 Jun;92(6 Suppl):119-126 [FREE Full text] [doi: [10.4269/ajtmh.14-0725](https://doi.org/10.4269/ajtmh.14-0725)] [Medline: [25897068](https://pubmed.ncbi.nlm.nih.gov/25897068/)]
13. World Health Organization. A study on the public health and socioeconomic impact of substandard and falsified medical products. World Health Organization. 2017. URL: <http://who.int/medicines/regulation/ssffc/publications/Layout-SEstudy-WEB.pdf> [accessed 2017-11-28]
14. Johnston A, Holt D. Substandard drugs: a potential crisis for public health. Br J Clin Pharmacol 2014 Aug;78(2):218-243 [FREE Full text] [doi: [10.1111/bcp.12298](https://doi.org/10.1111/bcp.12298)] [Medline: [24286459](https://pubmed.ncbi.nlm.nih.gov/24286459/)]
15. McManus D, Naughton B. A systematic review of substandard, falsified, unlicensed and unregistered medicine sampling studies: a focus on context, prevalence, and quality. BMJ Glob Health 2020 Aug;5(8):a [FREE Full text] [doi: [10.1136/bmjgh-2020-002393](https://doi.org/10.1136/bmjgh-2020-002393)] [Medline: [32859648](https://pubmed.ncbi.nlm.nih.gov/32859648/)]
16. Rojas-Cortés R. Substandard, falsified and unregistered medicines in Latin America, 2017-2018. Rev Panam Salud Publica 2020;10:44. [doi: [10.26633/rpsp.2020.125](https://doi.org/10.26633/rpsp.2020.125)]
17. Nayyar GM, Breman J, Mackey T. Falsified and substandard drugs: stopping the pandemic. Am J Trop Med Hyg 2019;8. [doi: [10.4269/ajtmh.18-098](https://doi.org/10.4269/ajtmh.18-098)]
18. Ozawa S, Evans DR, Bessias S, Haynie DG, Yemeke TT, Laing SK, et al. Prevalence and Estimated Economic Burden of Substandard and Falsified Medicines in Low- and Middle-Income Countries: A Systematic Review and Meta-analysis. JAMA Netw Open 2018 Aug 03;1(4):e181662 [FREE Full text] [doi: [10.1001/jamanetworkopen.2018.1662](https://doi.org/10.1001/jamanetworkopen.2018.1662)] [Medline: [30646106](https://pubmed.ncbi.nlm.nih.gov/30646106/)]
19. Hodges S, Garnett E. The ghost in the data: Evidence gaps and the problem of fake drugs in global health research. Global Public Health 2020 Mar 31;15(8):1103-1118. [doi: [10.1080/17441692.2020.1744678](https://doi.org/10.1080/17441692.2020.1744678)]
20. Garcia-Abreu A, Halperin S, Danel I. Public health surveillance toolkit: a guide for busy task managers. World Bank. 2002. URL: <https://openknowledge.worldbank.org/bitstream/handle/10986/20817/563720WP0REPLA00Box349502B00PUBLIC0.pdf> [accessed 2020-08-20]
21. Choi BCK. The Past, Present, and Future of Public Health Surveillance. Scientifica 2012;2012:1-26. [doi: [10.6064/2012/875253](https://doi.org/10.6064/2012/875253)]
22. Choi BC, Pak A. Lessons for surveillance in the 21st century: A historical perspective from the past five millennia. Soz Präventivmed 2001 Nov;46(6):361-368. [doi: [10.1007/bf01321662](https://doi.org/10.1007/bf01321662)]
23. Simonsen L, Gog JR, Olson D, Viboud C. Infectious Disease Surveillance in the Big Data Era: Towards Faster and Locally Relevant Systems. J Infect Dis 2016 Dec 01;214(suppl_4):S380-S385 [FREE Full text] [doi: [10.1093/infdis/jiw376](https://doi.org/10.1093/infdis/jiw376)] [Medline: [28830112](https://pubmed.ncbi.nlm.nih.gov/28830112/)]
24. Kirian ML, Weintraub JM. Prediction of gastrointestinal disease with over-the-counter diarrheal remedy sales records in the San Francisco Bay Area. BMC Med Inform Decis Mak 2010 Jul 20;10(1):39 [FREE Full text] [doi: [10.1186/1472-6947-10-39](https://doi.org/10.1186/1472-6947-10-39)] [Medline: [20646311](https://pubmed.ncbi.nlm.nih.gov/20646311/)]
25. Nsubuga P, White M, Thacker S. Public Health Surveillance: A Tool for Targeting Monitoring Interventions. In: Jamison DT, Breman JG, Measham AR, editors. Disease Control Priorities in Developing Countries. Washington DC: World Bank; 2006.
26. Schwartländer B, Ghys PD, Pisani E, Kiessling S, Lazzari S, Caraël M, et al. HIV surveillance in hard-to-reach populations. AIDS 2001;15:S1-S3. [doi: [10.1097/00002030-200104003-00001](https://doi.org/10.1097/00002030-200104003-00001)]

27. Pisani E, Lazzari S, Walker N, Schwartländer B. HIV surveillance: a global perspective. *J Acquir Immune Defic Syndr* 2003 Feb;32 Suppl 1:S3-11. [doi: [10.1097/00126334-200302011-00002](https://doi.org/10.1097/00126334-200302011-00002)] [Medline: [12571509](#)]
28. Newton PN, Lee S, Goodman C, Fernández FM, Yeung S, Phanouvong S, et al. Guidelines for field surveys of the quality of medicines: a proposal. *PLoS Med* 2009 Mar 24;6(3):e52 [FREE Full text] [doi: [10.1371/journal.pmed.1000052](https://doi.org/10.1371/journal.pmed.1000052)] [Medline: [19320538](#)]
29. WHO Expert Committee on Specifications for Pharmaceutical Preparations. Guidelines on the Conduct of Surveys of the Quality of Medicines. WHO Technical Report Series, No. 2016. URL: <http://apps.who.int/medicinedocs/documents/s22404en/s22404en.pdf> [accessed 2017-12-07]
30. Vida RG, Merczel S, Jahn E, Fittler A. Developing a framework regarding a complex risk based methodology in the evaluation of hazards associated with medicinal products sourced via the internet. *Saudi Pharm J* 2020 Dec;28(12):1733-1742 [FREE Full text] [doi: [10.1016/j.jsps.2020.10.018](https://doi.org/10.1016/j.jsps.2020.10.018)] [Medline: [33424264](#)]
31. General European OMCL Network (GEON). Incorporation of a risk based approach in market surveillance testing at OMCLs. General European OMCL Network. 2007. URL: https://www.edqm.eu/sites/default/files/omcl_incorporation_of_a_rb_approach_in_ms_testing_at_omcls.pdf [accessed 2020-08-24]
32. Aroca Á, Guzmán J. Model for a risk-focused approach to health inspection, surveillance, and control in Colombia. *Rev Panam Salud Publica Pan Am J Public Health* (41) 2017:e105. [doi: [10.26633/RPSP.2017.105](https://doi.org/10.26633/RPSP.2017.105)]
33. United States Pharmacopeial Convention, Babigumira J, Stegarchis A. A risk-based resource allocation framework for pharmaceutical quality assurance for medicines regulatory authorities in low- and middle-income countries. In: *USP Promoting Quality of Medicines*. Washington DC: United States Pharmacopeial Convention; Jun 2018.
34. United States Pharmacopeial Convention, Nkansah P, Smine K. Guidance for implementing risk-based post-marketing quality surveillance in low- and middle-income countries. In: *USP Promoting Quality of Medicines*. Washington DC: United States Pharmacopeial Convention; 2018:2018.
35. United States Food and Drug Administration. Predictive Risk-based Evaluation for Dynamic Import Compliance Targeting (PREDICT). FDA. 2014. URL: <https://www.fda.gov/media/83668/download> [accessed 2021-04-04]
36. FDA Center for Drug Evaluation Research. Drug Quality Sampling Testing Programs. FDA. 2021 Mar 02. URL: <https://www.fda.gov/drugs/science-and-research-drugs/drug-quality-sampling-and-testing-programs> [accessed 2021-06-09]
37. Taberner P, Mayxay M, Culzoni M, Dwivedi P, Swamidoss I, Allan EL, et al. A Repeat Random Survey of the Prevalence of Falsified and Substandard Antimalarials in the Lao PDR: A Change for the Better. *Am J Trop Med Hyg* 2015 Jun;92(6 Suppl):95-104 [FREE Full text] [doi: [10.4269/ajtmh.15-0057](https://doi.org/10.4269/ajtmh.15-0057)] [Medline: [25897062](#)]
38. World Health Organization. WHO List of Prequalified Quality Control Laboratories, 51st Edition. World Health Organization. 2021. URL: <https://extranet.who.int/pqweb/medicines/medicines-quality-control-laboratories-list> [accessed 2021-02-03]
39. US Pharmacopeia. USP Reference Standards Catalog. USP. 2020. URL: <https://static.usp.org/doc/referenceStandards/dailycatalog.pdf> [accessed 2020-12-21]
40. Pisani E, Nistor A, Hasnida A, Parmaksiz K, Xu J, Kok MO. Identifying market risk for substandard and falsified medicines: an analytic framework based on qualitative research in China, Indonesia, Turkey and Romania. *Wellcome Open Res* 2019;4:70 [FREE Full text] [doi: [10.12688/wellcomeopenres.15236.1](https://doi.org/10.12688/wellcomeopenres.15236.1)] [Medline: [31131333](#)]
41. Hamilton WL, Doyle C, Halliwell-Ewen M, Lambert G. Public health interventions to protect against falsified medicines: a systematic review of international, national and local policies. *Health Policy Plan* 2016 Dec;31(10):1448-1466. [doi: [10.1093/heapol/czw062](https://doi.org/10.1093/heapol/czw062)] [Medline: [27311827](#)]
42. Pisani E, Hasnida A, Rahmi M, Harsono S, Kok MO, Anggriani Y. Supplementary material for Pisani et al: Substandard and falsified medicines: proposed methods for case finding and sentinel surveillance. *Harvard Dataverse*, V2 2021 [FREE Full text] [doi: [10.7910/DVN/SELJ0Z](https://doi.org/10.7910/DVN/SELJ0Z)]
43. Hill AM, Barber MJ, Gotham D. Estimated costs of production and potential prices for the WHO Essential Medicines List. *BMJ Glob Health* 2018 Jan 29;3(1):e000571 [FREE Full text] [doi: [10.1136/bmjgh-2017-000571](https://doi.org/10.1136/bmjgh-2017-000571)] [Medline: [29564159](#)]
44. Vickers S, Bernier M, Zambrzycki S, Fernandez FM, Newton PN, Caillet C. Field detection devices for screening the quality of medicines: a systematic review. *BMJ Glob Health* 2018;3(4):e000725 [FREE Full text] [doi: [10.1136/bmjgh-2018-000725](https://doi.org/10.1136/bmjgh-2018-000725)] [Medline: [30233826](#)]
45. World Health Organization, UNAIDS. Guidelines on estimating the size of populations most at risk to HIV. In: *Guidelines on estimating the size of populations most at risk to HIV*. Geneva: World Health Organization; 2010.
46. Pisani E, Kok M. In the eye of the beholder: to make global health estimates useful, make them more socially robust. *Glob Health Action* 2017;10(sup1):1266180 [FREE Full text] [doi: [10.3402/gha.v9.32298](https://doi.org/10.3402/gha.v9.32298)] [Medline: [28532303](#)]
47. UNAIDS, World Health Organization. Case study on estimating HIV infection in a concentrated epidemic: lessons from Indonesia. In: *Case study on estimating HIV infection in a concentrated epidemic: lessons from Indonesia*. Geneva: UNAIDS; 2004.
48. Hasnida A, Kok MO, Pisani E. Challenges in maintaining medicine quality while aiming for universal health coverage: a qualitative analysis from Indonesia. *BMJ Glob Health* 2021 May 28;6(Suppl 3):e003663 [FREE Full text] [doi: [10.1136/bmjgh-2020-003663](https://doi.org/10.1136/bmjgh-2020-003663)] [Medline: [34049935](#)]
49. Pisani E. The wisdom of whores: Bureaucrats, brothels and the business of AIDS. In: *The wisdom of whores: Bureaucrats, brothels and the business of AIDS*. London: Granta; 2008.

50. Mahy M, Brown T, Stover J, Walker N, Stanecki K, Kirungi W, et al. Producing HIV estimates: from global advocacy to country planning and impact measurement. *Glob Health Action* 2017;10(sup1):1291169 [FREE Full text] [doi: [10.1080/16549716.2017.1291169](https://doi.org/10.1080/16549716.2017.1291169)] [Medline: [28532304](https://pubmed.ncbi.nlm.nih.gov/28532304/)]

Abbreviations

WHO: World Health Organization

Edited by Y Khader; submitted 01.04.21; peer-reviewed by H Akram, P Nguyen, R Bright; comments to author 13.05.21; revised version received 09.06.21; accepted 27.06.21; published 16.08.21.

Please cite as:

Pisani E, Hasnida A, Rahmi M, Kok MO, Harsono S, Anggriani Y

Substandard and Falsified Medicines: Proposed Methods for Case Finding and Sentinel Surveillance

JMIR Public Health Surveill 2021;7(8):e29309

URL: <https://publichealth.jmir.org/2021/8/e29309>

doi: [10.2196/29309](https://doi.org/10.2196/29309)

PMID: [34181563](https://pubmed.ncbi.nlm.nih.gov/34181563/)

©Elizabeth Pisani, Amalia Hasnida, Mawaddati Rahmi, Maarten Olivier Kok, Steven Harsono, Yusi Anggriani. Originally published in *JMIR Public Health and Surveillance* (<https://publichealth.jmir.org>), 16.08.2021. This is an open-access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in *JMIR Public Health and Surveillance*, is properly cited. The complete bibliographic information, a link to the original publication on <https://publichealth.jmir.org>, as well as this copyright and license information must be included.

Original Paper

Identifying Communities at Risk for COVID-19–Related Burden Across 500 US Cities and Within New York City: Unsupervised Learning of the Coprevalence of Health Indicators

Andrew Deonarine¹, MSc, MHSc, MD, PhD; Genevieve Lyons¹, MSc; Chirag Lakhani¹, PhD; Walter De Brouwer¹, PhD

XY.ai, Cambridge, MA, United States

Corresponding Author:

Andrew Deonarine, MSc, MHSc, MD, PhD

XY.ai

56 JFK Street

Cambridge, MA, 02138

United States

Phone: 1 8575000461

Email: andrew@xy.ai

Abstract

Background: Although it is well-known that older individuals with certain comorbidities are at the highest risk for complications related to COVID-19 including hospitalization and death, we lack tools to identify communities at the highest risk with fine-grained spatial resolution. Information collected at a county level obscures local risk and complex interactions between clinical comorbidities, the built environment, population factors, and other social determinants of health.

Objective: This study aims to develop a COVID-19 community risk score that summarizes complex disease prevalence together with age and sex, and compares the score to different social determinants of health indicators and built environment measures derived from satellite images using deep learning.

Methods: We developed a robust COVID-19 community risk score (COVID-19 risk score) that summarizes the complex disease co-occurrences (using data for 2019) for individual census tracts with unsupervised learning, selected on the basis of their association with risk for COVID-19 complications such as death. We mapped the COVID-19 risk score to corresponding zip codes in New York City and associated the score with COVID-19–related death. We further modeled the variance of the COVID-19 risk score using satellite imagery and social determinants of health.

Results: Using 2019 chronic disease data, the COVID-19 risk score described 85% of the variation in the co-occurrence of 15 diseases and health behaviors that are risk factors for COVID-19 complications among ~28,000 census tract neighborhoods (median population size of tracts 4091). The COVID-19 risk score was associated with a 40% greater risk for COVID-19–related death across New York City (April and September 2020) for a 1 SD change in the score (risk ratio for 1 SD change in COVID-19 risk score 1.4; $P<.001$) at the zip code level. Satellite imagery coupled with social determinants of health explain nearly 90% of the variance in the COVID-19 risk score in the United States in census tracts ($r^2=0.87$).

Conclusions: The COVID-19 risk score localizes risk at the census tract level and was able to predict COVID-19–related mortality in New York City. The built environment explained significant variations in the score, suggesting risk models could be enhanced with satellite imagery.

(*JMIR Public Health Surveill* 2021;7(8):e26604) doi:[10.2196/26604](https://doi.org/10.2196/26604)

KEYWORDS

COVID-19; satellite imagery; built environment; social determinants of health; machine learning; artificial intelligence; community; risk; United States; indicator; comorbidity; environment; population; determinant; mortality; prediction

Introduction

The COVID-19 pandemic has disrupted major world economies and overwhelmed hospital intensive care units worldwide [1]. In the United States alone, the virus has spread throughout urban and rural communities and killed over 300,000 Americans to date [2]. Case series and epidemiological surveillance data from the United States [3-6] and around the world [7-11] have implicated risk factors for COVID-19-related morbidity and mortality, including older age, male sex, impaired lung function, cardiometabolic-related diseases (eg, diabetes, heart disease, or stroke), and obesity. In the United States, comorbidities are known to cluster in geographies such as the southeast states and counties (eg, in chronic disease [12] and in COVID-19 [13-16]), and are partly mediated by built environment features, such as walkability [17]. Although race and ethnicity have been identified as risk factors, systemic racism and discrimination in the health care system play an important role in this relationship [18-20]. Additionally, racial and ethnic discrimination have influenced where individuals reside and has played a substantial role in the increased morbidity and mortality related to COVID-19 [21]. Other factors including the built environment and air pollution have been associated with COVID-19 infection and complications [22,23], but it has been unclear how to prioritize these associations to prevent complications. Both individual-level factors (eg, diabetes, smoking, and asthma [3,8,10,11]) and geographical-level social determinant factors (eg, census tract-level population density and increased household occupancy) are strong risk factors for COVID-19 infection and risk [24]. Social determinants of health are defined as “conditions in the environments where people are born, live, learn, work, play, worship, and age that affect a wide range of health, functioning, and quality-of-life outcomes and risks” [25]. Social determinants of health can be grouped into five domains, including economic stability, education access and quality, health care access and quality, neighborhood and built environment, and social and community context [25]. Recently, Maharana and Nsoesie [26] developed an approach to map the built environment to obesity prevalence using deep learning analysis of satellite imagery, highlighting a potentially novel method of using measurements of the built environment to quantify disease risk.

At the time of writing, New York emerged as a location with several COVID-19-related deaths spread across the 2141 census tracts in the city. Even within city hot spots like New York City, common chronic diseases and their risk factors for COVID-19 are geographically heterogeneous and vary per unit of geography, including within and across states, counties, and even cities. It is unclear how the heterogeneity of community-based risk or prevalence of diseases at a census tract level (median population sizes of ~3000-5000 individuals) is related to COVID-19 risk. Furthermore, analyses on coarser spatial resolutions will attenuate predictions and associations [27].

In this investigation, we sought to create a clinically focused risk score that could be used to predict COVID-19 cases and deaths within cities, identify hot spots at the subcounty (census tract) level, and identify potentially vulnerable communities,

and to determine how the social determinants of health and the built environment may explain the variance of this clinically focused risk score and whether the built environment explains statistically significant amounts of score variance even after accounting for the social determinants of health. To do this, we developed the COVID-19 community risk score (COVID-19 risk score) that summarizes the complex comorbidity and demographic patterns of small communities at the census tract level. Additionally, we examined how the social determinants of health (including the built environment, measured using satellite imagery methods [26]) explained score variance and validated the risk score by examining its relationship with zip code-level deaths during the late-May 2020 COVID-19 epidemic in New York City. Last, we deployed the COVID-19 risk score with an application programming interface and a browsable dashboard [28].

Methods

Study Data

We obtained geocoded disease prevalence data at the census tract level from the US Centers for Disease Control and Prevention (CDC) 500 Cities Project (the December 2019 release, which is based on data from 2016 to 2017 [29]; Figure 1A). The project 500 Cities contains disease and health indicator prevalence for 27,648 individual census tracts of the 500 largest cities in the United States, and these prevalences are estimated from the Behavioral Risk Factor Surveillance System [30].

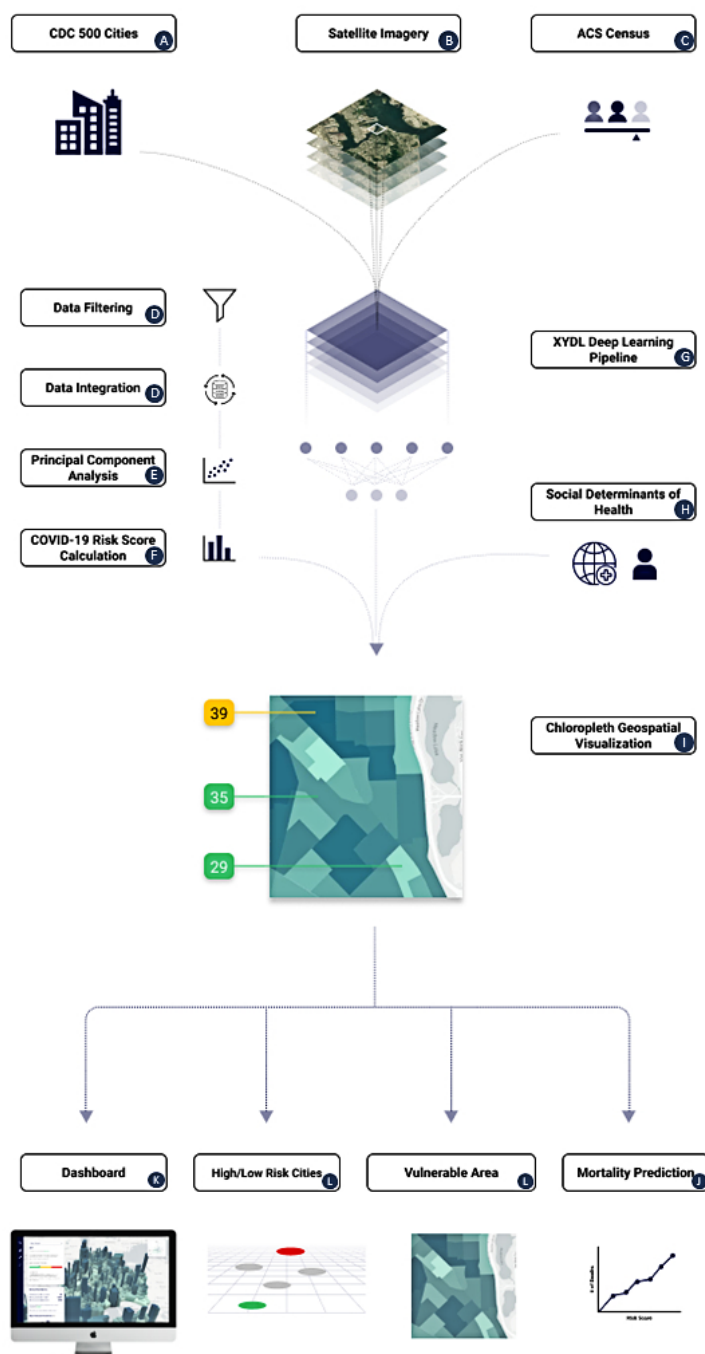
From the 500 Cities data, we chose 13 population-level health indicators that correspond to individual-level chronic disease risk factors associated with COVID-19-related hospitalization and death based on reports from China, Italy, and the United States (eg, [3,8,10,11]). Disease indicators include the prevalence among adults of diabetes, coronary heart disease, chronic kidney disease, asthma, arthritis, any cancer, or chronic obstructive pulmonary disorder. We also selected behavioral risk factors including smoking and obesity, and the prevalence of individuals on blood pressure medication. We chose these comorbidities and risk factors with guidance from the CDC because they were classed as among the strongest risk factors for COVID-19-related hospitalization, intensive care unit use, and death (eg, males and females older than 65 years, diabetes, heart disease, and stroke); were indicative of risk for cardiometabolic disease or impaired lung function, which are risk factors for COVID-19 (eg, smoking, obesity, high blood pressure, high cholesterol, kidney disease, asthma, or chronic pulmonary obstructive disorder); or involve pharmacological interventions that could result in an immunocompromised state (eg, certain antineoplastic, arthritis, and antihypertensive medications) [31].

We further obtained 5-year 2013-2017 American Community Survey (ACS) Census data [32], which contain sociodemographic prevalences and median values for census tracts (Figure 1C), and corresponded to the 2016-2017 CDC 500 Cities data. We also selected the total number of individuals in the tract, proportion of males and females older than 65 years, and proportion of individuals by race and ethnicity, which included African American, Mexican, Hispanic, Asian, and

White groups from the ACS. Race and ethnicity were examined to determine if there were different risks associated with these groups (where race is a socially constructed concept that can be used as a proxy for the complex interplay of institutional and individual-level racism and barriers to health care experienced by these different groups [33]). These data also included

information on socioeconomic indicators including median income, the proportion of individuals living in poverty, unemployment, cohabitation with more than one individual per room, and having no health insurance. These measures were previously identified as possible contributors to increased risk of infection or mortality associated with COVID-19 [20,34].

Figure 1. Overview of study. (A) CDC 500 Cities; (B) satellite imagery of 500 cities from OpenMapTiles; (C) ACS Census summary statistics for each census tract; (D) estimates of prevalence and coprevalence of disease and health indicators for risk of COVID-19 complications; (E) use of principal components analysis to reduce dimensionality of diseases and health indicators; (F) construction of COVID-19 score from principal components; (G) “XYDL” deep learning pipeline that inputs satellite imagery, social determinants of health indicators from ACS Census data to predict COVID-19 community risk score; (H) social determinants of health from ACS Census data; (I) visualization of the COVID-19 community risk score; (J) association of the COVID-19 risk score with mortality in New York City; (K) creation of a dashboard; (L) mapping highest and lowest risk cities and tracts as a function of the risk score. ACS: American Community Survey; CDC: Centers for Disease Control and Prevention.



Defining the COVID-19 Community Risk Score

Given the complex interplay between the social determinants of health, chronic disease, and the built environment, we sought to first examine how clinical comorbidities could be used to predict COVID-19 rates by developing a clinically focused risk score and then examine how these comorbidities relate to the built environment and social determinants of health. Understanding if the built environment and social determinants of health can explain the variance of a clinically focused risk score would show that more complex risk models could be built using this data in the future. To do this, we used the statistical programming language R (version 4.0.5; R Foundation for Statistical Computing) [35] to merge disease and behavior prevalence data from the CDC 500 Cities Project for each of the 27,648 census with ACS information and calculate their Pearson pairwise correlations (Figure 1D) to determine how the data were correlated with each other. We considered 15 variables in total, including 13 health indicators (eg, diseases and risk factors), and 2 demographic factors, the proportion of male and female individuals older than 65 years in the risk score. The disease prevalence included any form of cancer, arthritis, stroke, chronic asthma, chronic obstructive pulmonary disease (COPD), heart disease, diabetes, kidney disease, high blood pressure, and high cholesterol. Behavioral and lifestyle-related risk factors included smoking, obesity, and the rate of individuals on blood pressure medication. Finally, demographic factors included the prevalence of males and the prevalence of females older than 65 years.

Socioeconomic Correlates of the Community COVID-19 Risk Score

Next, we examined the relationship between the ACS-derived sociodemographic indicators with the COVID-19 risk score. This was done by calculating multivariate linear and random forest regressions to test the linear and nonlinear contribution of the sociodemographic indicators in the COVID-19 score (Figure 1H), and provide insight into the relationship of sociodemographic factors and the clinical indicators used in the COVID-19 score. This comparison to sociodemographic factors also serves as a form of validation, as the risk increases, one would expect certain sociodemographic indicators to also increase, such as poverty. Further details concerning the calculation of the linear and random forest regression can be found in Multimedia Appendix 1 [28,35-38].

Association of the COVID-19 Community Risk Score With Satellite Imagery

To correlate the COVID-19 risk score from satellite imagery (Figure 1B), millions of satellite images ($n=4,742,919$) were analyzed in an ensemble of an unsupervised deep learning algorithm and a supervised machine learning algorithm. The images are satellite raster tiles that were downloaded from the OpenMapTiles database. The images have a spatial resolution close to 20 meters per pixel, allowing a maximum zoom level of 13 [39]. Images were extracted in tiles from the OpenMapTiles database using the coordinate geometries of the census tracts. After extraction, images were digitally enlarged to achieve a zoom level of 18.

Many census tracts are large enough to contain multiple satellite images. The median number of images per tract was 94, and the number of images per census tract ranged from 1 image in the census tract to the largest geographical tract with 162,811 images (in Anchorage, Arkansas) with an IQR from 43 to 182 images. The geographical coverage of the images per census tract ranged from the smallest census tract covering 0.022 km^2 and the largest census tract covering 5679.52 km^2 , with an IQR from 0.93 km^2 to 3.89 km^2 and a median of 1.92 km^2 per census tract.

First, using the Python 3.7.7 programming language [40], we passed images through AlexNet [41], a pretrained convolutional neural network, in an unsupervised deep learning approach called feature extraction [42] (Figure 1G). The resulting vector from this process is a *latent space feature* representation of the image comprising 4096 features. This latent space representation is essentially an encoded (non-human readable) version of the visual patterns found in the satellite images, which, when coupled with machine learning approaches, is used to model the built environment of a given census tract [26]. For each census tract, we calculated the mean of the latent space feature representation. We performed feature extraction on a NVIDIA Tesla T4 GPU using the PyTorch package in Python. Finally, the latent space feature representation was regressed against the COVID-19 risk score variance using gradient boosted decision trees [43]. We deployed existing AlexNet deep learning models originally trained on images from the internet and fine-tuned [44] them to predict the variance associated with the COVID-19 risk score, framing the analysis as a regression task. To do this, we split the census tract data set (with the split being fully randomized) into 80:20 and 50:50 training and testing groups to get a conservative estimate of variance explained and predictive capability of the sociodemographic variables in the COVID-19 risk score while not overfitting the data. To train the model, we used a maximum tree depth of 5, a subsample of 80% of the features per tree, a learning rate (ie, feature weight shrinkage for each boosting step) of 0.1, and used threefold cross-validation to determine the optimal number of boosted trees. Training was completed on a NVIDIA Tesla T4 GPU using Python 3.7.7 and the XGBoost package. In a separate analysis, both satellite image features and the social determinants of health features (previously mentioned) were regressed against the COVID-19 risk score variance. We reported R^2 for the predictions in the test data (Figure 1G, 1H).

Association of the COVID-19 Community Risk Score With Zip Code–Level COVID-19–Attributed Mortality

We downloaded case and death count data on a zip code tabulation area (ZCTA) of New York City, a hot spot of the US COVID-19 epidemic as of May 20, 2020, and then again on September 20, 2020 (Figure 1J). We used 2010 census crossover files to map census tracts to ZCTAs. We mapped the COVID-19 risk score to each ZCTA in New York City in April and September 2020. Each ZCTA had information on the total number of COVID-19 tests, positive cases, and COVID-19–related deaths. We computed the average COVID-19 risk score for the ZCTA, weighting the average by population size of the census tract. As previously mentioned,

we estimated the ZCTA-level socioeconomic values and proportions. We associated the COVID-19 risk score with the death rate using a negative binomial model. We set the offset term as the logarithm of the total population size of a zip code. The exponentiated coefficients are interpreted as the incidence rate ratio for a unit change (eg, 1 SD increase) in the variable (vs no change). We also examined multicollinearity, calculating the variance inflation factor (VIF) using the VIF function in the *regclass* package in R.

Data Availability Through the COVID-19 Risk Score Application Programming Interface and Dashboard

Finally, the COVID-19 risk score was made publicly available through an application programming interface and online web dashboard (see [Multimedia Appendix 1](#)).

Ethics Approval

Ethics approval was not required for this investigation as the study did not involve any human participants, and all of the data used were obtained from publicly available data sets.

Results

Prevalence and Heterogeneity of COVID-19–Associated Comorbidities and Risk Factors Across 500 Cities of the United States

We present summary statistics of the prevalence of the 15 COVID-19 comorbidities and risk factors for 27,648 census

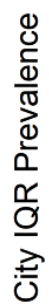
tracts across the United States using the 2019 release of the CDC 500 Cities data (derived from data obtained in 2017) and ACS data collected between 2013 and 2017 ([Figure 1A, 1C](#)). Census tracts represent small *communities* that have a median population size of 4091 (total range of 15–51,536). From the 500 cities analyzed, there was a median number of 28 (IQR 20–47) census tracts, with the most tracts found in New York (2141 tracts, with a population of $n=8,440,712$), Los Angeles (992 tracts, with a population of $n=3,961,681$), and Chicago (794 tracts, with a population of $n=2,726,431$), while Meridian, Idaho (4 tracts, with a population of $n=53,442$) has the fewest number of tracts. There was a wide range of prevalence values (ranging from 6% to 100%; [Figure 1D, Figure 2](#)) for the different prevalence measures, and a wide range of IQR values within cities was noted ([Figure 3](#) and Tables S1 and S2 in [Multimedia Appendices 2 and 3](#)).

Atlanta had the greatest IQR for obesity (22%–40%), high blood pressure (20%–44%), and COPD (4%–9%), while Gainesville had the highest variation in prevalence of high cholesterol (18%–34%) and blood pressure medication (51%–74%).

Figure 2. Per census tract prevalence for health indicators (y-axis). BP: blood pressure; COPD: chronic obstructive pulmonary disease.



	Kidney Disease	Stroke	Heart Disease	Cancer	COPD
1	0.000	0.000	0.000	0.000	0.000
2	0.000	0.000	0.000	0.000	0.000
3	0.000	0.000	0.000	0.000	0.000
4	0.000	0.000	0.000	0.000	0.000
5	0.000	0.000	0.000	0.000	0.000
6	0.000	0.000	0.000	0.000	0.000
7	0.000	0.000	0.000	0.000	0.000
8	0.000	0.000	0.000	0.000	0.000
9	0.000	0.000	0.000	0.000	0.000
10	0.000	0.000	0.000	0.000	0.000
11	0.000	0.000	0.000	0.000	0.000
12	0.000	0.000	0.000	0.000	0.000
13	0.000	0.000	0.000	0.000	0.000
14	0.000	0.000	0.000	0.000	0.000
15	0.000	0.000	0.000	0.000	0.000
16	0.000	0.000	0.000	0.000	0.000
17	0.000	0.000	0.000	0.000	0.000
18	0.000	0.000	0.000	0.000	0.000
19	0.000	0.000	0.000	0.000	0.000
20	0.000	0.000	0.000	0.000	0.000
21	0.000	0.000	0.000	0.000	0.000
22	0.000	0.000	0.000	0.000	0.000
23	0.000	0.000	0.000	0.000	0.000
24	0.000	0.000	0.000	0.000	0.000
25	0.000	0.000	0.000	0.000	0.000
26	0.000	0.000	0.000	0.000	0.000
27	0.000	0.000	0.000	0.000	0.000
28	0.000	0.000	0.000	0.000	0.000
29	0.000	0.000	0.000	0.000	0.000
30	0.000	0.000	0.000	0.000	0.000
31	0.000	0.000	0.000	0.000	0.000
32	0.000	0.000	0.000	0.000	0.000
33	0.000	0.000	0.000	0.000	0.000
34	0.000	0.000	0.000	0.000	0.000
35	0.000	0.000	0.000	0.000	0.000
36	0.000	0.000	0.000	0.000	0.000
37	0.000	0.000	0.000	0.000	0.000
38	0.000	0.000	0.000	0.000	0.000
39	0.000	0.000	0.000	0.000	0.000
40	0.000	0.000	0.000	0.000	0.000
41	0.000	0.000	0.000	0.000	0.000
42	0.000	0.000	0.000	0.000	0.000
43	0.000	0.000	0.000	0.000	0.000
44	0.000	0.000	0.000	0.000	0.000
45	0.000	0.000	0.000	0.000	0.000
46	0.000	0.000	0.000	0.000	0.000
47	0.000	0.000	0.000	0.000	0.000
48	0.000	0.000	0.000	0.000	0.000
49	0.000	0.000	0.000	0.000	0.000
50	0.000	0.000	0.000	0.000	0.000
51	0.000	0.000	0.000	0.000	0.000
52	0.000	0.000	0.000	0.000	0.000
53	0.000	0.000	0.000	0.000	0.000
54	0.000	0.000	0.000	0.000	0.000
55	0.000	0.000	0.000	0.000	0.000
56	0.000	0.000	0.000	0.000	0.000
57	0.000	0.000	0.000	0.000	0.000
58	0.000	0.000	0.000	0.	



The Pearson correlations between the 15 different prevalence values was calculated using census tract-level data (Figure 1D, Figure 4), with a median absolute value of correlation of 0.63 (IQR 0.35-0.78) noted with disease prevalences. The mean pairwise correlation between cardiometabolic diseases (diabetes, stroke, and heart disease) was 0.92, for cardiovascular risk factors (obesity, high blood pressure, and high cholesterol) was 0.62, and for smoking and respiratory conditions (asthma and COPD) was 0.69. An average correlation of 0.78 existed for diseases like diabetes, stroke, and heart disease, with obesity highly correlated with all of them (mean correlation 0.54), and

The first two principal components of the 15 COVID-19 health indicators and risk factors described 85% of the total variation (61% and 24% for component 1 and 2, respectively, see [Figure 5](#)) of the variation over all 27,648 census tracts ([Figure 1E](#)). The first principal component had equal contribution from all 15 health indicators and risk factors, except for cancer and males and females older than 65 years; the second principal component was dominated by cancer and age (Table S3 in [Multimedia Appendix 2](#)). This pattern of health indicator and risk factor contribution to principal components was also noted when the COVID-19 risk score was calculated at the city and county level (Table S3 in [Multimedia Appendix 2](#)).

Figure 4. Pearson correlation of health indicators across 27,648 census tracts (legend value corresponds to Pearson correlation value). BP: blood pressure; COPD: chronic obstructive pulmonary disease.

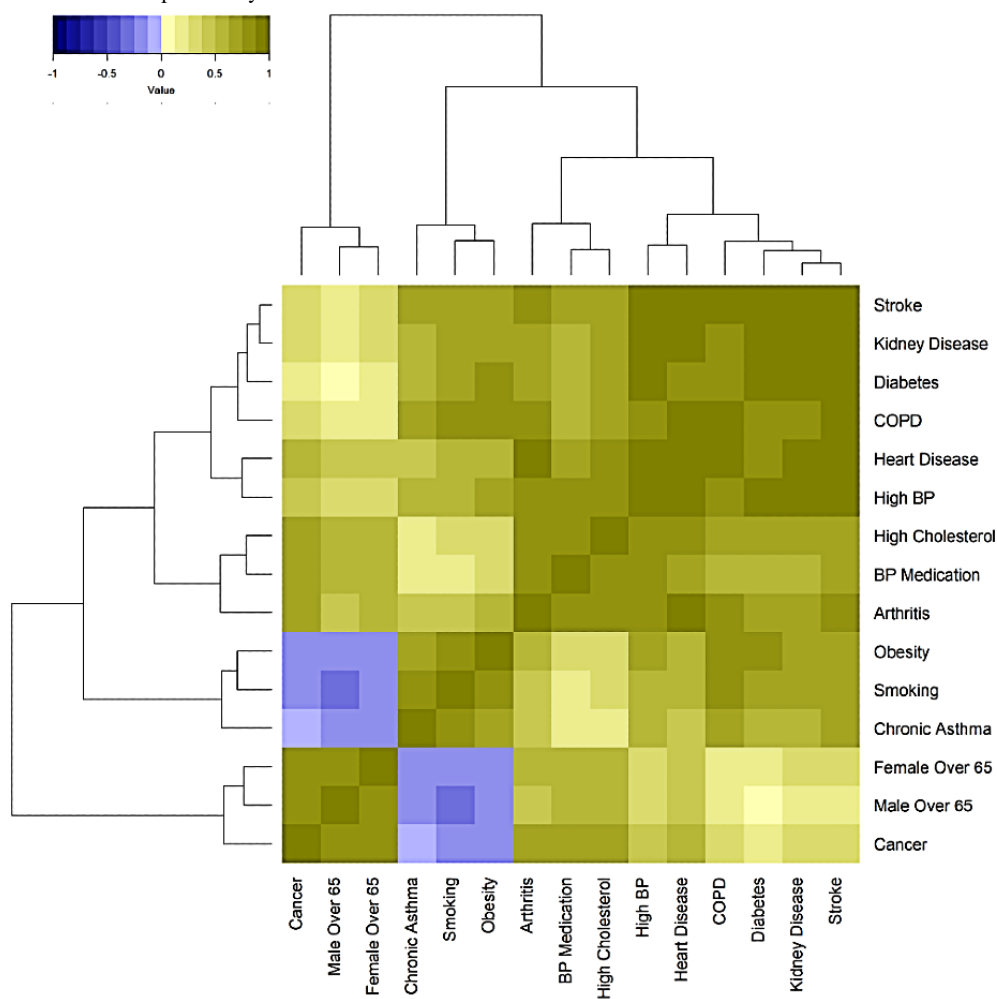
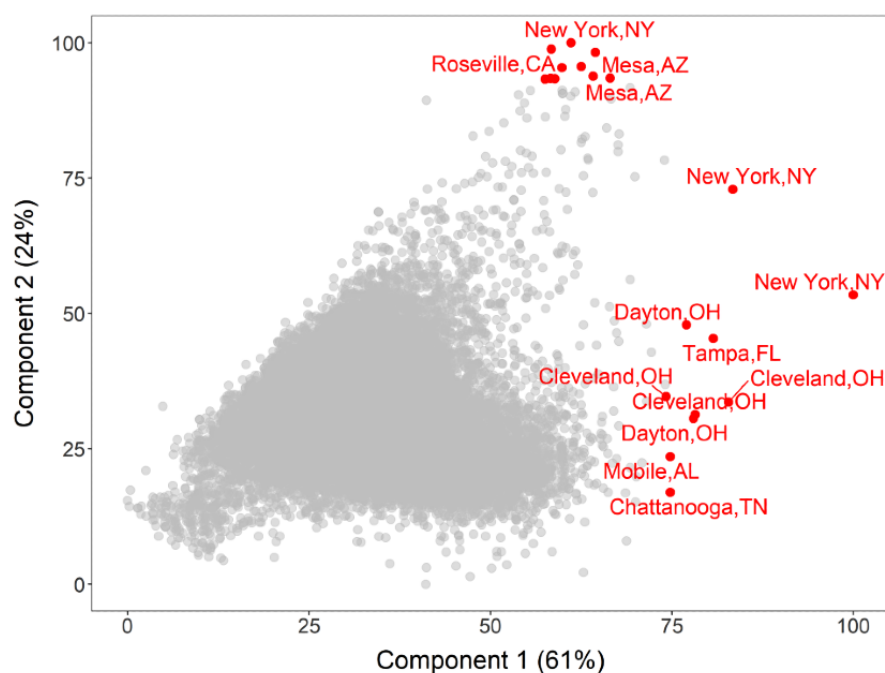


Figure 5. Scatterplot showing the relationship between the first and second principal components from principal component analysis, with each point indicating a city or census tract in the United States (top 10 cities/tracts by principal component 1 or 2 are highlighted in red).



Calculating a Robust COVID-19 Community Risk Score

The COVID-19 risk score was calculated using the 15 disease and health indicators for 27,648 included census tracts. The

average score was 33.7 (SD 8.6); the median was 33.32 (IQR 28-38). [Table 1](#) shows the communities with the highest variation of scores in the United States. The average error of the COVID-19 risk score across the census tracts was 1.25 (SD 0.85).

Table 1. Cities with the largest variation of COVID-19 risk score.

City, State	Median	Min	Max	25th percentile	75th percentile	SD	IQR
Athens, GA	32.2	6.3	42.6	20.7	35.7	10.0	15.0
Atlanta, GA	33.7	4.7	53.7	23.8	41.4	11.0	17.6
Boynton Beach, FL	41.7	21.6	81.3	35.5	52.1	16.2	16.6
Champaign, IL	29.8	3.4	45.2	17.9	34.1	12.8	16.2
Gainesville, FL	27.1	2.2	75.2	16.6	38.5	15.0	21.9
Hemet, CA	41.2	30.1	67.9	36.4	53.4	11.1	17.0
Mesa, AZ	32.2	7.7	84.6	29.0	43.2	14.7	14.2
Montgomery, AL	41.0	15.0	61.3	33.5	47.6	9.8	14.1
St. Louis, MO	36.9	22.0	57.3	31.7	46.3	8.6	14.6
Surprise, AZ	30.2	24.1	77.9	26.1	58.7	20.6	32.6
Birmingham, AL	43.9	18.8	57.9	36.4	49.3	9.6	12.9
Cape Coral, FL	43.4	30.3	63.3	37.3	49.3	8.2	12.0
Clearwater, FL	42.9	28.7	66.4	39.6	51.5	8.0	12.0
Cleveland, OH	42.4	18.3	78.3	38.0	49.6	9.1	11.6
Dayton, OH	43.1	6.0	78.1	38.6	49.8	11.8	11.2
Huntsville, AL	42.7	22.1	56.3	32.9	45.6	8.4	12.7
Lake Charles, LA	41.5	27.4	54.0	36.4	46.5	7.5	10.1
Lakeland, FL	43.9	18.0	65.8	38.3	49.2	10.9	10.9
Largo, FL	45.0	26.1	75.3	41.1	53.2	11.6	12.2
Palm Coast, FL	46.8	33.9	58.1	43.7	54.8	7.5	11.0
Pompano Beach, FL	43.6	27.1	64.3	37.0	48.5	9.7	11.5
Shreveport, LA	42.8	21.7	64.0	37.7	49.7	8.6	12.1
Gary, IN	50.8	42.5	61.8	47.1	54.6	5.2	7.6

COVID-19 Community Risk Score Variance Can Be Explained by Social Determinants of Health and Satellite Images of the Built Environment

The social determinants of health measures (excluding built environment) and demographic characteristics of a community ([Figure 1C](#), 1H) explain 54% of the total additive variation calculated using multiple linear regression ($r^2=0.54$; $P<.001$) of the COVID-19 risk score in the testing data set (when using a 50:50 and 80:20 fully randomized training:testing split). In this regression analysis, low to moderate multicollinearity was found with VIFs ranging from 1.41 for the variable *not employed* to 4.71 for *less than high school*. We found an additional 11% of variation attributed to nonlinear relationships, or a total of 65% between social determinants and the COVID-19 risk score, in the testing data using random forest-based regression ($r^2=0.65$; $P<.001$). The built environment features captured by satellite images contributed to 27% of the variation in the COVID-19 risk score. In total, combining both social

determinants and satellite imagery explained 87% of the variation of the COVID-19 risk score when using an 80:20 training:testing split ([Figure 1G](#), 1H).

Concerning important features, all 13 sociodemographic variables correlated with the COVID-19 risk score (linear regression $P<.001$ for 11 out of 13 variables) illustrated in [Table 2](#). The variables that had the largest additive contribution included the proportion of the community that was nonemployed (for a 1 SD change in proportion of nonemployed was associated with a 5.3 unit increase in the COVID-19 score; $P<.001$). A 1 SD increase in the increase of individuals with less than a high school education was associated with a 2 unit increase in the score. However, a 1 SD change in the increase of those at or below the poverty level was associated with a 3.3 unit decrease in the COVID-19 risk score. We found low to moderate VIFs associated with each sociodemographic variable ([Table 2](#)).

When assessing the explained variance using nonlinear regression (random forest) methods, the *most important*

variables in the training data (ascertained through a permutation of each variable sequentially) included the proportion of the tract that was not employed (273% increase of mean squared error [MSE] when permuted), of Asian ethnicity (93% increase of MSE), at or below poverty (91% increase of MSE), Hispanic

(78% increase MSE), and less than high school (78% increase MSE). The rank order of the importance of these variables was similar to the strength of their association in the linear model (Table 2). The same results were observed when the training:testing split was 50:50 and 80:20.

Table 2. Multivariate coefficients and CIs for linear regression and random forest regression of the COVID-19 risk score.

Variable	Linear coefficient	P value	Low (95% CI)	High (95% CI)	MSE ^a	Node purity ^b	VIF ^{c,d}
Median income	-1.34	<.001	-1.53	-1.16	42	59,736	3.68
Median home value	-0.13	.07	-0.27	0.01	39	33,163	2.21
At or below poverty (%)	-3.24	<.001	-3.42	-3.07	61	78,890	3.04
Unemployment (%)	0.73	<.001	0.60	0.86	87	68,364	1.69
Nonemployed (%)	5.38	<.001	5.26	5.50	285	316,903	1.42
Less than high school (%)	2.12	<.001	1.90	2.33	71	63,048	4.71
No health insurance (%)	0.69	<.001	0.55	0.83	50	34,818	2.18
More than 1 occupant (%)	-0.89	<.001	-1.04	-0.73	59	41,387	2.46
African American (%)	0.73	<.001	0.59	0.87	68	84,497	2.09
Hispanic (%)	-2.30	<.001	-2.49	-2.10	78	63,847	4.12
Asian (%)	-1.14	<.001	-1.25	-1.02	91	93,675	1.42
Other race (%)	-0.51	<.001	-0.67	-0.36	69	45,301	2.45

^aMSE: mean standard error.

^bNode impurity: residual sum of squares for the random forest model.

^cVIF: variance inflation factor.

^dFor the linear regression model.

COVID-19 Community Risk Score Was Associated With COVID-19 Death Rate in New York City

A 1 SD increase in the COVID-19 risk score was associated with a 40% increase in the incident rate ratio (IRR 1.40 per 1 SD increase; $P<.001$; Figure 6 and Table 3) in both May and

September 2020. For zip codes (eg, Figure 6 annotated zip codes) that had COVID-19 risk scores greater than 40, there was an almost twofold increase in death rates (IRR 1.98, 95% CI 1.43-2.77; $P<.001$). Additionally, we assessed multicollinearity by calculating the VIFs for each variable and found moderate to high multicollinearity.

Figure 6. COVID-19 deaths as a function of the COVID-19 risk score in New York City for each zip code (middle panel). The zip codes with the highest and lowest death rates are annotated. Blue points denote data on the epidemic death counts in September 2020. Red points denote epidemic death counts in May 2020.



Table 3. Multivariate incidence rate ratios (for 1 SD change in the variable) for zip code–level deaths in New York City in May and September 2020.

Variable (per 1 SD unit)	May IRR ^a (95% CI)	May <i>P</i> value	VIF ^b	September IRR (95% CI)	September <i>P</i> value	VIF
COVID-19 risk score	1.40 (1.27-1.55)	<.001	2.20	1.40 (1.27-1.53)	<.001	2.20
Median income	1.02 (0.84-1.22)	.80	9.06	0.99 (0.82-1.18)	.90	9.12
Less than high school	0.81 (0.008-1.81)	.10	19.80	0.81 (0.62-1.06)	.10	19.64
College educated	0.93 (0.26-1.92)	.50	10.83	0.93 (0.76-1.14)	.50	10.77
African American	1.14 (1.03-2.78)	.03	3.91	1.16 (1.03-1.30)	.01	3.95
Mexican	0.9 (0.87-1.08)	.60	3.72	0.97 (0.87-1.07)	.50	3.73
Hispanic	1.27 (1.19-1.46)	<.001	5.60	1.29 (1.12-1.47)	<.001	5.60
Asian	1.12 (1.00-1.26)	.05	4.34	1.15 (1.02-1.28)	.02	4.34
At or below poverty	1.04 (0.87-1.25)	.60	8.94	0.99 (0.83-1.17)	.90	8.92
More than 1 occupant per room	1.12 (1.00-1.27)	.06	4.83	1.10 (0.98-1.23)	.10	4.71
No health insurance	1.02 (0.91-1.16)	.70	4.66	1.03 (0.91-1.16)	.70	4.68
Unemployment	1.01 (0.91-1.13)	.80	3.36	1.02 (0.91-1.14)	.70	3.36
COVID-19 case count	1.08 (0.97-1.21)	.10	2.94	1.09 (0.98-1.21)	.10	2.88

^aIRR: incidence rate ratio.^bVIF: variance inflation factor.

Discussion

Principal Results

In this multi-scale analysis integrating and comparing spatial disease information from gold standard disease prevalence sources such as the US CDC, social determinants of health information from the US census, and satellite imagery data, we demonstrate an approach to identify characteristics of communities at risk for COVID-19 complications. We used the tools of unsupervised learning to develop a COVID-19 risk score that provides a single interpretable number that summarizes a communities' (census tract) aggregate risk. The constituents of the COVID-19 risk score included census tract–level chronic disease risk factors that corresponded to previously identified individual-level risk factors for COVID-19, such as age, obesity, diabetes, and heart disease.

Others have deployed similar risk scores to identify communities at risk for COVID-19 [16] and have used social determinants of health to identify this risk [45-47]. Furthermore, we were inspired by the work of others that demonstrate how remote sensing images predict obesity prevalence [26]. However, to our knowledge, this is the first study to examine the relationship between COVID-19 risk in neighborhoods (quantified using the COVID-19 risk score) and the social determinants of health and satellite image information. We found that, by combining established social determinants, information measured on earth with the built environment from space can explain most of the variation in the COVID-19 risk score, with a mere 13 sociodemographic variables explaining 50% of variation and, when combined with satellite images, could explain ~90% of variation. As more COVID-19 data becomes available, this finding suggests that future risk models for COVID-19 could incorporate satellite imagery together with social determinants of health to better model risk. Currently, comprehensive

measurement of the built environment is not typically used in the public health response to outbreaks, and COVID-19 pandemic risk models are typically modeled at the county level [46,47], a coarse geographical resolution that can obscure local hot spots or areas of need. Building models using the approach outlined here could help facilitate precision public health responses down to the local community (census tract) or subcensus tract level, thereby facilitating more precise allocations of resources to areas that need it.

Although it could be argued that the deep learning analysis of satellite imagery is simply a measurement of population density, this approach also measures several other factors that may contribute to COVID-19 infection and death rates independent of population density, such as built environment features that contribute to the development of COVID-19 risk factors and features that may put individuals at risk of contracting COVID-19. Examples of features that may put individuals at risk for developing risk factors include walkability (which contributes to obesity [48]) and road proximity (which can increase risk for heart disease [49]). Additionally certain architectural and built environment features that might put individuals at risk of COVID-19 infection, such as the configuration of pedestrian traffic in an urban area [50], can be partly quantified with this approach.

We believe that the COVID-19 risk score can be a tool in the growing armamentarium for public health and health care companies' toolbox to enable communities to prepare for the potential onslaught of cases in the coming winter months, ultimately helping to "flatten the curve" [51] and achieve precision public health goals of improving local health. Notably, we found that the zip code–level COVID-19 risk score for New York City and surrounding areas predicted risk for COVID-19 complications such as death. Zip codes with the highest COVID-19 scores (in the top 5%) had double the risk of COVID-19 death versus zip codes with the lowest scores. As

of this writing, New York City is contemplating another lockdown due to a surge in the same zip codes we identified as high risk [52]. Given the heterogeneity of various census tracts and neighborhoods across the United States and the range of COVID-19 rates and deaths, a more comprehensive national analysis will need to be performed using nationally representative comorbidity data and satellite data before extending the conclusions from the New York City analysis to similar jurisdictions in the United States or across the whole country.

As a byproduct of developing a risk score for communities, we observed that there is substantial variation of chronic disease prevalence within cities and across cities in the United States. With the exception of New York City and a few other places in the United States, public health agencies mostly collect COVID-19 case and death records at the county level across the country. However, the findings in our study implicate that smaller populations are at risk, and counties are heterogeneous.

We demonstrated how COVID-19 rates can be modeled using the COVID-19 risk score and how social determinants of health and the built environment can explain most of the score variance. Through simulations of the coprevalences of each of the 27,648 census tracts, we found that the point estimates for the community risk scores were robust to simulated sampling error. Many cities in the southwest and southeast demonstrated wide ranges in the COVID-19 risk score values. For example, Surprise, Arizona had a COVID-19 risk score with an IQR of 26 to 59. Atlanta, Georgia had an IQR of 24 to 41 (Figure S1 in [Multimedia Appendix 2](#)). Social determinants of health are hierarchical in structure and distributed over both geographic space and time whose measurement can occur on both the individual level (exposure of a person) or area level (exposure levels of a place). Satellite images provide a microscope into the area-level built environment, a concept that encapsulates the physical structures of how humans live, such as the city layout, resource presence, and landscape. A total of 65% of COVID-19 community risk score variance was explained by demographics and the social determinants of health, and 87% explained when the built environment was included. Given the large proportion of variance explained by the built environment, future precision public health strategies like hot spot identification and vaccine prioritization could be quickly improved by including measurements of the built environment to identify geographical areas in need of assistance.

This large proportion of COVID-19-associated risk variance explained by the social determinants of health and built

environment may be partly due to how discrimination affects where people live, their built environment, and access to health care [15,53,54]. Since the built environment and social determinants of health were found to play an important role in explaining the variance associated with the COVID-19 risk score, we plan to integrate this information into future COVID-19 risk score calculations that can be extended across the United States beyond the 500 Cities data set. We found that ~90% of the variation of prevalence of the 15 disease and health indicator prevalences (eg, diabetes, obesity, cardiovascular disease, populations that take blood pressure medication, and average age) can be explained by just two dimensions.

Limitations

The following are limitations of this study. First, we relied on disease and health indicator prevalence from the 500 largest cities in the United States but missed out on less urban areas whose populations are at risk for COVID-19 complications. In the future, we aim to task satellite imaging technology to locations that cannot be covered by resource-limited public surveillance programs. Second, although the CDC 500 Cities data are reflective of the diversity of individuals who live in a census tract, they are updated every 2 years and are dated to the latest collection (2019 data release reflects disease prevalence in 2017). Relatedly, neither individual-level disease nor COVID-19 status of individuals from these communities are measured. Third, satellite image data are captured at a resolution of approximately 20 m per pixel. It is not clear from our study if higher resolution images (that can theoretically capture more human-visible details of the built environment) would lead to better predictions of the COVID-19 risk score. Finally, interpretations of the New York-related data is limited due to the fact that it is aggregated to the zip code level. It is clear that COVID-19 is a disease of disparity; however, we cannot make a causal claim between the instruments such as the COVID-19 risk score, satellite imagery, and census tract-level sociodemographic factors, and eventual individual-level COVID-19-related complications.

Conclusions

Although it is clear that individual-level comorbidities are associated with risk for COVID-19, here we show that communities' clinical coprevalence structure are predictive of risk quantified by the COVID-19 risk score, and the variance of that score can be explained using the social determinants of health and the built environment measured from satellite imagery. We provide all our tools to monitor COVID-19 risk and related data in an interactive web-based dashboard.

Acknowledgments

We thank Emmanuel Coloma and Sumeet Parekh for their help in crafting the figures. This study is funded by XY Health, Inc, a company that develops machine learning approaches for prediction of health outcomes using satellite and land sensor data.

Conflicts of Interest

This study was funded by XY Health, Inc, and AD, GL, CL, and WDB completed this study while employed at XY Health Inc.

Multimedia Appendix 1

Supplementary methods.

[DOCX File, 16 KB - [publichealth_v7i8e26604_app1.docx](#)]

Multimedia Appendix 2

Supplementary information containing figures, tables, and application programming interface specification for the COVID-19 community risk score.

[DOCX File, 735 KB - [publichealth_v7i8e26604_app2.docx](#)]

Multimedia Appendix 3

Table S2.

[XLSX File (Microsoft Excel File), 489 KB - [publichealth_v7i8e26604_app3.xlsx](#)]

References

- Munshi L, Hall JB. Respiratory support during the COVID-19 pandemic: is it time to consider using a helmet? JAMA 2021 May 04;325(17):1723-1725. [doi: [10.1001/jama.2021.4975](#)] [Medline: [33764370](#)]
- Nuzzo J, Moss B, Watson C, Rutkow L. COVID-19 Dashboard by the Center for Systems Science and Engineering (CSSE) at Johns Hopkins University (JHU). Johns Hopkins Coronavirus Resource Center. 2021. URL: <https://coronavirus.jhu.edu/map.html> [accessed 2021-07-06]
- CDC COVID-19 Response Team. Severe outcomes among patients with coronavirus disease 2019 (COVID-19) - United States, February 12-March 16, 2020. MMWR Morb Mortal Wkly Rep 2020 Mar 27;69(12):343-346. [doi: [10.15585/mmwr.mm6912e2](#)] [Medline: [32214079](#)]
- Garg S, Kim L, Whitaker M, O'Halloran A, Cummings C, Holstein R, et al. Hospitalization rates and characteristics of patients hospitalized with laboratory-confirmed coronavirus disease 2019 - COVID-NET, 14 States, March 1-30, 2020. MMWR Morb Mortal Wkly Rep 2020 Apr 17;69(15):458-464. [doi: [10.15585/mmwr.mm6915e3](#)] [Medline: [32298251](#)]
- Gold JAW, Wong KK, Szablewski CM, Patel PR, Rossow J, da Silva J, et al. Characteristics and clinical outcomes of adult patients hospitalized with COVID-19 - Georgia, March 2020. MMWR Morb Mortal Wkly Rep 2020 May 08;69(18):545-550. [doi: [10.15585/mmwr.mm6918e1](#)] [Medline: [32379729](#)]
- Petrilli C, Jones S, Yang J, Rajagopalan H, O'Donnell L, Chernyak Y, et al. Factors associated with hospital admission and critical illness among 5279 people with coronavirus disease 2019 in New York City: prospective cohort study. BMJ 2020 May 22;369:m1966 [FREE Full text] [doi: [10.1136/bmj.m1966](#)] [Medline: [32444366](#)]
- Nandi A, Balasubramanian R, Laxminarayan R. Who is at the highest risk from COVID-19 in India? Analysis of health, healthcare access, and socioeconomic indicators at the district level Internet. medRxiv. Preprint posted online on June 9, 2020. [doi: [10.1101/2020.04.25.20079749](#)]
- Zhang X, Tan Y, Ling Y, Lu G, Liu F, Yi Z, et al. Viral and host factors related to the clinical outcome of COVID-19. Nature 2020 Jul;583(7816):437-440. [doi: [10.1038/s41586-020-2355-0](#)] [Medline: [32434211](#)]
- Yang X, Yu Y, Xu J, Shu H, Xia J, Liu H, et al. Clinical course and outcomes of critically ill patients with SARS-CoV-2 pneumonia in Wuhan, China: a single-centered, retrospective, observational study. Lancet Respir Med 2020 May;8(5):475-481 [FREE Full text] [doi: [10.1016/S2213-2600\(20\)30079-5](#)] [Medline: [32105632](#)]
- Grasselli G, Zangrillo A, Zanella A, Antonelli M, Cabrini L, Castelli A, COVID-19 Lombardy ICU Network. Baseline characteristics and outcomes of 1591 patients infected with SARS-CoV-2 admitted to ICUs of the Lombardy Region, Italy. JAMA 2020 Apr 28;323(16):1574-1581 [FREE Full text] [doi: [10.1001/jama.2020.5394](#)] [Medline: [32250385](#)]
- Zhou F, Yu T, Du R, Fan G, Liu Y, Liu Z, et al. Clinical course and risk factors for mortality of adult inpatients with COVID-19 in Wuhan, China: a retrospective cohort study. Lancet 2020 Mar 28;395(10229):1054-1062 [FREE Full text] [doi: [10.1016/S0140-6736\(20\)30566-3](#)] [Medline: [32171076](#)]
- Barker LE, Kirtland KA, Gregg EW, Geiss LS, Thompson TJ. Geographic distribution of diagnosed diabetes in the U.S.: a diabetes belt. Am J Prev Med 2011 Apr;40(4):434-439. [doi: [10.1016/j.amepre.2010.12.019](#)] [Medline: [21406277](#)]
- Krieger N, Chen JT, Waterman PD, Soobader M, Subramanian SV, Carson R. Geocoding and monitoring of US socioeconomic inequalities in mortality and cancer incidence: does the choice of area-based measure and geographic level matter?: the Public Health Disparities Geocoding Project. Am J Epidemiol 2002 Sep 01;156(5):471-482. [doi: [10.1093/aje/kwf068](#)] [Medline: [12196317](#)]
- Chin T, Kahn R, Li R, Chen JT, Krieger N, Buckee CO, et al. U.S. county-level characteristics to inform equitable COVID-19 response. medRxiv. Preprint posted online on April 11, 2020. [doi: [10.1101/2020.04.08.20058248](#)] [Medline: [32511610](#)]
- Figueroa JF, Wadhwa RK, Lee D, Yeh RW, Sommers BD. Community-level factors associated with racial and ethnic disparities in COVID-19 rates in Massachusetts. Health Aff (Millwood) 2020 Nov;39(11):1984-1992. [doi: [10.1377/hlthaff.2020.01040](#)] [Medline: [32853056](#)]
- Jin J, Agarwala N, Kundu P, Harvey B, Zhang Y, Wallace E, et al. Individual and community-level risk for COVID-19 mortality in the United States. Nat Med 2021 Feb;27(2):264-269. [doi: [10.1038/s41591-020-01191-8](#)]

17. Jones AC, Chaudhary NS, Patki A, Howard VJ, Howard G, Colabianchi N, et al. Neighborhood walkability as a predictor of incident hypertension in a national cohort study. *Front Public Health* 2021;9:611895. [doi: [10.3389/fpubh.2021.611895](https://doi.org/10.3389/fpubh.2021.611895)] [Medline: [33598444](https://pubmed.ncbi.nlm.nih.gov/33598444/)]
18. Lopez L, Hart LH, Katz MH. Racial and ethnic health disparities related to COVID-19. *JAMA* 2021 Feb 23;325(8):719-720. [doi: [10.1001/jama.2020.26443](https://doi.org/10.1001/jama.2020.26443)] [Medline: [33480972](https://pubmed.ncbi.nlm.nih.gov/33480972/)]
19. Webb Hooper M, Nápoles AM, Pérez-Stable EJ. COVID-19 and racial/ethnic disparities. *JAMA* 2020 Jun 23;323(24):2466-2467. [doi: [10.1001/jama.2020.8598](https://doi.org/10.1001/jama.2020.8598)] [Medline: [32391864](https://pubmed.ncbi.nlm.nih.gov/32391864/)]
20. Khazanchi R, Evans CT, Marcelin JR. Racism, not race, drives inequity across the COVID-19 continuum. *JAMA Netw Open* 2020 Sep 01;3(9):e2019933 [FREE Full text] [doi: [10.1001/jamanetworkopen.2020.19933](https://doi.org/10.1001/jamanetworkopen.2020.19933)] [Medline: [32975568](https://pubmed.ncbi.nlm.nih.gov/32975568/)]
21. Berkowitz RL, Gao X, Michaels EK, Mujahid MS. Structurally vulnerable neighbourhood environments and racial/ethnic COVID-19 inequities. *Cities Health* 2020 Jul 29;1-4. [doi: [10.1080/23748834.2020.1792069](https://doi.org/10.1080/23748834.2020.1792069)]
22. Wu X, Nethery RC, Sabath MB, Braun D, Dominici F. Exposure to air pollution and COVID-19 mortality in the United States: a nationwide cross-sectional study. *medRxiv*. Preprint posted online on April 27, 2020. [doi: [10.1101/2020.04.05.20054502](https://doi.org/10.1101/2020.04.05.20054502)]
23. Travaglio M, Yu Y, Popovic R, Selley L, Leal NS, Martins LM. Links between air pollution and COVID-19 in England. *Environ Pollut* 2021 Jan 01;268(Pt A):115859 [FREE Full text] [doi: [10.1016/j.envpol.2020.115859](https://doi.org/10.1016/j.envpol.2020.115859)] [Medline: [33120349](https://pubmed.ncbi.nlm.nih.gov/33120349/)]
24. Cromer S, Lakhani C, Wexler D, Burnett-Bowie S, Udler M, Patel C. Geospatial analysis of individual and community-level socioeconomic factors impacting SARS-CoV-2 prevalence and outcomes. *medRxiv*. Preprint posted online on September 30, 2020. [doi: [10.1101/2020.09.30.20201830](https://doi.org/10.1101/2020.09.30.20201830)] [Medline: [33024982](https://pubmed.ncbi.nlm.nih.gov/33024982/)]
25. Social determinants of health. Office of Disease Prevention and Health Promotion. URL: <https://health.gov/healthypeople/objectives-and-data/social-determinants-health> [accessed 2021-05-12]
26. Maharana A, Nsoesie EO. Use of deep learning to examine the association of the built environment with prevalence of neighborhood adult obesity. *JAMA Netw Open* 2018 Aug 03;1(4):e181535 [FREE Full text] [doi: [10.1001/jamanetworkopen.2018.1535](https://doi.org/10.1001/jamanetworkopen.2018.1535)] [Medline: [30646134](https://pubmed.ncbi.nlm.nih.gov/30646134/)]
27. Krieger N, Waterman P, Chen JT, Soobader M, Subramanian SV, Carson R. Zip code caveat: bias due to spatiotemporal mismatches between zip codes and US census-defined geographic areas--the Public Health Disparities Geocoding Project. *Am J Public Health* 2002 Jul;92(7):1100-1102. [doi: [10.2105/ajph.92.7.1100](https://doi.org/10.2105/ajph.92.7.1100)] [Medline: [12084688](https://pubmed.ncbi.nlm.nih.gov/12084688/)]
28. XY.ai COVID-19 Community Risk Score Dashboard. URL: <https://covid19satellite.org/> [accessed 2021-02-01]
29. Centers for Disease Control and Prevention. URL: <https://www.cdc.gov/places/about/500-cities-2016-2019/index.html> [accessed 2020-10-04]
30. Zhang X, Holt JB, Lu H, Wheaton AG, Ford ES, Greenlund KJ, et al. Multilevel regression and poststratification for small-area estimation of population health outcomes: a case study of chronic obstructive pulmonary disease prevalence using the behavioral risk factor surveillance system. *Am J Epidemiol* 2014 Apr 15;179(8):1025-1033. [doi: [10.1093/aje/kwu018](https://doi.org/10.1093/aje/kwu018)] [Medline: [24598867](https://pubmed.ncbi.nlm.nih.gov/24598867/)]
31. People with certain medical conditions. Centers for Disease Control and Prevention. 2020. URL: <https://www.cdc.gov/coronavirus/2019-ncov/need-extra-precautions/people-with-medical-conditions.html> [accessed 2020-12-12]
32. American Community Survey Data. United States Census Bureau. URL: <https://www.census.gov/programs-surveys/acs/data.html> [accessed 2020-10-04]
33. Boyd RW, Lindo EG, Weeks LD, McLemore MR. On racism: a new standard for publishing on racial health inequities. *Health Affairs*. URL: <https://www.healthaffairs.org/doi/10.1377/hblog20200630.939347> [accessed 2021-05-07]
34. Maroko AR, Nash D, Pavilonis BT. COVID-19 and inequity: a comparative spatial analysis of New York City and Chicago hot spots. *J Urban Health* 2020 Aug;97(4):461-470 [FREE Full text] [doi: [10.1007/s11524-020-00468-0](https://doi.org/10.1007/s11524-020-00468-0)] [Medline: [32691212](https://pubmed.ncbi.nlm.nih.gov/32691212/)]
35. The R Project for Statistical Computing. 2017. URL: <http://www.R-project.org/> [accessed 2020-10-10]
36. CSSEGISandData / COVID-19. GitHub. URL: <https://github.com/CSSEGISandData/COVID-19> [accessed 2020-10-04]
37. Cartographic boundary files - shapefile. United States Census Bureau. URL: <https://www.census.gov/geographies/mapping-files/time-series/geo/carto-boundary-file.html> [accessed 2020-12-17]
38. Patel CJ. xyhealth / covid_comorbidity_score. GitHub. URL: https://github.com/xyhealth/covid_comorbidity_score [accessed 2020-12-16]
39. Klokkan Technologies GmbH. URL: <https://openmaptiles.com/> [accessed 2020-10-04]
40. Python. URL: <https://www.python.org/> [accessed 2021-05-14]
41. Krizhevsky A. One weird trick for parallelizing convolutional neural networks. *arXiv*. 2014. URL: <http://arxiv.org/abs/1404.5997> [accessed 2020-09-14]
42. Krizhevsky A, Sutskever I, Hinton G. ImageNet classification with deep convolutional neural networks. *Commun ACM* 2017 May 24;60(6):84-90. [doi: [10.1145/3065386](https://doi.org/10.1145/3065386)]
43. Chen T, Guestrin C. XGBoost: a scalable tree boosting system. In: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 2016 Presented at: KDD '16; August 13-17, 2016; San Francisco, CA p. 785-794. [doi: [10.1145/2939672.2939785](https://doi.org/10.1145/2939672.2939785)]

44. Razavian A, Azizpour H, Sullivan J, Carlsson S. CNN features off-the-shelf: an astounding baseline for recognition. 2014 Presented at: 2014 IEEE Conference on Computer Vision and Pattern Recognition Workshops; June 23-28, 2014; Columbus, OH. [doi: [10.1109/cvprw.2014.131](https://doi.org/10.1109/cvprw.2014.131)]
45. Kim SJ, Bostwick W. Social vulnerability and racial inequality in COVID-19 deaths in Chicago. *Health Educ Behav* 2020 Aug;47(4):509-513 [FREE Full text] [doi: [10.1177/1090198120929677](https://doi.org/10.1177/1090198120929677)] [Medline: [32436405](https://pubmed.ncbi.nlm.nih.gov/32436405/)]
46. Marvel SW, House JS, Wheeler M, Song K, Zhou Y, Wright FA, et al. The COVID-19 Pandemic Vulnerability Index (PVI) Dashboard: monitoring county-level vulnerability using visualization, statistical modeling, and machine learning. *Environ Health Perspect* 2021 Jan;129(1):17701 [FREE Full text] [doi: [10.1289/EHP8690](https://doi.org/10.1289/EHP8690)] [Medline: [33400596](https://pubmed.ncbi.nlm.nih.gov/33400596/)]
47. Mehta M, Julaiti J, Griffin P, Kumara S. Early stage machine learning-based prediction of US county vulnerability to the COVID-19 pandemic: machine learning approach. *JMIR Public Health Surveill* 2020 Sep 11;6(3):e19446 [FREE Full text] [doi: [10.2196/19446](https://doi.org/10.2196/19446)] [Medline: [32784193](https://pubmed.ncbi.nlm.nih.gov/32784193/)]
48. Yang S, Chen X, Wang L, Wu T, Fei T, Xiao Q, et al. Walkability indices and childhood obesity: a review of epidemiologic evidence. *Obes Rev* 2021 Feb;22 Suppl 1:e13096 [FREE Full text] [doi: [10.1111/obr.13096](https://doi.org/10.1111/obr.13096)] [Medline: [33185012](https://pubmed.ncbi.nlm.nih.gov/33185012/)]
49. Gan WQ, Tamburic L, Davies HW, Demers PA, Koehoorn M, Brauer M. Changes in residential proximity to road traffic and the risk of death from coronary heart disease. *Epidemiology* 2010 Sep;21(5):642-649. [doi: [10.1097/EDE.0b013e3181e89f19](https://doi.org/10.1097/EDE.0b013e3181e89f19)] [Medline: [20585255](https://pubmed.ncbi.nlm.nih.gov/20585255/)]
50. Yao Y, Shi W, Zhang A, Liu Z, Luo S. Examining the diffusion of coronavirus disease 2019 cases in a metropolis: a space syntax approach. *Int J Health Geogr* 2021 Apr 29;20(1):17 [FREE Full text] [doi: [10.1186/s12942-021-00270-4](https://doi.org/10.1186/s12942-021-00270-4)] [Medline: [33926460](https://pubmed.ncbi.nlm.nih.gov/33926460/)]
51. COVID-19 contact tracing training and resources. Centers for Disease Control and Prevention. 2020. URL: <https://www.cdc.gov/coronavirus/2019-ncov/php/contact-tracing/strategies-to-reduce-spread.html> [accessed 2020-10-04]
52. Cook L. 9 NYC zip codes with coronavirus clusters face shutdown of schools, businesses Internet. PIX11. 2020. URL: <https://pix11.com/news/coronavirus/9-nyc-zip-codes-with-coronavirus-clusters-face-shutdown-of-schools-businesses/> [accessed 2021-05-12]
53. Alcendor DJ. Racial disparities-associated COVID-19 mortality among minority populations in the US. *J Clin Med* 2020 Jul 30;9(8):2442 [FREE Full text] [doi: [10.3390/jcm9082442](https://doi.org/10.3390/jcm9082442)] [Medline: [32751633](https://pubmed.ncbi.nlm.nih.gov/32751633/)]
54. Poteat T, Millett GA, Nelson LE, Beyrer C. Understanding COVID-19 risks and vulnerabilities among black communities in America: the lethal force of syndemics. *Ann Epidemiol* 2020 Jul;47:1-3 [FREE Full text] [doi: [10.1016/j.annepidem.2020.05.004](https://doi.org/10.1016/j.annepidem.2020.05.004)] [Medline: [32419765](https://pubmed.ncbi.nlm.nih.gov/32419765/)]

Abbreviations

ACS: American Community Survey
CDC: Centers for Disease Control and Prevention
COPD: chronic obstructive pulmonary disease
IRR: incident rate ratio
MSE: mean squared error
VIF: variance inflation factor
ZCTA: zip code tabulation area

Edited by T Sanchez; submitted 18.12.20; peer-reviewed by A Maharan, R Berkowitz; comments to author 13.01.21; revised version received 14.05.21; accepted 15.07.21; published 26.08.21.

Please cite as:

Deonarine A, Lyons G, Lakhani C, De Brouwer W
Identifying Communities at Risk for COVID-19–Related Burden Across 500 US Cities and Within New York City: Unsupervised Learning of the Copevalence of Health Indicators
JMIR Public Health Surveill 2021;7(8):e26604
URL: <https://publichealth.jmir.org/2021/8/e26604>
doi: [10.2196/26604](https://doi.org/10.2196/26604)
PMID: [34280122](https://pubmed.ncbi.nlm.nih.gov/34280122/)

©Andrew Deonarine, Genevieve Lyons, Chirag Lakhani, Walter De Brouwer. Originally published in *JMIR Public Health and Surveillance* (<https://publichealth.jmir.org>), 26.08.2021. This is an open-access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in *JMIR Public Health and Surveillance*, is properly

cited. The complete bibliographic information, a link to the original publication on <https://publichealth.jmir.org>, as well as this copyright and license information must be included.

Original Paper

Online News Coverage of the Sugar-Sweetened Beverages Tax in Malaysia: Content Analysis

Muhammad Faiz Mohd Hanim^{1*}, BDS, MSc; Budi Aslinie Md Sabri^{1*}, BDS, MSc, PhD, DDPHRCS; Norashikin Yusof^{1*}, BDS, DDPHRCS, MPA

Centre of Population Oral Health and Clinical Prevention Studies, Faculty of Dentistry, Universiti Teknologi MARA (UiTM), Sungai Buloh, Malaysia

* all authors contributed equally

Corresponding Author:

Norashikin Yusof, BDS, DDPHRCS, MPA

Centre of Population Oral Health and Clinical Prevention Studies

Faculty of Dentistry

Universiti Teknologi MARA (UiTM)

UiTM Sg. Buloh Campus

Sungai Buloh, 47000

Malaysia

Phone: 60 3 6126 6621

Email: nyshikin@gmail.com

Abstract

Background: In Malaysia, the Sugar-Sweetened Beverages (SSBs) tax was announced during the parliament's 2019 Budget Speech. The tax was slated to be enforced by April 2019 but was later postponed to July 2019. The announcement has since generated significant media coverage and public feedback.

Objective: This study presents a qualitative and quantitative cross-sectional study using netnography to examine how Malaysian online news articles responded to the SSBs tax after the announcement and postimplementation.

Methods: Online news articles published on popular online news platforms from November 2018 to August 2019 were downloaded using NCapture and imported into NVivo for analysis using the inductive approach and thematic content analysis following the initial SSBs implementation announcement.

Results: A total of 62 news articles were analyzed. Most of the articles positively portrayed the SSBs tax (46.8%) and highlighted its health impacts (76%). There were 7 key framing arguments identified in the articles. The positive arguments revolved around incentivizing manufacturers to introduce healthier products voluntarily, positive health consequences, the tax's impact on government revenue, and the use of the generated revenue toward beneficial social programs. The opposing arguments included increased operating costs to the manufacturer, the increased retail price of drinks, and how the SSBs tax is not a robust solution to obesity. The top priority sector considered in introducing the tax was the health perspective, followed by economic purposes and creating policies such as regulating the food and drinks industry.

Conclusions: The majority of online news articles positively reported the implementation of the SSBs tax in Malaysia. This suggests media played a role in garnering support for the health policy. As such, relevant bodies can use negative findings to anticipate and reframe counteracting arguments opposing the SSBs tax.

(*JMIR Public Health Surveill* 2021;7(8):e24523) doi:[10.2196/24523](https://doi.org/10.2196/24523)

KEYWORDS

sugar-sweetened beverages; obesity; taxes; media content analysis; public health policy; media content; public health; netnography; malaysia; budget

Introduction

Most noncommunicable diseases (NCD) such as cancer, cardiovascular diseases, dental caries, and diabetes share a common risk factor that can be attributed to high BMI and

obesity [1-5]. Previously, NCDs were primarily a problem in high-income countries but have become prevalent in low-income and middle-income countries [6]. One of the main factors contributing to the rising prevalence of NCDs is the increased consumption of sugar-sweetened beverages (SSBs) [6]. The

amount of sugar added to SSBs is quite substantial, and a 330 ml or 12 oz portion of a sugar-sweetened carbonated drink usually contains about 35 grams of sugar, which is equivalent to almost 9 teaspoons of sugar, offering little nutritional value [7].

The widespread availability of SSBs and their convenience have led to increased consumption worldwide [8,9]. Globally, consumption of SSBs is highest among young adults aged between 20-39 years [9]. Studies from Han [8] and Hazam [10] found that 40% to 70% of young adults in low-income countries consumed SSBs daily, especially in the form of energy drinks and soda or soft drinks.

One way to curb increased SSBs consumption is to introduce a tax or excise duty [11]. Most policymakers view this tax as one of the tools to tackle the obesity epidemic by limiting the choices available to consumers, encouraging manufacturers to limit the sugar content of their products, and creating a revenue stream to combat health-related problems caused by the SSBs [11,12].

Several countries have introduced SSBs taxation, but its implementation has received mixed reactions. Some groups oppose it, with health advocates who lobbied for SSBs taxation supporting it [11-14]. Opponents of SSBs taxation question its effectiveness in addressing the obesity problem and do not believe SSBs taxes alone can tackle obesity and overweight problems [13-16]. Nonetheless, SSBs taxation is gaining momentum because of its relative ease of implementation compared to other food or nutrition policy options [17]. Even though the United States has not implemented a nationwide tax, there have been some successes among several local and city governments that have adopted a similar tax [13]. For example, in Berkeley, California, the sugar tax resulted in a 21% decrease in soft drink consumption among the low-income neighborhoods [13,18]. While in Philadelphia, Pennsylvania, the SSBs consumption fell by 26% 2 months after the implementation of a beverage excise tax [18]. In the United Kingdom, a 3-tiered levy on SSBs was implemented in April 2018, in which the amount of levy imposed depends on the sugar content per 100 ml [19]. This has led to soft drink manufacturers reducing the sugar content of their products and a significant decrease (approximately 50%) in sales of soft drinks subjected to the levy [19].

In Australia, the final report by the "Select Committee into the Obesity Epidemic in Australia" recommended that the Australian government introduce a tax on SSBs with the objectives of reducing consumption and accelerating the reformulation of products [14]. Despite this recommendation, there was opposition, especially from the beverage industry, and there was even a consideration to reject its implementation. However, in 2018, the Australian Beverages Council announced it would seek to reduce the sugar content in beverages by 10% by 2020 and a further 10% by 2025 [13].

In Malaysia, the SSBs tax proposal levied 40 cents per 1 liter for beverages containing more than 5 grams of sugar per 100 ml and fruit juices with 12 grams of sugar per 100 ml. The tax was scheduled to be enforced in July 2019 [20].

Research has suggested that the media can influence public perception and the potential acceptance of public health policies [21-23]. Mass print or online media provide a comprehensive platform to reach out to the public, and they can effectively put health issues on the public agenda and determine how these issues are framed [23]. This framing process involves emphasizing some aspects of a debate while downplaying others. Framing can include representations of the severity of societal problems, their causes, and the potential effectiveness of proposed solutions. Thus, framing can influence the audience's perception and evaluation of different approaches to address public health problems [23]. Media framing of the SSBs tax within the context of sugar and SSBs consumption issues could influence the public acceptance of this upstream solution via taxation.

Policy articles posted on online media platforms play an important role in identifying supporting and opposing messages, and the challenges stakeholders may face when advocating for policy agendas such as SSBs taxation and tobacco control. A range of studies on the content analyses of media coverage on SSBs taxes and tobacco-related taxes have been published [21-28]. The strategic messaging by proponents of SSBs taxes can influence public opinion and shape policy development [22,25]. In a study conducted by Niederdeppe, Jeff et al [26] examined the news coverage of public debates on SSBs to illuminate how the news media frames these debates. Most of the debates framed the issue in favorable ways [26]. For example, the comprehensive media coverage generated by the Australian government's announcement of legislation mandating that tobacco products be sold in plain packaging provided an opportunity for tobacco control advocates to anticipate and counteract arguments opposed to the legislation [28]. Likewise, the media coverage of sugar-based taxes and SSBs consumption has helped shape policy to favor fiscal solutions that curb SSBs consumption and drive greater public acceptance of the sugar levy in the United Kingdom [22].

In Malaysia, many articles regarding the SSBs tax have been written in online news media. Media coverage of any particular issue is generally linked to the agenda or priorities of the government in power. However, during our study period, the media, especially the newspapers, were not linked to either government or the opposition, as reported in the Star Online, April 18, 2019. The article stated, "Malaysia has jumped 22 places to 123rd in the latest World Press Freedom Index, compiled by Reporters Without Borders (RSF)," becoming a top-ranked country among other ASEAN (Association of Southeast Asian Nations) countries [29]. In 2020, Malaysia's rank (101/180, showed that the Malaysian media were still enjoying their freedom. This showed that the articles written by the media are less likely to be biased but rather featured the media's accurate sentiment as reported through their online platforms [29]. Furthermore, the content of Malaysian-based online news platforms typically mirrors their paper or print-based counterparts, allowing for the generalizability of online news content analyses to the general media.

Since its announcement, no study has been conducted to examine the online news articles' responses regarding Malaysia's

SSBs tax. Thus, this study examines the response of Malaysian online news articles and how the media framed the arguments during the 10 months after its announcement and up to its implementation.

We also documented any sector-related and health-related issues used to justify the SSBs tax in the news reporting. All news outlets were captured to build a picture of which key framing arguments and messages were featured in the news articles.

Methods

Collection of Data

According to the reporting website, Malaysia's total number of online news platforms varies from 29 to 35 platforms [30,31]. In this study, the top 10 Malaysia-based online news platforms with total unique visitors in July 2019 (the latest information available) reported by the Malaysian Digital Association were purposefully selected to ensure a sample coverage representing various Malaysian online news readers [32]. As a result, the top Malaysia-based online news platform visited was a Malay version of the Harian Metro (Metro Daily), with 4.2 million unique visitors. Conversely, the online news platform with the lowest unique visitor rates included in the sample was the English newspaper version of The New Straits Time, which recorded 2.2 million unique visitors [32].

Data collection included news articles written between November 2018 to August 2019. Separate searches were conducted in the respective online news platforms. Suggestive keywords used were “sugar tax,” “sugar-sweetened beverages tax,” and “soda tax.” Suggestive keywords in other languages such as Malay, Mandarin, and Tamil were translated using Google translate. The news articles that appeared in the search results were scanned for relevance and downloaded using either the web browser extension NCapture (version 1.0.290.0; QSR International) and then imported into NVivo (version 12, QSR International). The articles were read, and to avoid duplication, the articles were excluded if the Malaysian online news platform did not produce the original articles.

Content Analysis

All news articles were read in full, and 2 researchers conducted the coding. The coding variables are shown in [Textbox 1](#). After the researchers separately coded the articles, the results were discussed and calibrated until a consensus was reached. News articles were coded and analyzed for the topic, framing arguments, overall slant, related sectors, health-related issues used to justify the SSBs tax, and direct quotes or position statements. Coding categories were developed iteratively, adapting a methodological approach similar to the one taken by Christina Watts and Becky Freeman [21]. Each article was coded differently, representing the primary message of the article.

Textbox 1. The coding variables used to code the arguments in the articles.

- **Topic:** Overall, what is the news article about? (1 topic per article coded).
- **Framing argument:** What is the argument presented concerning the introduction of the sugar-sweetened beverages (SSBs) tax in Malaysia? The framing argument was determined by identifying the argument presented most frequently within the articles.
- **Slant:** Is the article presenting the SSBs tax as a positive or negative policy? Or is the article neutral toward the tax? A positive slant is defined by a framing argument that favors SSBs taxation, and a negative slant is defined by a framing argument that is opposed to SSBs taxation. A neutral slant is defined by the absence of a framing argument and the presentation of a neutral debate.
- **Related sectors:** Which related sectors were mentioned or highlighted in the article as justification for SSBs taxation?
- **Health-related issues:** What type of health-related issues used as justification for SSBs taxation were mentioned or highlighted in the article?
- **Direct quotes or position statements:** Who is quoted or paraphrased in the article with a position or opinion on the SSBs tax?

Articles were also content analyzed and coded according to the most frequently presented argument in the article. Each framing argument was categorized as either in favor of, against, or neutral toward the SSBs tax. Articles were also analyzed for the justification used in introducing the tax with references to related sectors, health-related issues, and direct quotes or position statements about implementing the SSBs tax.

Results

Overview

Out of the 79 online news articles, 17 (21.5%) were excluded due to duplications (n=4, 23.5%), readers' letters, editors' opinions, and infographics (n=13, 76.5%). The final content analysis consisted of 62 news articles.

Topics and Overall Slant

The topics and overall slants are summarized in [Table 1](#). The most frequently highlighted topic concerned the impact of SSBs taxation on people's health (40.3%), followed by the implementation announcement (21%), and the contribution of the SSBs tax toward government revenue (6.5%). There was an overall positive slant towards SSBs taxation (46.8%), primarily highlighting positive health outcomes as the main impact of the SSBs tax ([Table 1](#)). On the other hand, there was a 16.1% negative slant, which includes the manufacturers' response towards the SSBs tax (33.3%) and its impact on the consumers (50%). However, online newspaper articles were neutral when reporting the SSBs tax announcement in Malaysia.

Table 1. Article themes by overall slant.

Topic	Brief description	Slant, n (%)			
		Positive	Negative	Neutral	Total
Announcement or implementation of the SSBs ^a tax	Articles report on the announcement or implementation of the SSBs tax	0 (0.0)	0 (0.0)	13 (100)	13 (100)
Health effects	Articles explore the impact of SSBs taxation on health	19 (76)	2 (8)	4 (16)	25 (100)
Government revenue	Taxation on SSBs will generate another source of income for the government to reduce the country's deficit	2 (50)	0 (0.0)	2 (50)	4 (100)
Manufacturer response	Articles report on how the food industry or beverage companies react to SSBs taxation	7 (58.4)	4 (33.3)	1 (8.3)	12 (100)
Consumer response	Articles report on the impact of the SSBs tax on the consumers	1 (12.5)	4 (50)	3 (37.5)	8 (100)
Total	N/A ^b	29 (46.8)	10 (16.1)	23 (37.1)	62 (100)

^aSSBs: sugar-sweetened beverages.

^bN/A: not applicable.

Framing Arguments

There were 7 key framing arguments used in the online news articles (Table 2). Framing arguments primarily supported taxation (46.8%), with 16.1% of the arguments opposing the taxation, and 37.1% framed as balanced arguments as they did not have any primary framing arguments. The majority of arguments supporting SSBs taxation reported the positive health gains of reducing SSBs consumption in the general public (21.7%). Some examples include nongovernmental organizations (NGOs) who applauded the SSBs tax because it will encourage healthy lifestyle behaviors and help reduce the incidence of NCDs (eg, diabetes) in the country [25,26].

The articles also argued that the tax incentivizes SSBs manufacturers to reduce the sugar content and introduce healthier products (14.5%). About 8% of the arguments reported that the tax would help to increase the government's revenues,

and the extra resources collected could be used to treat diseases that list sugar consumption as one of the risk factors. It was also noted that the extra resources could also provide a free, nutritious, and healthy breakfast to the school children. Articles that presented arguments opposing the SSBs tax argued the tax would increase manufacturers' operating costs due to reformulating their products, resulting in profit reduction (8.1%). Approximately 4.8% of online articles reported the tax was not an appropriate solution for diabetes control. However, only 3.2% opposed SSBs taxation due to the perception that the price for all the drinks will increase. An example of a statement regarding the price hike included "implementation of excise duty on sugar-sweetened beverages today has resulted in a price increase for most of the products between RM0.20 to RM0.70... few sundry shops have already adhered to the price adjustment of their goods...the products affected are carbonated drinks, ready-to-drink coffee, milk, juices, and canned and packet drinks" [33].

Table 2. Framing arguments presented in online news articles.

Framing argument	Articles, n (%)
Argument supporting the SSBs^a tax	29 (46.8)
Incentive for the manufacturer to introduce a healthier product voluntarily	9 (14.5)
The collection of the SSBs taxes will increase the government revenue and help to reduce the burden of future treatment costs in public healthcare facilities	5 (8.1)
Positive health consequences of reducing SSB consumption	13 (21.7)
Revenue generated will be used to help children (eg, providing free breakfasts for all school children)	2 (3.2)
Argument opposing the SSBs Tax	10 (16.1)
SSBs taxation will increase the operating costs of reformulating products to avoid taxation, resulting in decreased profit margins	5 (8.1)
SSBs taxation is not the solution to obesity	3 (4.8)
The tax might increase the overall price as the SSBs will be taxed twice (SST and SSBs Tax)	2 (3.2)
Balanced argument	23 (37.1)

^aSSBs: sugar-sweetened beverages.

Sector-Related Purposes

Online news articles mentioned a few sectors used to justify the introduction of the SSBs tax. The top priority sector used as consideration for the tax was presented from the health perspective (35.5%), followed by economic purposes and creating healthy policies (22.6%), and regulating the food and drinks industry or manufacturers (9.7%). However, few sectors, including health, were mentioned in the online news articles discussing the implementation of the tax (9.7%).

Health-Related Issues

The primary consideration for introducing the SSBs taxation was health-related issues. The main highlight for this argument was that sugar is a high-risk factor for diabetes (21%), followed by general health (42.9%) and obesity (10.7%). However, there was no mention of caries or dental problems as one of the reasons justifying SSBs taxation.

Quotes or Position Statement

Most of the quotes or position statements were from the manufacturers (17.7%). The main concern of the beverage industry was the increased cost that they had to bear once the tax is imposed. The finance ministry, which made up 16.1% of the quotes, was concerned about how the tax would be implemented, while the health ministry hoped that the tax would lead to reduced consumption of SSBs by consumers and improved health gains (14.5%). However, 1 NGO argued that the tax would not solve any health-related issues. Instead, it will have a negative impact on the lower-income populations [34] due to consumers incurring the additional cost imposed by the SSBs tax. In contrast, 7 out of 8 NGOs agreed SSBs taxation could alleviate the government's financial burden in the long term as fewer expenses may be required to treat sugar-related diseases [35].

Online news articles also quoted statements from other government stakeholders such as the Director-General of Royal Malaysia Customs Department (8.1%), the Deputy Prime Minister (6.5%), the WHO (World Health Organization; 6.5%), economic analysts from local institutions (3.2%), political parties (1.6%), and the Ministry of Domestic Trade and Consumer Affairs (1.6%). The remaining articles either quoted more than 1 organization (6.5%) or did not quote any organization (4.8%).

Discussion

Principal Findings

This study aimed to explore and generate an in-depth understanding of the media response to the SSBs taxation policy. Most of the online news articles were positive towards the implementation of SSBs taxation. They primarily addressed the benefits of SSBs taxation to the general society and the individual by advocating that reduced consumption of SSBs will improve one's overall health by reducing the risk of sugar-related diseases [36-38]. According to Vermeen et al [39], an additional 20% tax on SSBs would result in a modest reduction in BMI, translating into positive health gains adding up to approximately 170,000 healthy life years over the lifetime of the Australian adult population [39].

The Malaysian government also hoped that the SSBs manufacturers could reformulate or reinvent their products to make healthier beverages without imposing any policy. The suggestion from the finance minister stated that manufacturers need to lower the sugar content to avoid the imposed SSBs tax [37,38,40-45]. Out of the 62 articles analyzed, 29 (46.8%) articles supported the taxation. This is a positive result for public health advocates, indicating that the online news articles on the SSBs tax were generally accepted as reputable, newsworthy, and fundamental to reducing the overconsumption of sugar.

The articles also highlighted that taxation would increase government revenue and reduce the burden of future treatment costs on public healthcare facilities, which are heavily subsidized by the government [36,38]. Another finding supported the SSBs tax as it can reduce the disease burden and health care costs associated with sugar consumption [39].

The articles also emphasized using the tax revenue to benefit the school children by providing them with healthy and nutritional breakfasts. The articles aimed to garner public support, especially from the parents, as the SSBs tax would benefit their children rather than cause harm [46]. This is also in line with the Sustainable Development Goal 2 "to ensure that every child, young person, and woman received a nutritious, safe, affordable and sustainable diet that they need to reach their full potential" [47]. It is crucial for children to grow and learn and participate in their communities during their school learning days. The MyBreakfast study on breakfast consumption among Malaysian primary and secondary children highlighted that 1 out of 4 children skipped breakfast, concluding a regular breakfast is associated with healthier body weight and should be encouraged [48].

However, despite the positive slant of most online news articles, articles regarding the manufacturer and consumer response contributed to the negative slant. There is a possibility that these 2 groups were the most affected by the SSBs tax. One of the manufacturers' concerns was the increase in operating costs to reformulate their products, resulting in a price hike on SSBs that is passed to the consumers [49,50]. Another respondent stated that the SSBs tax is not a solution to diabetes and associated the ineffectiveness of SSBs taxation with the tobacco tax, which did not solve society's smoking addiction [51]. A systematic review by Escobar et al [52] suggested that an increase in the price of SSBs was associated with a decrease in consumption. The higher the price increase, the more significant the reduction in consumption. However, the argument against the imposition of the SSBs tax is that it is regressive and could negatively impact lower-income households who spent a considerable portion of their income on inexpensive, prepackaged consumable goods compared to the higher-income households [52]. However, a study by Bourke and Veerman [53] in Indonesia suggested that while an excise tax on SSBs could decrease the incidence of NCDs in all groups, the health benefits will accrue primarily in the high-income groups as they consume more sugary drinks and pay more of the tax than the lower-income group and thus the tax is not regressive [53].

Most of the articles portrayed potential health outcomes as the government's justification in implementing the SSBs tax.

Evidence supports that such taxes could substantially reduce consumption and reduce the incidence of diabetes and obesity [6,52]. It is also well-known that excessive sugar intake is a significant risk factor for caries development. Findings from the latest Malaysian oral health surveys showed that even though caries prevalence has reduced over the years, it needs to be addressed as it is still high among certain age groups, particularly children aged 6 years (71.3% in the 2015 National Oral Health Survey of Preschool Children) and adults (88.9% in the 2010 National Oral Health Survey of Adults) [54]. Jevdjevic et al [24] showed that SSBs taxation might reduce the caries-related burden and improve oral health, especially among the younger age groups [24].

Given the high media reliance on politicians and key decision-makers to portray social problems, there is a potential for widespread public acceptance or rejection of SSBs taxation [55]. Quoting the appropriate spokesperson to avoid inaccuracies in news media articles regarding SSBs taxation is crucial [55]. A study by Bødker et al [56] demonstrated that active industry lobbying and consequent judicial actions could undermine policy support from all stakeholders [56]. Another study found that SSBs manufacturers did have substantial coverage in the Malaysian media, allowing them to express their perspective on the tax. The issues they put forward mainly involved the impact of taxation on their businesses and the consumers. They attempted to portray that the taxation will increase their operation costs, which will ultimately result in higher costs for the consumers due to the increased retail price of SSBs [57,58]. This perspective may result in public outcry, which will negatively impact policy and ultimately undermine public health efforts to reduce SSBs consumption.

This study analyzed a current debate in the media that supported reducing the obesity and diabetes prevalence by introducing the SSBs tax as a sugar intake reduction policy. The findings have helped to frame evidentiary public health policy more broadly. The total number of articles written regarding SSBs taxation showed how public health studies and facts could promote media agenda setting. Public health proponents will welcome the substantial portrayal of sugar and SSBs as a social health issue, primarily driven by the food and drink industries and best

tackled by policy measures. Moreover, while the framing arguments produced by online news outlets were overwhelmingly favorable to public health advocacy and supporting coverage outnumbered opposition coverage, public health advocates should be mindful of the predominance of opposition outlets around the SSBs tax announcement. A concerted media advocacy campaign could have mitigated the wave of SSBs taxation opposition.

The findings of this study are significant because readers' comments and public views on online news coverage can also influence decisions about the final form of the SSBs tax, which may contribute to its successful implementation and overall efficacy. Thus, public health advocates should understand and study public opinion to ensure the effective implementation of the levy [16].

Study Limitations

It is also important to note that the findings from this study are subject to a range of limitations. For example, news articles were not analyzed for readership numbers or the number of times each article was shared on social media; therefore, the articles' overall reach could not be determined. Another limitation of the content analysis was that the coding was based on the subjective assessment of the 2 independent researchers.

Conclusions

The findings showed that most online news articles were written with a favorable slant towards implementing the SSBs tax. This suggests the media played an important role in supporting the health policy. Besides, the policymakers should also complement the SSBs tax with other strategies such as adequate health education and promotion as there may be a lack of awareness in the general public. Hence, the potential benefit of the SSBs tax will go beyond merely influencing price-based purchasing behavior and extend to the normalization of sugar-free beverage consumption, including plain drinking water. Furthermore, the opposing arguments towards the SSBs tax found through this study could be used by relevant bodies to anticipate opposition and assist in reframing and counteracting arguments opposed to the SSBs tax.

Acknowledgments

Special gratitude and appreciation go to the Universiti Teknologi MARA (UiTM) Shah Alam for providing the NVivo software license used to complete this study. We would also like to thank the Director-General of Health Malaysia for his permission to publish this article.

Conflicts of Interest

None declared.

References

1. Thow AM, Jan S, Leeder S, Swinburn B. The effect of fiscal policy on diet, obesity and chronic disease: a systematic review. *Bull World Health Organ* 2010 Feb 22;88(8):609-614. [doi: [10.2471/blt.09.070987](https://doi.org/10.2471/blt.09.070987)]
2. Malik VS, Popkin BM, Bray GA, Després J, Willett WC, Hu FB. Sugar-sweetened beverages and risk of metabolic syndrome and type 2 diabetes: a meta-analysis. *Diabetes Care* 2010 Nov;33(11):2477-2483 [FREE Full text] [doi: [10.2337/dc10-1079](https://doi.org/10.2337/dc10-1079)] [Medline: [20693348](https://pubmed.ncbi.nlm.nih.gov/20693348/)]

3. Malik VS, Schulze MB, Hu FB. Intake of sugar-sweetened beverages and weight gain: a systematic review. *Am J Clin Nutr* 2006 Aug;84(2):274-288 [FREE Full text] [doi: [10.1093/ajcn/84.1.274](https://doi.org/10.1093/ajcn/84.1.274)] [Medline: [16895873](https://pubmed.ncbi.nlm.nih.gov/16895873/)]
4. Malik V, Willett W, Hu F. Sugar-sweetened beverages and BMI in children and adolescents: reanalyses of a meta-analysis. *Am J Clin Nutr* 2009 Jan;89(1):438-439. [doi: [10.3945/ajcn.2008.26980](https://doi.org/10.3945/ajcn.2008.26980)] [Medline: [19056589](https://pubmed.ncbi.nlm.nih.gov/19056589/)]
5. Imamura F, O'Connor L, Ye Z, Mursu J, Hayashino Y, Bhupathiraju SN, et al. Consumption of sugar sweetened beverages, artificially sweetened beverages, and fruit juice and incidence of type 2 diabetes: systematic review, meta-analysis, and estimation of population attributable fraction. *BMJ* 2015 Jul 21;351:h3576 [FREE Full text] [doi: [10.1136/bmj.h3576](https://doi.org/10.1136/bmj.h3576)] [Medline: [26199070](https://pubmed.ncbi.nlm.nih.gov/26199070/)]
6. Taxes on sugary drinks: Why do it? World Health Organization. 2017. URL: <https://apps.who.int/iris/handle/10665/260253> [accessed 2020-04-08]
7. Lobstein T. Reducing consumption of sugar-sweetened beverages to reduce the risk of childhood overweight and obesity. World Health Organization. 2014 Sep. URL: https://www.who.int/elena/titles/bbc/ssbs_childhood_obesity/en/ [accessed 2020-04-08]
8. Han E, Kim TH, Powell LM. Beverage consumption and individual-level associations in South Korea. *BMC Public Health* 2013 Mar 06;13:195 [FREE Full text] [doi: [10.1186/1471-2458-13-195](https://doi.org/10.1186/1471-2458-13-195)] [Medline: [23497024](https://pubmed.ncbi.nlm.nih.gov/23497024/)]
9. Han E, Powell LM. Consumption patterns of sugar-sweetened beverages in the United States. *J Acad Nutr Diet* 2013 Jan;113(1):43-53 [FREE Full text] [doi: [10.1016/j.jand.2012.09.016](https://doi.org/10.1016/j.jand.2012.09.016)] [Medline: [23260723](https://pubmed.ncbi.nlm.nih.gov/23260723/)]
10. Al Otaibi HH. Sugar Sweetened Beverages Consumption Behavior and Knowledge among University Students in Saudi Arabia Risk on Internet Banking Acceptance from the User Perspective. *JOEBM* 2017;5(4):173-176. [doi: [10.18178/joebm.2017.5.4.507](https://doi.org/10.18178/joebm.2017.5.4.507)]
11. Wan L, Watson E, Arthur R. Sugar taxes: The global picture in 2017. *BeverageDaily*. 2017 Dec 20. URL: <https://www.beveragedaily.com/Article/2017/12/20/Sugar-taxes-The-global-picture-in-2017> [accessed 2019-06-18]
12. Sowa P, Keller E, Stormon N, Laloo R, Ford P. The impact of a sugar-sweetened beverages tax on oral health and costs of dental care in Australia. *Eur J Public Health* 2019 Feb 01;29(1):173-177. [doi: [10.1093/eurpub/cky087](https://doi.org/10.1093/eurpub/cky087)] [Medline: [29796599](https://pubmed.ncbi.nlm.nih.gov/29796599/)]
13. Allen W, Allen KJ. Should Australia tax sugar-sweetened beverages? *J Paediatr Child Health* 2020 Jan;56(1):8-15 [FREE Full text] [doi: [10.1111/jpc.14666](https://doi.org/10.1111/jpc.14666)] [Medline: [31782574](https://pubmed.ncbi.nlm.nih.gov/31782574/)]
14. Di Natale R, Singh L, Georgiou P, Kitching K, Paterson J, Stoker A, et al. Select Committee into the Obesity Epidemic in Australia - Final Report. Parliament of Australia. Canberra: Parliament House, Canberra; 2018 Dec 05. URL: https://www.aph.gov.au/Parliamentary_Business/Committees/Senate/Obesity_epidemic_in_Australia/Obesity/Final_Report [accessed 2020-08-06]
15. Curry LE, Rogers T, Williams P, Homs G, Willett J, Schmitt CL. Public Attitudes and Support for a Sugar-Sweetened Beverage Tax in America's Heartland. *Health Promot Pract* 2018 May;19(3):418-426. [doi: [10.1177/1524839917709759](https://doi.org/10.1177/1524839917709759)] [Medline: [28587533](https://pubmed.ncbi.nlm.nih.gov/28587533/)]
16. Thomas-Meyer M, Mytton O, Adams J. Public responses to proposals for a tax on sugar-sweetened beverages: A thematic analysis of online reader comments posted on major UK news websites. *PLoS One* 2017;12(11):e0186750 [FREE Full text] [doi: [10.1371/journal.pone.0186750](https://doi.org/10.1371/journal.pone.0186750)] [Medline: [29166399](https://pubmed.ncbi.nlm.nih.gov/29166399/)]
17. Popkin B, Ng SW. Sugar-sweetened beverage taxes: Lessons to date and the future of taxation. *PLoS Med* 2021 Jan;18(1):e1003412 [FREE Full text] [doi: [10.1371/journal.pmed.1003412](https://doi.org/10.1371/journal.pmed.1003412)] [Medline: [33411708](https://pubmed.ncbi.nlm.nih.gov/33411708/)]
18. Lee MM, Falbe J, Schillinger D, Basu S, McCulloch CE, Madsen KA. Sugar-Sweetened Beverage Consumption 3 Years After the Berkeley, California, Sugar-Sweetened Beverage Tax. *Am J Public Health* 2019 Apr;109(4):637-639. [doi: [10.2105/AJPH.2019.304971](https://doi.org/10.2105/AJPH.2019.304971)] [Medline: [30789776](https://pubmed.ncbi.nlm.nih.gov/30789776/)]
19. Bandy LK, Scarborough P, Harrington RA, Rayner M, Jebb SA. Reductions in sugar sales from soft drinks in the UK from 2015 to 2018. *BMC Med* 2020 Jan 13;18(1):20 [FREE Full text] [doi: [10.1186/s12916-019-1477-4](https://doi.org/10.1186/s12916-019-1477-4)] [Medline: [31931800](https://pubmed.ncbi.nlm.nih.gov/31931800/)]
20. Belanjawan 2019. Kementerian Kewangan Malaysia. Kuala Lumpur: Pencetakan Nasional Malaysia Berhad; 2019 Nov 02. URL: <https://www.mof.gov.my/arkib/belanjawan/2019/ub19.pdf> [accessed 2020-08-06]
21. Watts C, Freeman B. "Where There's Smoke, There's Fire": A Content Analysis of Print and Web-Based News Media Reporting of the Philip Morris-Funded Foundation for a Smoke-Free World. *JMIR Public Health Surveill* 2019 Jun 06;5(2):e14067 [FREE Full text] [doi: [10.2196/14067](https://doi.org/10.2196/14067)] [Medline: [31172959](https://pubmed.ncbi.nlm.nih.gov/31172959/)]
22. Buckton CH, Patterson C, Hyseni L, Katikireddi SV, Lloyd-Williams F, Elliott-Green A, et al. The palatability of sugar-sweetened beverage taxation: A content analysis of newspaper coverage of the UK sugar debate. *PLoS One* 2018;13(12):e0207576 [FREE Full text] [doi: [10.1371/journal.pone.0207576](https://doi.org/10.1371/journal.pone.0207576)] [Medline: [30517133](https://pubmed.ncbi.nlm.nih.gov/30517133/)]
23. Otten AL. The influence of the mass media on health policy. *Health Aff (Millwood)* 1992;11(4):111-118. [doi: [10.1377/hlthaff.11.4.111](https://doi.org/10.1377/hlthaff.11.4.111)] [Medline: [1483630](https://pubmed.ncbi.nlm.nih.gov/1483630/)]
24. Jevdjevic M, Trescher A, Rovers M, Listl S. The caries-related cost and effects of a tax on sugar-sweetened beverages. *Public Health* 2019 Apr;169:125-132. [doi: [10.1016/j.puhe.2019.02.010](https://doi.org/10.1016/j.puhe.2019.02.010)] [Medline: [30884363](https://pubmed.ncbi.nlm.nih.gov/30884363/)]
25. Jou J, Niederdeppe J, Barry CL, Gollust SE. Strategic messaging to promote taxation of sugar-sweetened beverages: lessons from recent political campaigns. *Am J Public Health* 2014 May;104(5):847-853. [doi: [10.2105/AJPH.2013.301679](https://doi.org/10.2105/AJPH.2013.301679)] [Medline: [24625177](https://pubmed.ncbi.nlm.nih.gov/24625177/)]

26. Niederdeppe J, Gollust SE, Jarlenski MP, Nathanson AM, Barry CL. News coverage of sugar-sweetened beverage taxes: pro- and antitax arguments in public discourse. *Am J Public Health* 2013 Jun;103(6):e92-e98. [doi: [10.2105/AJPH.2012.301023](https://doi.org/10.2105/AJPH.2012.301023)] [Medline: [23597354](https://pubmed.ncbi.nlm.nih.gov/23597354/)]
27. Essman M, Stoltze FM, Carpentier FD, Swart EC, Taillie LS. Examining the news media reaction to a national sugary beverage tax in South Africa: a quantitative content analysis. *BMC Public Health* 2021 Mar 06;21(1):454 [FREE Full text] [doi: [10.1186/s12889-021-10460-1](https://doi.org/10.1186/s12889-021-10460-1)] [Medline: [33676468](https://pubmed.ncbi.nlm.nih.gov/33676468/)]
28. Freeman B. Tobacco plain packaging legislation: a content analysis of commentary posted on Australian online news. *Tob Control* 2011 Sep;20(5):361-366. [doi: [10.1136/tc.2011.042986](https://doi.org/10.1136/tc.2011.042986)] [Medline: [21527406](https://pubmed.ncbi.nlm.nih.gov/21527406/)]
29. Bedi RS. Malaysia jumps up 22 places in latest Press Freedom index. *The Star*. 2019 Apr 18. URL: <https://www.thestar.com.my/news/nation/2019/04/18/malaysia-jumps-up-22-places-in-latest-press-freedom-index> [accessed 2020-06-15]
30. Malaysia Central. 2016. URL: <http://www.mycen.com.my/malaysia/news.html> [accessed 2020-01-20]
31. Malaysia Complete News Online Listing. Malaysia Service Centre. 2017. URL: <http://www.malaysiaservicecentre.com/news-online.html> [accessed 2020-01-20]
32. July 2019 - MDA Release Rankings of Top Web Entities in Malaysia. Malaysian Digital Association. 2019 May 27. URL: <https://www.malaysiandigitalassociation.org.my/wp-content/uploads/2019/05/MDA-Releases-Rankings-of-Top-Web-Entities-in-Malaysia-for-Jan-2019.pdf> [accessed 2020-12-20]
33. Khan KNA. Minuman bergula naik harga [METROTV]. *Harian Metro*. 2019 Jul 01. URL: <https://www.hmetro.com.my/mutakhir/2019/07/471085/minuman-bergula-naik-harga-metrotv> [accessed 2020-05-18]
34. Halid S. Cukai soda tak selesaikan masalah kesihatan. *Harian Metro*. 2018 Aug 30. URL: <https://www.hmetro.com.my/mutakhir/2018/08/372711/cukai-soda-tak-selesaikan-masalah-kesihatan> [accessed 2020-05-20]
35. Pelaksanaan cukai ke atas minuman bergula dapat kurangkan kos perubatan masa depan. *BH Online*. 2019 Jul 01. URL: <https://www.bharian.com.my/berita/nasional/2019/07/579990/pelaksanaan-cukai-ke-atas-minuman-bergula-dapat-kurangkan-kos> [accessed 2020-05-20]
36. Bernama. Bantu kurangkan kos perubatan pada masa depan. *Sinar Harian*. 2019 Jul 01. URL: <https://www.sinarharian.com.my/article/35329/BERITA/Nasional/Bantu-kurangkan-kos-perubatan-pada-masa-depan> [accessed 2020-05-20]
37. Bernama. Cukai minuman bergula mampu kurangkan obesiti, penyakit berkaitan - Kementerian Kewangan. *Astro Awani*. 2019 Aug 02. URL: <https://www.astroawani.com/gaya-hidup/cukai-minuman-bergula-mampu-kurangkan-obesiti-penyakit-berkaitan-kementerian-kewangan-213994> [accessed 2020-05-25]
38. Bernama. Sugar tax can reduce diabetes and related conditions, says MoF. *New Straits Times*. 2019 Aug 02. URL: <https://tinyurl.com/ke28zw4y> [accessed 2020-05-23]
39. Veerman JL, Sacks G, Antonopoulos N, Martin J. The Impact of a Tax on Sugar-Sweetened Beverages on Health and Health Care Costs: A Modelling Study. *PLoS One* 2016;11(4):e0151460 [FREE Full text] [doi: [10.1371/journal.pone.0151460](https://doi.org/10.1371/journal.pone.0151460)] [Medline: [27073855](https://pubmed.ncbi.nlm.nih.gov/27073855/)]
40. Bernama. Belanjawan 2019: Langkah percukaian tepat demi masa depan lebih mantap. *Astro Awani*. 2018 Nov 02. URL: <https://www.astroawani.com/berita-malaysia/belanjawan-2019-langkah-percukaian-tepat-demi-masa-depan-lebih-mantap-190034> [accessed 2020-05-23]
41. Bernama. Cukai ke atas produk minuman bergula. *Sinar Harian*. 2019 Aug 06. URL: <https://www.sinarharian.com.my/article/41846/INFOGRAFIK/Cukai-ke-atas-produk-minuman-bergula> [accessed 2020-05-23]
42. Customs: Two-month transition period as sugar tax comes into effect. *The Star*. 2019 Jun 30. URL: <https://www.thestar.com.my/news/nation/2019/06/30/customs-two-month-transition-period-as-sugar-tax-comes-into-effect> [accessed 2020-05-23]
43. Karim LAA. Duti eksais gula mulai April 2019. *Berita Harian*. 2018 Nov 02. URL: <https://www.bharian.com.my/berita/nasional/2018/11/493722/duti-eksais-gula-mulai-april-2019> [accessed 2020-05-23]
44. Rohman MAP. Kurangkan gula boleh elak cukai. *Berita Harian*. 2019 Jul 11. URL: <https://www.bharian.com.my/berita/nasional/2019/07/583726/kurangkan-gula-boleh-elak-cukai-guan-eng> [accessed 2020-05-23]
45. Carvalho M, Sivanandam H, Rahim R, Tan T. Lower sugar content to avoid paying tax, SMEs told. *The Star*. 2019 Jul 12. URL: <https://www.thestar.com.my/news/nation/2019/07/12/lower-sugar-content-to-avoid-paying-tax-smes-told/> [accessed 2020-05-23]
46. Bernama. Hasil cukai gula biaya program sarapan percuma. *Harian Metro*. 2019 Mar 19. URL: <https://www.hmetro.com.my/bisnes/2019/03/435454/hasil-cukai-gula-biaya-program-sarapan-percuma> [accessed 2020-05-24]
47. Key asks for 2020 SDG Voluntary National Reviews: SDG 2. UNICEF. 2015. URL: <https://www.unicef.org/documents/sdg-issue-brief-2> [accessed 2020-08-11]
48. Tee E, Nurliyana A, Norimah A, Mohamed HBJ, Tan SY, Appukutty M, et al. Breakfast consumption among Malaysian primary and secondary school children and relationship with body weight status - Findings from the MyBreakfast Study. *Asia Pac J Clin Nutr* 2018;27(2):421-432 [FREE Full text] [doi: [10.6133/apjcn.062017.12](https://doi.org/10.6133/apjcn.062017.12)] [Medline: [29384332](https://pubmed.ncbi.nlm.nih.gov/29384332/)]
49. Bernama. Cukai gula perkembangan negatif untuk F and N. *Berita Harian*. 2019 Nov 09. URL: <https://www.bharian.com.my/bisnes/korporat/2018/11/496249/cukai-gula-perkembangan-negatif-untuk-fn> [accessed 2020-06-21]

50. Ishak SR. Cukai soda beri impak jangka pendek perniagaan CSR. Berita Harian. 2018 Nov 08. URL: <https://www.bharian.com.my/bisnes/korporat/2018/11/495799/cukai-soda-beri-impak-jangka-pendek-perniagaan-csr> [accessed 2020-07-14]
51. Halid S. Cukai soda tak selesaikan masalah kesihatan. Harian Metro. 2018 Aug 30. URL: <https://www.hmetro.com.my/mutakhir/2018/08/372711/cukai-soda-tak-selesaikan-masalah-kesihatan> [accessed 2020-07-15]
52. Cabrera Escobar MA, Veerman JL, Tollman SM, Bertram MY, Hofman KJ. Evidence that a tax on sugar sweetened beverages reduces the obesity rate: a meta-analysis. BMC Public Health 2013 Nov 13;13:1072 [FREE Full text] [doi: [10.1186/1471-2458-13-1072](https://doi.org/10.1186/1471-2458-13-1072)] [Medline: [24225016](https://pubmed.ncbi.nlm.nih.gov/24225016/)]
53. Bourke EJ, Veerman JL. The potential impact of taxing sugar drinks on health inequality in Indonesia. BMJ Glob Health 2018;3(6):e000923 [FREE Full text] [doi: [10.1136/bmjgh-2018-000923](https://doi.org/10.1136/bmjgh-2018-000923)] [Medline: [30555724](https://pubmed.ncbi.nlm.nih.gov/30555724/)]
54. The Oral Health Status of Malaysians. Oral Health Division, Ministry of Health Malaysia. URL: <http://ohd.moh.gov.my/index.php/en/allcategories-en-gb/2-uncategorised/124-oral-healthcare-today> [accessed 2020-07-04]
55. Einwiller SA, Carroll CE, Korn K. Under What Conditions Do the News Media Influence Corporate Reputation? The Roles of Media Dependency and Need for Orientation. Corp Reputation Rev 2010 Jan 20;12(4):299-315 [FREE Full text] [doi: [10.1057/crr.2009.28](https://doi.org/10.1057/crr.2009.28)]
56. Bødker M, Pisinger C, Toft U, Jørgensen T. The rise and fall of the world's first fat tax. Health Policy 2015 Jun;119(6):737-742. [doi: [10.1016/j.healthpol.2015.03.003](https://doi.org/10.1016/j.healthpol.2015.03.003)] [Medline: [25840733](https://pubmed.ncbi.nlm.nih.gov/25840733/)]
57. Excise duty on sweetened beverages negative for Fraser and Neave. The Star. 2018 Nov 09. URL: <https://www.thestar.com.my/business/business-news/2018/11/09/excise-duty-on-sweetened-beverages-negative-for-fraser-and-neave/> [accessed 2020-08-11]
58. Breaking down the chain: a guide to the soft drink industry. National Policy & Legal Analysis Network to Prevent Childhood Obesity. Newark, NJ: ChangeLabSolution; 2012. URL: https://www.changelabsolutions.org/sites/default/files/ChangeLab-Beverage_Industry_Report-FINAL_201109.pdf [accessed 2020-08-06]

Abbreviations

SSBs: sugar-sweetened beverages
NCD: noncommunicable disease
NGO: nongovernmental organization

Edited by T Sanchez; submitted 23.09.20; peer-reviewed by DTI Bt. Rosli, M Amini; comments to author 16.12.20; revised version received 01.03.21; accepted 07.06.21; published 18.08.21.

Please cite as:

Mohd Hanim MF, Md Sabri BA, Yusof N

Online News Coverage of the Sugar-Sweetened Beverages Tax in Malaysia: Content Analysis

JMIR Public Health Surveill 2021;7(8):e24523

URL: <https://publichealth.jmir.org/2021/8/e24523>

doi: [10.2196/24523](https://doi.org/10.2196/24523)

PMID: [34406125](https://pubmed.ncbi.nlm.nih.gov/34406125/)

©Muhammad Faiz Mohd Hanim, Budi Aslinie Md Sabri, Norashikin Yusof. Originally published in JMIR Public Health and Surveillance (<https://publichealth.jmir.org>), 18.08.2021. This is an open-access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR Public Health and Surveillance, is properly cited. The complete bibliographic information, a link to the original publication on <https://publichealth.jmir.org>, as well as this copyright and license information must be included.

Corrigenda and Addenda

Correction: Census Tract Patterns and Contextual Social Determinants of Health Associated With COVID-19 in a Hispanic Population From South Texas: A Spatiotemporal Perspective

Cici Bauer¹, PhD; Kehe Zhang¹, MS; Miryoung Lee², PhD; Susan Fisher-Hoch², MD; Esmeralda Guajardo³, MA; Joseph McCormick², MD; Isela de la Cerda², MS; Maria E Fernandez⁴, PhD; Belinda Reininger⁵, DrPH

¹Department of Biostatistics and Data Science, School of Public Health, The University of Texas Health Science Center at Houston, Houston, TX, United States

²Department of Epidemiology, Human Genetics and Environmental Science, School of Public Health, The University of Texas Health Science Center at Houston, Brownsville, TX, United States

³Cameron County Public Health, San Benito, TX, United States

⁴Department of Health Promotion and Behavior Sciences, School of Public Health, The University of Texas Health Science Center at Houston, Houston, TX, United States

⁵Department of Health Promotion and Behavior Sciences, School of Public Health, The University of Texas Health Science Center at Houston, Brownsville, TX, United States

Corresponding Author:

Cici Bauer, PhD

Department of Biostatistics and Data Science

School of Public Health

The University of Texas Health Science Center at Houston

1200 Pressler Street

Houston, TX, 77030

United States

Phone: 1 713 500 9581

Email: cici.x.bauer@uth.tmc.edu

Related Article:

Correction of: <https://publichealth.jmir.org/2021/8/e29205>

(*JMIR Public Health Surveill* 2021;7(8):e32870) doi:[10.2196/32870](https://doi.org/10.2196/32870)

In “Census Tract Patterns and Contextual Social Determinants of Health Associated With COVID-19 in a Hispanic Population From South Texas: A Spatiotemporal Perspective” (*JMIR Public Health Surveill* 2021;7(8):e29205), one error was noted.

Due to a system error, the name of one author, Joseph McCormick, was replaced with the name of another author on the paper, Isela de la Cerda. In the originally published paper, the order of authors was listed as follows:

Cici Bauer, Kehe Zhang, Miryoung Lee, Susan Fisher-Hoch, Esmeralda Guajardo, Isela de la Cerda, Isela de la Cerda, Maria E Fernandez, Belinda Reininger

This has been corrected to:

Cici Bauer, Kehe Zhang, Miryoung Lee, Susan Fisher-Hoch, Esmeralda Guajardo, Joseph

McCormick, Isela de la Cerda, Maria E Fernandez, Belinda Reininger

In the originally published paper, the ORCID of author Isela de la Cerda was incorrectly published as follows:

0000-0002-5844-8102

This has been corrected to:

0000-0003-3625-8954

The correction will appear in the online version of the paper on the JMIR Publications website on August 18, 2021, together with the publication of this correction notice. Because this was made after submission to PubMed, PubMed Central, and other full-text repositories, the corrected article has also been resubmitted to those repositories.

Submitted 12.08.21; this is a non-peer-reviewed article; accepted 12.08.21; published 18.08.21.

Please cite as:

Bauer C, Zhang K, Lee M, Fisher-Hoch S, Guajardo E, McCormick J, de la Cerda I, Fernandez ME, Reininger B
Correction: Census Tract Patterns and Contextual Social Determinants of Health Associated With COVID-19 in a Hispanic Population
From South Texas: A Spatiotemporal Perspective
JMIR Public Health Surveill 2021;7(8):e32870
URL: <https://publichealth.jmir.org/2021/8/e32870>
doi: [10.2196/32870](https://doi.org/10.2196/32870)
PMID: [34406965](https://pubmed.ncbi.nlm.nih.gov/34406965/)

©Cici Bauer, Kehe Zhang, Miryoung Lee, Susan Fisher-Hoch, Esmeralda Guajardo, Joseph McCormick, Isela de la Cerda, Maria E Fernandez, Belinda Reininger. Originally published in JMIR Public Health and Surveillance (<https://publichealth.jmir.org>), 18.08.2021. This is an open-access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR Public Health and Surveillance, is properly cited. The complete bibliographic information, a link to the original publication on <https://publichealth.jmir.org>, as well as this copyright and license information must be included.

Original Paper

Forecasting COVID-19 Hospital Census: A Multivariate Time-Series Model Based on Local Infection Incidence

Hieu M Nguyen^{1*}, MSc; Philip J Turk^{1*}, MSc, PhD; Andrew D McWilliams¹, MPH, MD

Center for Outcomes Research and Evaluation, Atrium Health, Charlotte, NC, United States

*these authors contributed equally

Corresponding Author:

Hieu M Nguyen, MSc

Center for Outcomes Research and Evaluation

Atrium Health

1300 Scott Ave

Charlotte, NC, 28204

United States

Phone: 1 9706914892

Email: hieu.nguyen@atriumhealth.org

Abstract

Background: COVID-19 has been one of the most serious global health crises in world history. During the pandemic, health care systems require accurate forecasts for key resources to guide preparation for patient surges. Forecasting the COVID-19 hospital census is among the most important planning decisions to ensure adequate staffing, number of beds, intensive care units, and vital equipment.

Objective: The goal of this study was to explore the potential utility of local COVID-19 infection incidence data in developing a forecasting model for the COVID-19 hospital census.

Methods: The study data comprised aggregated daily COVID-19 hospital census data across 11 Atrium Health hospitals plus a virtual hospital in the greater Charlotte metropolitan area of North Carolina, as well as the total daily infection incidence across the same region during the May 15 to December 5, 2020, period. Cross-correlations between hospital census and local infection incidence lagging up to 21 days were computed. A multivariate time-series framework, called the vector error correction model (VECM), was used to simultaneously incorporate both time series and account for their possible long-run relationship. Hypothesis tests and model diagnostics were performed to test for the long-run relationship and examine model goodness of fit. The 7-days-ahead forecast performance was measured by mean absolute percentage error (MAPE), with time-series cross-validation. The forecast performance was also compared with an autoregressive integrated moving average (ARIMA) model in the same cross-validation time frame. Based on different scenarios of the pandemic, the fitted model was leveraged to produce 60-days-ahead forecasts.

Results: The cross-correlations were uniformly high, falling between 0.7 and 0.8. There was sufficient evidence that the two time series have a stable long-run relationship at the .01 significance level. The model had very good fit to the data. The out-of-sample MAPE had a median of 5.9% and a 95th percentile of 13.4%. In comparison, the MAPE of the ARIMA had a median of 6.6% and a 95th percentile of 14.3%. Scenario-based 60-days-ahead forecasts exhibited concave trajectories with peaks lagging 2 to 3 weeks later than the peak infection incidence. In the worst-case scenario, the COVID-19 hospital census can reach a peak over 3 times greater than the peak observed during the second wave.

Conclusions: When used in the VECM framework, the local COVID-19 infection incidence can be an effective leading indicator to predict the COVID-19 hospital census. The VECM model had a very good 7-days-ahead forecast performance and outperformed the traditional ARIMA model. Leveraging the relationship between the two time series, the model can produce realistic 60-days-ahead scenario-based projections, which can inform health care systems about the peak timing and volume of the hospital census for long-term planning purposes.

(*JMIR Public Health Surveill* 2021;7(8):e28195) doi:[10.2196/28195](https://doi.org/10.2196/28195)

KEYWORDS

COVID-19; forecasting; time-series model; vector error correction model; hospital census; hospital resource utilization; infection incidence

Introduction

SARS-CoV-2 is a novel member of the coronavirus family, and infections in humans can result in the disease COVID-19. The virus is transmitted primarily through droplets from coughing and sneezing and is highly infectious. Its basic reproduction rate is estimated to be in the low to mid 2s based on different models [1], compared to 2 for severe acute respiratory syndrome (SARS) and 1.3 for the 2009 swine flu [2]. Moderate to severe disease typically manifests with acute hypoxemia, and can progress to acute respiratory distress syndrome, multiorgan dysfunction, and death. Furthermore, an estimated 25%-30% of patients admitted to hospitals require intensive care admission [2]. In December 2019, the first cases were recorded in Wuhan, China, with subsequent spread across the world. In early 2020, the World Health Organization declared COVID-19 to be a global health emergency [3]. At the end of December 2020, SARS-CoV-2 had resulted in over 82 million documented cases and nearly 2 million deaths [4].

Our work is motivated by the need of hospital leaders to have timely and accurate forecasts to guide planning for surges in hospital demands due to the pandemic. Adequate preparation can help prevent or mitigate strains on hospital resources that result when hospitals exceed their historical capacity. On the contrary, being caught off-guard under a pandemic can devastate the population and health care systems. For example, previous models in India suggested falsely that it had reached herd immunity, encouraging complacency and insufficient preparation; however, on May 4, 2021, there was still a reported rolling average of 378,000 cases a day, which overwhelmed hospitals and health workers and resulted in a national health crisis [5]. Thus, to a health care system, an essential tool is a model that provides short- and long-range forecasting of the number of COVID-19-positive patients who will be admitted. This COVID-19 hospital census plays a central role in planning decisions that frequently require considerable lead time, such as increasing staff, creating physical beds and rooms, and procuring vital equipment (eg, ventilators and personal protective equipment).

Prior research has demonstrated the utility of forecasting hospital demands (eg, hospital admissions, intensive care unit census, and hospital overall census) using univariate time-series models such as the autoregressive integrated moving average (ARIMA), the seasonal autoregressive integrated moving average (SARIMA), and exponential smoothing [6-8]. Another approach is to use ensemble-based modeling. For example, a hybrid of a SARIMA model and a nonlinear autoregression artificial neural network model has been used to forecast hospital admissions [9]. In another example, two separate models, a time-series model for hospital admission and a patient-level logistic regression model for hospital discharge, were combined to predict the hospital census [10]. While these examples demonstrate the powerful potential of univariate time-series and ensemble modeling, neither incorporate factors inherent to the behavior of the pandemic, which may serve as important leading indicators of hospital census, especially at times when infection rates become increasingly dynamic (eg, on the approach or descent of peak infection prevalence). To

incorporate pandemic indicators into modeling requires recognition that such indicators are typically nonstationary. Consequently, while a stationary multivariate time-series model, called vector autoregression (VAR), has been successfully employed to forecast emergency department patient census by including other hospital resource indicators [11], it cannot be used in this situation. Rather, our problem will require nonstationary multivariate time-series models like the vector error correction model (VECM).

Recently, VECM has been used to forecast the demand for intensive care units during the COVID-19 pandemic by including hospital admission as a leading indicator [12]. Although hospital admission is a natural choice as a leading indicator, it has a short period of lead time (ie, hours to days) and thus, limited predictive power. A more powerful indicator for planning purposes would lead by days to weeks. We have previously used VECM to forecast COVID-19 hospital census using leading indicators from Google relative search volumes for COVID-19 testing-related terms combined with the number of people flagged as having possible COVID-19 when using an internet-based virtual health screening bot [13]. However, these COVID-19 indicators, which are based on symptoms, have limitations. For example, the symptoms of COVID-19 cannot be easily separated from other common conditions, such as the seasonal flu, and search patterns may change due to other external factors over time.

During the COVID-19 pandemic, many papers have been devoted to developing predictive models for the volume of new cases (ie, infection incidence) using various methods from time-series analyses [14-16] to advanced machine learning [17,18]. However, virtually no effort was focused on developing statistical models linking infection incidence to hospitalization. Because hospital admission typically follows the symptoms or exposure that may provoke a person to be tested by roughly 1 week, we hypothesize that at a local population level, infection incidence rates may have a stable relationship with and serve as a reliable leading indicator for the COVID-19 hospital census. In this paper, our main objective is to explore whether the local COVID-19 infection incidence and the COVID-19 hospital census can be successfully incorporated within a VECM to delivery satisfactory 7-days-ahead forecast performance and examine the application of this model to scenario-based long-term forecasting. From our experience, since there can be systematic changes due to the day of the week in a hospital time series, we will need to account for weekly seasonal effects and examine implications on short-term resource planning.

Methods

Time-Series Data

Atrium Health is a large, integrated health care system operating in North Carolina, South Carolina, and Georgia. In this paper, the COVID-19 hospital census (census) refers to the daily aggregate number of beds occupied by patients with COVID-19 at midnight across the subset of 11 Atrium Health hospitals in the greater Charlotte metropolitan area of North Carolina, plus a virtual hospital (Atrium Health Hospital at Home). The virtual hospital uses telemedicine to treat patients who require only a

minimal level of care. The local COVID-19 infection incidence (incidence) is the aggregate daily count of new COVID-19–positive cases from 11 local counties belonging to the Cities Readiness Initiative (CRI) region, as designated by the North Carolina Department of Health and Human Services. The CRI region roughly approximates the market catchment area of these hospitals.

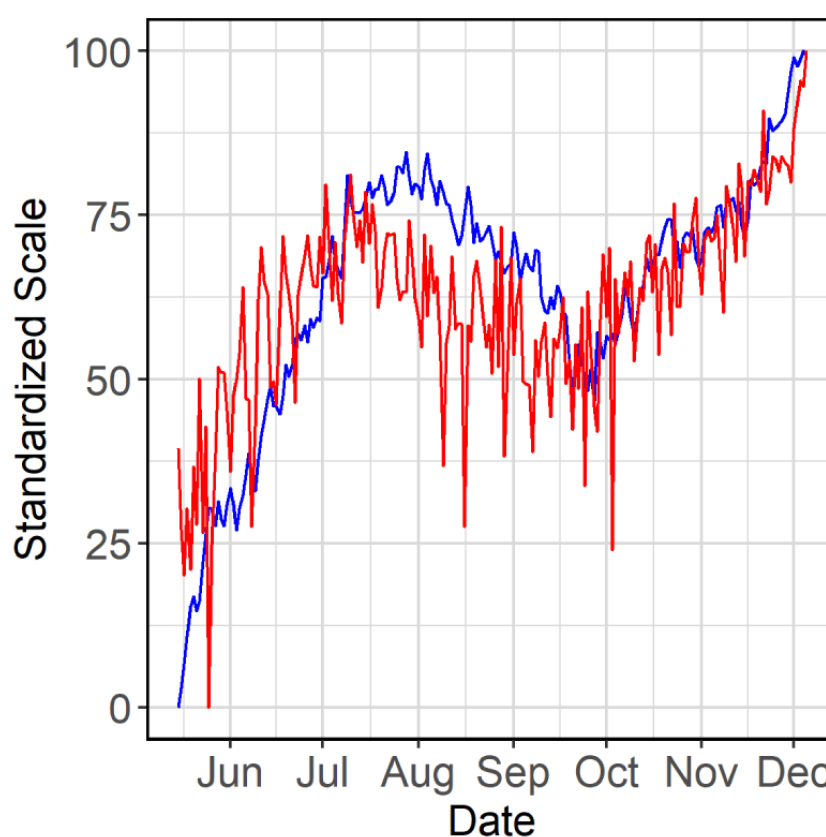
Using STL (seasonal and trend decomposition using Loess) time-series decomposition [19], we observed that the two time series had multiplicative weekly seasonality. We transformed both time series to achieve additive seasonality and linearize their relationship. The usual log transformation was applied to incidence. For operational purposes, the health system had

previously decided to place an upper bound of 1000 patients with COVID-19 on the hospital time-series range, so we applied the following constrained log transformation so that the back-transformed census forecasts would satisfy the constraint:

$$\log(x)$$

The forecast model described in the following sections was developed for these transformed time series. Figure 1 shows a plot of transformed census and incidence on a standardized scale for the period from May 15 to December 5, 2020. To affirm the association between the two transformed time series, we computed the Pearson cross-correlations between census and values of incidence at lags 0, -1, ..., -21.

Figure 1. Scaled time series for COVID-19 hospital census and local COVID-19 infection incidence in the Cities Readiness Initiative region for the period from May 15 to December 5, 2020. Transformed census (blue) and incidence (red) are linearly standardized to the 0-100 scale.



VECM

A VECM is a vector autoregressive model used for nonstationary multivariate time series and accounts for stable long-run relationships, that is, cointegration, between the time series. A $k \times 1$ time-series vector y_t is said to be cointegrated if there is at least one nonzero $k \times 1$ vector β_r , such that the linear combination $\beta_r' y_t$ is trend-stationary. If r such linearly independent vectors β_i ($i=1, \dots, r$) exist, we say y_t is cointegrated with cointegration rank r [20].

Following Pfaff [20], we first describe the VAR representation of order p of the VECM:

$$y_t = \mu + \Phi_1 y_{t-1} + \dots + \Phi_p y_{t-p} + \epsilon_t$$

for time $t=1, \dots, T$, where Π_i (for $i=1, \dots, p$) are $k \times k$ coefficient matrices of the lagged series at lag i , μ is a $k \times 1$ vector of constants, D_t is a 6×1 vector of weekly seasonal indicators, Φ is a $k \times 6$ coefficient matrix for seasonal indicators, and ϵ_t is a $k \times 1$ vector of random errors.

The VECM specification can be formulated as an algebraic rearrangement of the VAR representation as:

$$\Delta y_t = \mu + \Phi_1 \Delta y_{t-1} + \dots + \Phi_p \Delta y_{t-p} + \epsilon_t$$

where Δy_t is a $k \times 1$ vector of the differenced series Δy_t and Δy_{t-1} .

The model has the following assumptions:

- Assumption 1: The components of \mathbf{y}_t are at most $I(1)$, that is, an integrated of order 1
- Assumption 2: $0 \leq r = \text{rank}(\Pi) \leq k$
- Assumption 3: ε_t are identically and independently distributed $N(\mathbf{0}, \Sigma)$ random vectors with covariance matrix Σ .

We now discuss the implications of the assumptions. For assumption 2, if $r=k$, then it can be shown that the VECM becomes a standard VAR model. If $r=0$, then Π is the zero matrix and there is no cointegration relationship between the series. The VECM then becomes a VAR model for differenced time series. If $0 < r < k$, then Π can be factored into $\Pi = \alpha\beta^T$, where α and β are both $k \times r$ matrices. From assumption 1, the differenced series $\Delta\mathbf{y}_t$, and its lags $\Delta\mathbf{y}_{t-1}, \dots, \Delta\mathbf{y}_{t-p+1}$ are stationary. It follows that $\Pi\mathbf{y}_{t-1} = \alpha\beta^T\mathbf{y}_{t-1}$, as well as $\beta^T\mathbf{y}_{t-1}$, also called the error correction term, is (trend-)stationary, depending on the specification of the deterministic components. The r linearly independent columns of β are the cointegrating vectors, and the rank r is equal to the cointegration rank of the system of time series.

Estimation and Inference

The VECM was specified and fitted with the steps below.

First, to choose the order p of the VAR representation, we fitted a VAR model to the data and made the decision based on the Akaike information criterion (AIC) [21].

Second, we determined the number of cointegration relationships ($r=0$ or $r=1$) using the Johansen trace test [22].

Third, we needed to decide where to place the constant μ in the model. One option was to leave μ as shown previously to account for linear trend in the data. Another option was to restrict $\mu = \alpha\rho$. The constant would be absorbed into the cointegration relationship as an intercept, and the data would not exhibit linear trend.

We made our decision about whether to restrict μ based on a likelihood ratio test for linear trend, as described elsewhere [23,24].

Fourth, we used maximum likelihood estimation to fit the model, reported parameter estimates, the corresponding T tests, and the omnibus F tests with a significance level of .05, following Johansen [23].

Finally, we computed the 7-days-ahead forecasts and the 80% forecast intervals. Once the forecasts of the transformed census were made with the VECM, they were back-transformed to the original scale of census. We created 80% forecast intervals for the transformed census using a bootstrap procedure [25]. Then, the lower and upper bound of the forecast intervals were also back-transformed.

The model was fitted to the data between May 15 and December 5, 2020. All the data analysis was done using R statistical software, version 4.0.3 (R Core Team). The implementation of the VECM was done with the *tsDyn*, *vars*, and *urca* R packages. Since there were no packages to make bootstrapped forecast intervals for the VECM, we coded our own implementation.

The data and code used in the data analysis are publicly available on GitHub [26].

Model Diagnostics

We examined the omnibus F tests to look for signs of lack of fit and also performed the multivariate Portmanteau test for the existence of serial correlation in the errors. Autocorrelation function and cross-correlation function plots were also generated for visual inspection. We performed the univariate and multivariate Jarque-Bera normality test on the errors [27] and also checked whether the cointegration relationship was stable, that is, stationary, using the Augmented Dickey-Fuller (ADF) test [28] and the Kwiatkowski-Phillips-Schmidt-Shin (KPSS) test [29]. Finally, we checked the stability of the estimated VAR representation. To do so, we looked at the companion matrix of the VAR representation and checked whether the maximum eigenvalue modulus was strictly smaller than 1, which, if true, would imply the stability of the VAR representation [30]. We also generated a trace plot of the maximum eigenvalue modulus, where the model was repeatedly fitted on a daily rolling basis, to check for the consistency of this value over time.

Forecast Performance

We used mean absolute percentage error (MAPE) to evaluate the 7-days-ahead forecasts of census:

$$\text{MAPE} = \frac{1}{n} \sum_{t=1}^n \left| \frac{F_t - A_t}{A_t} \right|$$

where F_t is the forecast value and A_t is the actual value.

In order to approximate the sampling distribution of MAPE, we performed time-series cross-validation. From June 16 to November 28, 2020, for each day, we iteratively fitted the model, made 7-days-ahead forecasts, and computed the MAPE. Eventually, we obtained 166 values of MAPE, plotted the distribution, and computed the median as well as the 95th percentile. We will consider a median MAPE below 10% to be satisfactory, based on the practical effect of a peak surge on bed capacity at our health care system.

Scenario-Based Long-Term Forecasting

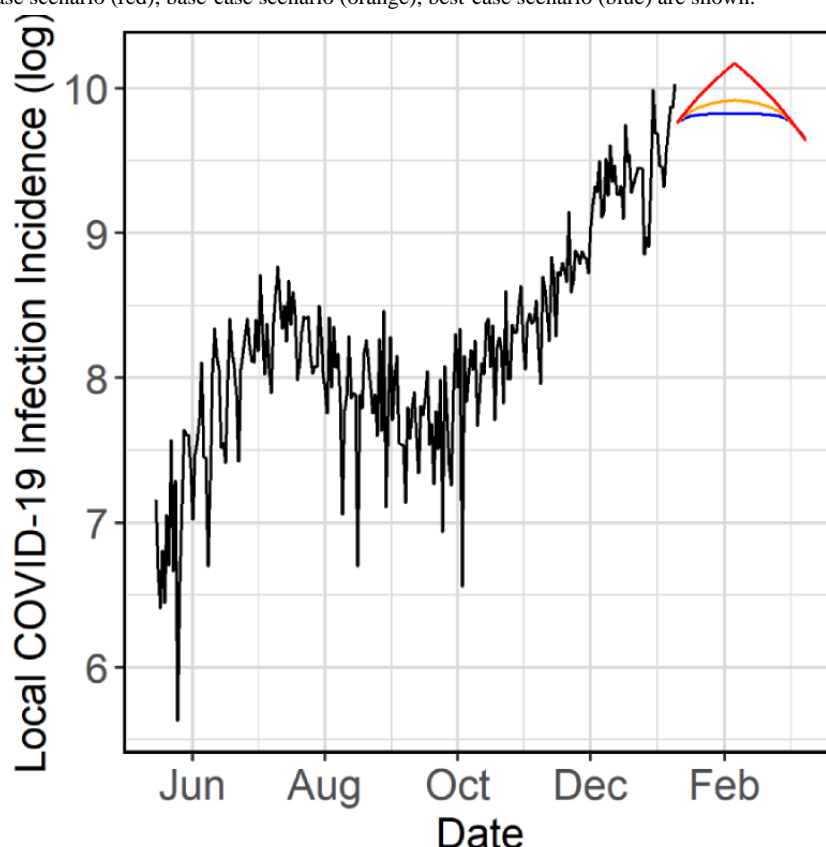
Leading up to and at the peak of infection prevalence, there can be high anxiety and uncertainty about how much more incidence and, in particular, census may increase. Furthermore, traditional univariate time-series models may give linear forecasts for census that do not accurately represent pandemic behavior. However, cointegration allows for census forecasts that leverage subtle, but critical, changes in incidence (eg, concavity). This suggests, if not necessitates, the forecasting of census under different pandemic scenarios. For resource planning, hospital leaders will want to understand the implications associated with a worst-case scenario.

For our health care system, besides routine 7-days-ahead census forecasts, we also deployed our model for 60-days-ahead census forecasts, considering 3 different scenarios of what could happen with incidence (ie, best case, base case, and worst case). On January 9, 2021, we expected the winter surge to reach peak infection prevalence around February 5, 2021, based on an extension of an epidemiological model called the susceptible-infected-removed model [31]. While peak infection

incidence typically leads peak infection prevalence, in the absence of definitively knowing either peak date, we took a conservative approach and linearly extrapolated incidence with a positive trend up to the expected pandemic peak. The severity of a scenario was controlled by a trend-dampening parameter [32]. After the peak, the descent path was initially symmetric to its ascent and then eventually became linear (Figure 2).

Using our model refitted on January 9, 2021, with an increased capacity of 1250 patients, we generated forecasts iteratively forward for 60 days using the past census forecasts together with projected incidence under each scenario. To account for uncertainty in future census and incidence, we also simulated 1000 conditional sample paths of the two time series under each scenario using the bootstrap procedure mentioned earlier and computed the 10th and 90th percentile at each horizon to obtain the 80% forecast intervals.

Figure 2. The 60-day projected local COVID-19 infection incidence in the Cities Readiness Initiative region on the log scale, as of January 9, 2021. Past values (black), worst-case scenario (red), base-case scenario (orange), best-case scenario (blue) are shown.



Ethical Review

Our research protocol was submitted to the Atrium Health Institutional Review Board (IRB) prior to execution, and the study was deemed exempt from IRB oversight. In compliance with HIPAA (Health Insurance Portability and Accountability Act) regulations, individual patient information was not disclosed, and all data have been deidentified and reported as aggregates. The procedures set out in this protocol, pertaining to the conduct, evaluation, and documentation of this study, were designed to ensure that the investigators abide by Good Clinical Practice guidelines and under the guiding principles detailed in the Declaration of Helsinki.

Results

Estimation and Inference

Our model was specified as a VECM with 7 lags in its VAR representation ($p=7$), 1 cointegration relationship ($r=1$), and a restricted constant parameter μ so that the series would not have

linear trend. The AIC scores of VAR models with a varying number of lags from 2 to 14 were inconclusive. However, we found that 7 lags were sufficient to account for all the correlation in the data, as evidenced by the autocorrelation function and cross-correlation function plots of the residuals (Figure 3). The Johansen trace test indicated that there was 1 cointegration relationship (significant at 1%, based on tabulated critical values). Finally, the likelihood ratio test for linear trend indicated that there was no linear trend in the data ($P=.32$). Furthermore, the restricted model had a lower AIC score than the unrestricted model (the AIC scores were -1519 and -1516 , respectively).

The output from the maximum likelihood estimation showed that the cointegration relationship, that is, the error correction term, had a significant negative effect on census change ($P<.001$); no significant effect was observed for incidence change ($P=.26$) (Table 1). The long-run cointegration relationship was estimated as:

$$ect_{t-1} = census_{t-1} - 0.8013incidence_{t-1} + 7.8266$$

where ect_{t-1} was the (lagged) error correction term. Table 1 also shows that past changes in census and incidence also had meaningful effects on current census change. Past census changes had significant effect at lag 2 ($P=.002$). Past incidence changes had significant effects at lag 1 ($P=.005$), lag 2 ($P=.04$), lag 4 ($P=.02$), lag 5 ($P=.03$), and lag 6 ($P=.02$).

From Table 2, there were some significant seasonal effects, that is, differences in both census and incidence changes among days of the week. Compared to Thursday, census change was higher on Monday and incidence change was lower on Sunday, with significant differences ($P=.01$ and $P=.002$, respectively).

Figure 3. Autocorrelation functions and cross-correlation functions of the residuals: (A) census residuals, (B) lagged census residuals and incidence residuals, (C) census residuals and lagged incidence residuals, and (D) incidence residuals.

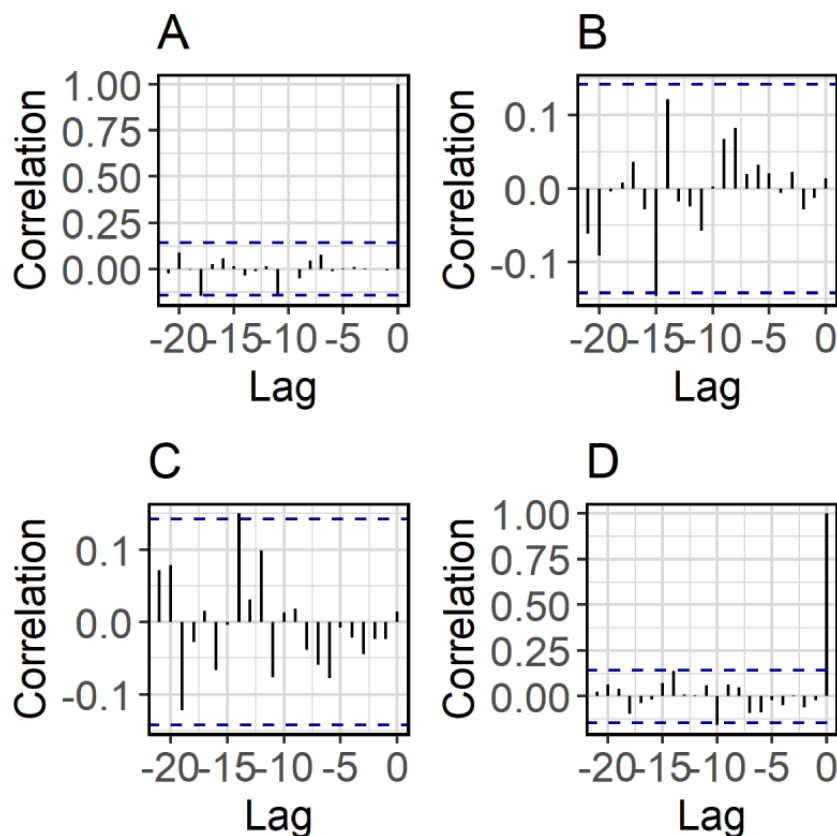


Table 1. Parameter estimates and *T* tests for nonseasonal effects.

Predictor	$\Delta Census_t$			$\Delta Incidence_t$		
	Estimate	<i>T</i> statistics	<i>P</i> value	Estimate	<i>T</i> statistics	<i>P</i> value
ect_{t-1}	-0.1265	-5.6993	<.001	-0.1216	-1.1323	.26
$\Delta Census_{t-1}$	-0.0489	-0.7143	.48	0.5487	1.6555	.10
$\Delta Incidence_{t-1}$	-0.0665	-2.8222	.005	-0.9808	-8.6067	<.001
$\Delta Census_{t-2}$	-0.2220	-3.2277	.002	-0.0614	-0.1844	.85
$\Delta Incidence_{t-2}$	-0.0532	-2.0881	.04	-0.6955	-5.6431	<.001
$\Delta Census_{t-3}$	-0.0700	-0.9949	.32	0.0643	0.1890	.85
$\Delta Incidence_{t-3}$	-0.0472	-1.9094	.06	-0.6428	-5.3755	<.001
$\Delta Census_{t-4}$	-0.0785	-1.1224	.26	0.9769	2.8871	.004
$\Delta Incidence_{t-4}$	-0.0567	-2.4165	.02	-0.5564	-4.8999	<.001
$\Delta Census_{t-5}$	-0.0499	-0.7140	.48	-0.0792	-0.2341	.82
$\Delta Incidence_{t-5}$	-0.0465	-2.1907	.03	-0.4589	-4.4634	<.001
$\Delta Census_{t-6}$	0.0077	0.1107	.91	0.4533	1.3404	.18
$\Delta Incidence_{t-6}$	-0.0373	-2.4015	.02	-0.2384	-3.1739	.002

Table 2. Parameter estimates and *T* tests for day-of-the-week effects, in comparison with Thursday being the reference.

Predictor	$\Delta Census_t$			$\Delta Incidence_t$		
	Estimate	<i>T</i> statistics	<i>P</i> value	Estimate	<i>T</i> statistics	<i>P</i> value
Friday	-0.0213	-1.1120	.27	0.0095	0.1024	.92
Saturday	0.0083	0.3980	.69	-0.1528	-1.5176	.13
Sunday	0.0030	0.1330	.89	-0.3340	-3.0744	.002
Monday	0.0585	2.6205	.01	-0.1939	-1.7950	.07
Tuesday	0.0291	1.3896	.17	-0.1284	-1.2655	.21
Wednesday	-0.0037	-0.1895	.85	0.0343	0.3672	.71

Model Diagnostics

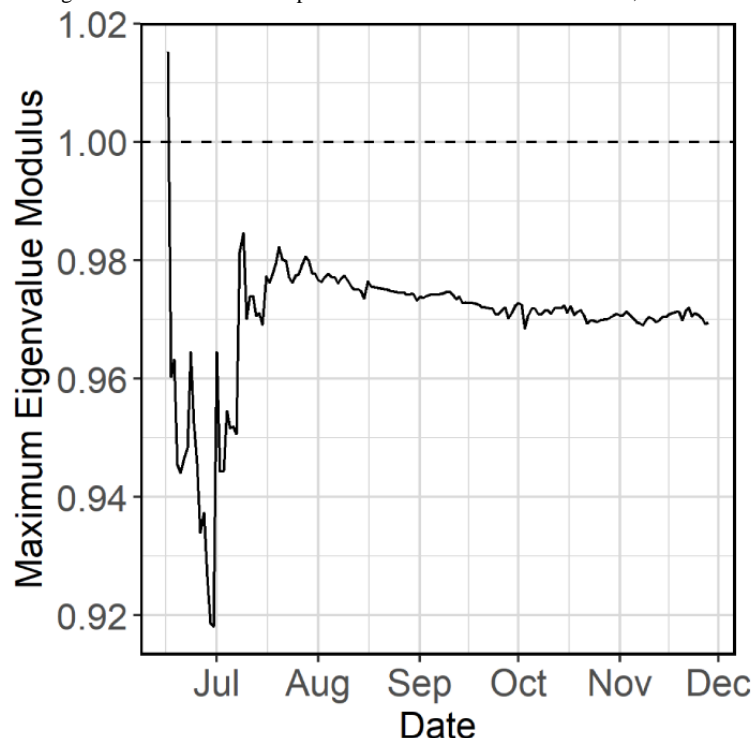
The omnibus *F* tests were significant for both census ($P<.001$) and incidence ($P<.001$) components.

The Portmanteau test did not show sufficient evidence that the errors were autocorrelated ($P=.19$). From the residual autocorrelation function and cross-correlation function plots, the correlations were within the 95% confidence band (Figure 3). The Jarque-Bera normality tests failed to reject the normality null hypothesis for the census errors ($P=.71$) but did for incidence ($P<.001$). Specifically, the incidence residuals were moderately left-skewed. The Jarque-Bera multivariate test also rejected the multivariate normality null hypothesis ($P<.001$).

The Augmented Dickey-Fuller test for stationarity of the error correction term rejected the unit root null hypothesis at the 10% significance level but failed to reject the null hypothesis at the 5% significance level (based on tabulated critical values). The KPSS test failed to reject the stationarity null hypothesis ($P=.10$). Examination of the time plot of the predicted error correction term showed no obvious departure from stationarity.

The companion matrix of the VAR representation had a maximum eigenvalue modulus of 0.97, strictly less than 1. Although this value was close to 1, the trace plot showed that this value had been slowly declining and below 1 across time when the model was fitted repeatedly in a daily rolling basis from June 16 to November 28 (Figure 4).

Figure 4. Trace plot of the maximum eigenvalue modulus for the period from June 16 to November 28, 2020.



Forecast Performance

We obtained the approximate sampling distribution of the out-of-sample MAPE from the time-series cross-validation (Figure 5). The typical value (median) of MAPE was 5.9% and the 95th percentile of MAPE was 13.4%. For the sake of comparison, the corresponding values from an ARIMA model

using the COVID-19 hospital census only were 6.6% and 14.3%. Additionally, after fitting the data from May 15 to December 5, we forecasted the census out to 7 days. Subsequently, the actual values were accurately forecasted with a MAPE of 1.9% and were all within the 80% bootstrapped forecast intervals (Figure 6).

Figure 5. Distribution of the 7-days-ahead mean absolute percentage error from the time-series cross-validation for the period from June 16 to November 28, 2020. Median (blue) and 95th percentile (red) are shown.

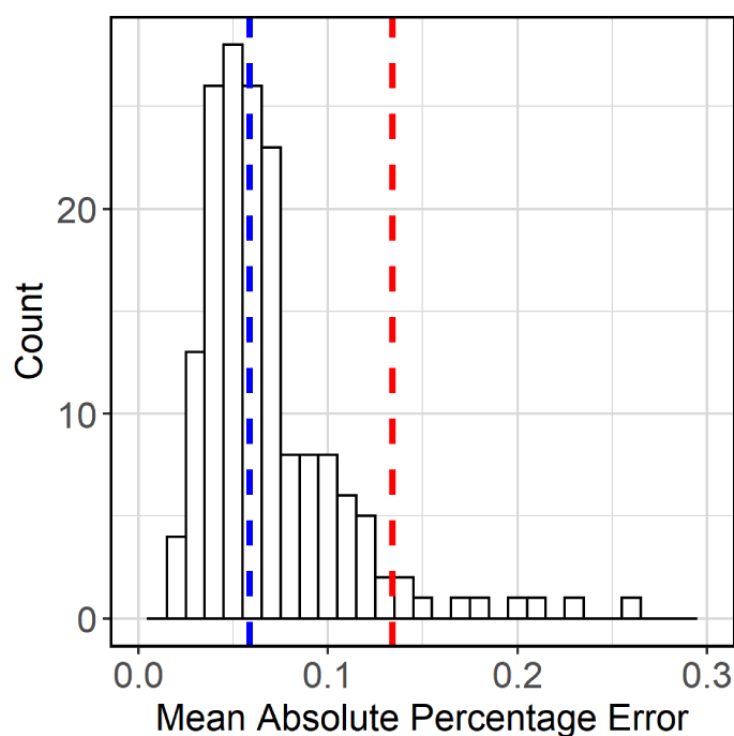
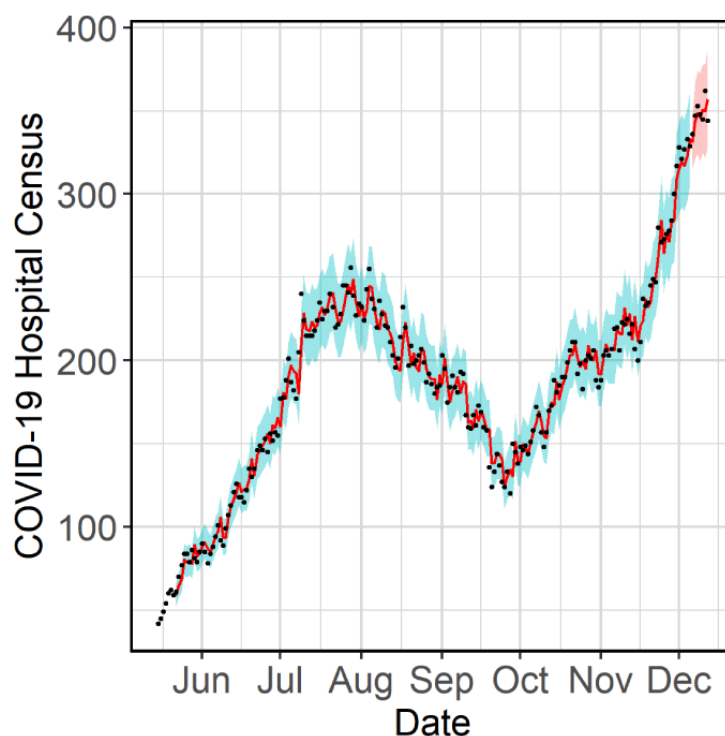


Figure 6. One-step-ahead in-sample and 7-days-ahead out-of-sample predictions for COVID-19 hospital census in the Cities Readiness Initiative region. True values (black), in-sample and out-of-sample predictions (red line), 95% prediction intervals (blue band), 80% forecast intervals (red band) are shown. The model is fitted on data from May 15 to December 5, 2020.

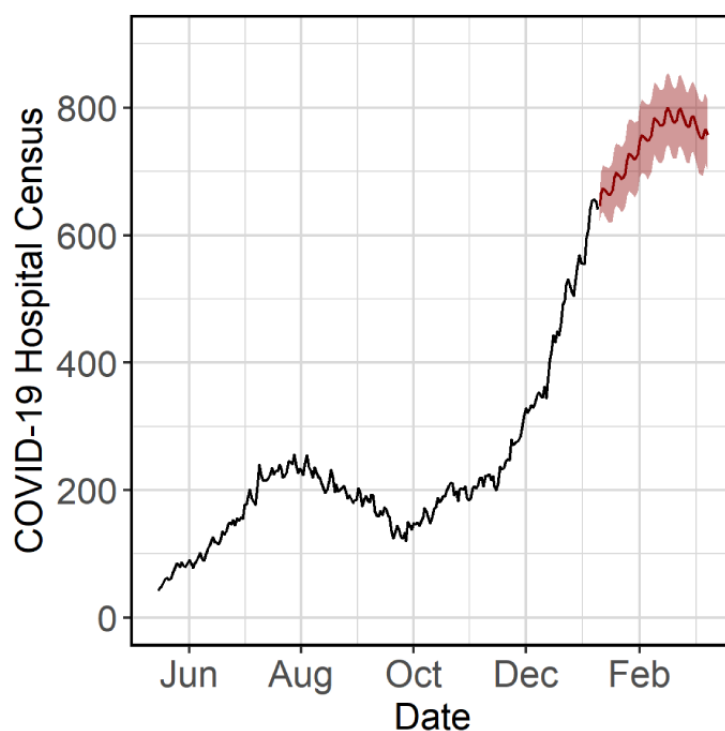


Scenario-Based Long-Term Forecasting

In all scenarios, due to cointegration, census followed corresponding concave trajectories with peaks occurring approximately 2 to 3 weeks later than incidence depending on

the scenario. In the worst-case scenario, census was projected to peak on February 16, 2021 (11 days later than incidence), with approximately 850 patients at the 80% forecast interval upper bound (Figure 7).

Figure 7. Worst-case-scenario, 60-day forecasts for COVID-19 hospital census in the Cities Readiness Initiative region, as of January 9, 2021. Past values (black), forecasts (red line), and 80% forecast intervals (red band) are shown.



Discussion

Principal Results

Our VECM provides a very good fit to the data and outperforms models with no or other leading indicators. Significant omnibus *F* tests showed that the model fit was better than that of a reduced VECM representation with no predictors (ie, a bivariate random walk model). When we examined model diagnostics, there was no sign of any serious departure from model assumptions. From the Portmanteau test, the errors were not different from white noise (ie, the errors do not exhibit serial correlation). Although the normality assumption (for incidence) was not met, the asymptotic properties of our estimation and hypothesis tests in the VECM would not be affected [33]. To address the possible effect of this violation on the forecast intervals, we implemented a bootstrap procedure for the forecast intervals. Both the ADF test and KPSS test showed reasonable evidence that the long-run relationship was stable. With the maximum eigenvalue modulus of the VAR representation consistently below 1 across time, the model itself was quite stable. Examining the day-of-the-week effects, we observed a higher increase in census at the beginning of the week. This agrees with our observations of hospital operations and suggests higher resource allocation when starting the week, as is also reflected in the forecasts (Figure 7). In terms of forecast performance, the VECM yielded a smaller MAPE, in terms of the median and the 95th percentile, when compared to an ARIMA model using the COVID-19 hospital census only. Our VECM also performed better than another VECM that uses two internet-based leading indicators (median MAPE of 10.5%), albeit on time domains that were partially overlapping [13].

The long-run relationship plays a crucial role in the model. Our model results show how future census responds to perturbations in the long-run cointegration relationship in the direction that would preserve the stability of the relationship. For instance, if incidence increases significantly and drives the error correction term below 0, the next-day census will tend to increase so that the error correction term will move back toward 0. Compared to short-run relationships between census change and past changes in incidence and census, the long-run relationship effect is also strongly significant and is a major driver in the model.

We observed that local infection incidence led the hospital census by about 2 weeks. The cross-correlations between incidence and census were uniformly high, between 0.7 and 0.8 at different lags, but the highest correlation was at lag 14. Clinically, we know that after someone is diagnosed with SARS-CoV-2, it can take several days before they become sick enough to be hospitalized. During the summer 2020 wave of the pandemic, incidence peaked 18 days earlier, on July 10, than when census peaked, on July 28. In the model, we also saw that past incidence changes at multiple lags have statistically significant effects on census. While previous studies have focused on other types of leading indicators [12,13], our model results and our observations demonstrate that local infection incidence can be a very effective leading indicator for COVID-19 hospital census.

Applying the model to scenario-based forecasting in a health care system is an important method for long-term forecasting when approaching an infection prevalence peak and helps determine the potential for resource capacity to be exceeded under a worst-case scenario. There are several advantages to our approach. With a scenario-based and epidemiologically informed approach, the VECM produces realistic, nonlinear, long-range trajectories of census. In contrast, an ARIMA model can have an upward linear trajectory even as we approach and arrive at the infection prevalence peak because it is agnostic to incidence. Hence, the VECM fit with scenario-based incidence will provide better accuracy since it is more reflective of pandemic behavior. Additionally, when the concern is a specific scenario, our approach is particularly useful at minimizing long-range forecast uncertainty, since the bootstrapped sample paths are constrained to fluctuate around the marginalized scenario-based census projection. Without such a constraint, 60-day forecasts can typically have wide forecast intervals that are of no practical utility.

Our study has mathematically ascertained the stable long-run relationship, that is, cointegration, between the COVID-19 hospital census and the local infection incidence, and we have developed a statistical incidence-based model to forecast the COVID-19 hospital census. In comparison, prior COVID-19 hospital capacity planning models that make use of infection incidence data rely on simplified assumptions about the incidence-census relationship. For example, in the COVID-19 Hospital Impact Model for Epidemics (CHIME) at the University of Pennsylvania [34], the ratio between hospital admissions and infection incidence is a scenario parameter defined by the user and is not time varying.

Limitations

Although our model has been thoroughly developed, it is not free of limitations. First, it is possible that we may lose the stable long-run relationship at some point in the future, either because it has run its course or due to structural changes in the time series. For instance, in the latter case, inadequate community-based testing might suddenly underestimate the actual local infection incidence, and there may be a level shift in the relationship that would have to be accounted for by a modified VECM [35,36]. In other cases, more complex structural changes may arise and be challenging to model. Second, in the future, other regions may find that the ratio between asymptomatic and symptomatic cases fluctuates considerably over time. Because case severity affects the time to hospitalization, this situation may require model revision. A potential remedy is to include both the number of asymptomatic and symptomatic cases as two leading indicators with census in a VECM in the hopes that some cointegration exists among the three variables. Third, it is relatively more difficult to fit a VECM. For univariate models such as ARIMA and exponential smoothing, well-developed R packages exist for automated model specification and estimation. With the VECM, more deliberate modeling decisions and careful checking of assumptions need to be made to fit a reliable model. Finally, the inclusion of seasonal effects in our model requires that the seasonality is deterministic. However, another health care system may find that their time-series data have stochastic seasonality

or multiple deterministic seasonality. If seasonality is not important, we potentially may resolve this by simply deseasonalizing the series. Otherwise, it may be possible to account for this with more advanced parameterization of the seasonal effects.

Conclusions

The construct presented here provides a framework in the context of a health care system for incorporating other leading indicators that may yield further increases in forecasting performance. For instance, the VECM that uses internet-based leading indicators [13] could potentially be improved by including incidence. It is also possible to incorporate other nested hospital-related time series, such as the number of intensive care units and the number of ventilators, into the

VECM if there was a need to simultaneously forecast other resources. Additionally, a VECM could be a valuable candidate for a model-averaged ensemble. This can be particularly useful if the ensemble consists only of agnostic univariate time-series models.

We have shown that infection incidence can be successfully tethered with hospital census in a multivariate time-series model to achieve accurate forecasting of COVID-19 hospital census. When coupled with scenario-based forecasting, the model helped our leaders evaluate resource capacity against different possible peak resource demands. In hindsight, our analyses correctly assured our leaders of our capability to handle a worst-case scenario, alleviated uncertainty, and effectively guided long-term planning of adequate staffing, bed capacity, and equipment supplies through the pandemic.

Authors' Contributions

HMN prepared the original draft. HMN and PJT were involved in study conceptualization, statistical analysis, and review and editing of the manuscript. ADM supervised the study and contributed to the review and editing of the manuscript.

Conflicts of Interest

ADM is an administrative member of iEnroll LLC.

References

1. Cheng ZJ, Shan J. 2019 Novel coronavirus: where we are and what we know. *Infection* 2020 Apr;48(2):155-163 [FREE Full text] [doi: [10.1007/s15010-020-01401-y](https://doi.org/10.1007/s15010-020-01401-y)] [Medline: [32072569](#)]
2. Singhal T. A Review of Coronavirus Disease-2019 (COVID-19). *Indian J Pediatr* 2020 Apr 13;87(4):281-286 [FREE Full text] [doi: [10.1007/s12098-020-03263-6](https://doi.org/10.1007/s12098-020-03263-6)] [Medline: [32166607](#)]
3. WHO Director-General's statement on IHR Emergency Committee on Novel Coronavirus (2019-nCoV). World Health Organization. 2020 Jan 30. URL: [https://www.who.int/director-general/speeches/detail/who-director-general-s-statement-on-ihr-emergency-committee-on-novel-coronavirus-\(2019-ncov\)](https://www.who.int/director-general/speeches/detail/who-director-general-s-statement-on-ihr-emergency-committee-on-novel-coronavirus-(2019-ncov)) [accessed 2020-12-29]
4. Covid-19 Map. Johns Hopkins Coronavirus Resource Center. URL: <https://coronavirus.jhu.edu/map.html> [accessed 2020-12-29]
5. The Lancet. India's COVID-19 emergency. *The Lancet* 2021 May;397(10286):1683. [doi: [10.1016/s0140-6736\(21\)01052-7](https://doi.org/10.1016/s0140-6736(21)01052-7)]
6. Earnest A, Chen MI, Ng D, Sin LY. Using autoregressive integrated moving average (ARIMA) models to predict and monitor the number of beds occupied during a SARS outbreak in a tertiary hospital in Singapore. *BMC Health Serv Res* 2005 May 11;5(1):36 [FREE Full text] [doi: [10.1186/1472-6963-5-36](https://doi.org/10.1186/1472-6963-5-36)] [Medline: [15885149](#)]
7. Jones S, Thomas A, Evans R, Welch S, Haug P, Snow G. Forecasting daily patient volumes in the emergency department. *Acad Emerg Med* 2008 Feb;15(2):159-170 [FREE Full text] [doi: [10.1111/j.1553-2712.2007.00032.x](https://doi.org/10.1111/j.1553-2712.2007.00032.x)] [Medline: [18275446](#)]
8. Capan M, Hoover S, Jackson E, Paul D, Locke R. Time Series Analysis for Forecasting Hospital Census: Application to the Neonatal Intensive Care Unit. *Appl Clin Inform* 2017 Dec 16;07(02):275-289. [doi: [10.4338/aci-2015-09-ra-0127](https://doi.org/10.4338/aci-2015-09-ra-0127)]
9. Zhou L, Zhao P, Wu D, Cheng C, Huang H. Time series model for forecasting the number of new admission inpatients. *BMC Med Inform Decis Mak* 2018 Jun 15;18(1):39 [FREE Full text] [doi: [10.1186/s12911-018-0616-8](https://doi.org/10.1186/s12911-018-0616-8)] [Medline: [29907102](#)]
10. Koestler DC, Ombao H, Bender J. Ensemble-based methods for forecasting census in hospital units. *BMC Med Res Methodol* 2013 May 30;13(1):67 [FREE Full text] [doi: [10.1186/1471-2288-13-67](https://doi.org/10.1186/1471-2288-13-67)] [Medline: [23721123](#)]
11. Jones SS, Evans RS, Allen TL, Thomas A, Haug PJ, Welch SJ, et al. A multivariate time series approach to modeling and forecasting demand in the emergency department. *J Biomed Inform* 2009 Feb;42(1):123-139 [FREE Full text] [doi: [10.1016/j.jbi.2008.05.003](https://doi.org/10.1016/j.jbi.2008.05.003)] [Medline: [18571990](#)]
12. Berta P, Paruolo P, Verzillo S, Lovaglio PG. A bivariate prediction approach for adapting the health care system response to the spread of COVID-19. *PLoS One* 2020 Oct 15;15(10):e0240150 [FREE Full text] [doi: [10.1371/journal.pone.0240150](https://doi.org/10.1371/journal.pone.0240150)] [Medline: [33057389](#)]
13. Turk P, Tran T, Rose G, McWilliams A. A Predictive Internet-Based Model for COVID-19 Hospitalization Census. medRxiv. Preprint posted online November 18, 2020. [FREE Full text] [doi: [10.1101/2020.11.15.20231845](https://doi.org/10.1101/2020.11.15.20231845)]
14. Lynch CJ, Gore R. Short-Range Forecasting of COVID-19 During Early Onset at County, Health District, and State Geographic Levels Using Seven Methods: Comparative Forecasting Study. *J Med Internet Res* 2021 Mar 23;23(3):e24925 [FREE Full text] [doi: [10.2196/24925](https://doi.org/10.2196/24925)] [Medline: [33621186](#)]

15. Singh RK, Rani M, Bhagavathula AS, Sah R, Rodriguez-Morales AJ, Kalita H, et al. Prediction of the COVID-19 Pandemic for the Top 15 Affected Countries: Advanced Autoregressive Integrated Moving Average (ARIMA) Model. *JMIR Public Health Surveill* 2020 May 13;6(2):e19115 [[FREE Full text](#)] [doi: [10.2196/19115](#)] [Medline: [32391801](#)]
16. Zeng C, Zhang J, Li Z, Sun X, Olatosi B, Weissman S, et al. Spatial-Temporal Relationship Between Population Mobility and COVID-19 Outbreaks in South Carolina: Time Series Forecasting Analysis. *J Med Internet Res* 2021 Apr 13;23(4):e27045 [[FREE Full text](#)] [doi: [10.2196/27045](#)] [Medline: [33784239](#)]
17. Yeung AY, Roewer-Despres F, Rosella L, Rudzicz F. Machine Learning-Based Prediction of Growth in Confirmed COVID-19 Infection Cases in 114 Countries Using Metrics of Nonpharmaceutical Interventions and Cultural Dimensions: Model Development and Validation. *J Med Internet Res* 2021 Apr 23;23(4):e26628 [[FREE Full text](#)] [doi: [10.2196/26628](#)] [Medline: [33844636](#)]
18. Mehta M, Julaiti J, Griffin P, Kumara S. Early Stage Machine Learning-Based Prediction of US County Vulnerability to the COVID-19 Pandemic: Machine Learning Approach. *JMIR Public Health Surveill* 2020 Sep 11;6(3):e19446 [[FREE Full text](#)] [doi: [10.2196/19446](#)] [Medline: [32784193](#)]
19. Cleveland R, Cleveland W, McRae J, Terpenning I. STL: A Seasonal-Trend Decomposition Procedure Based on Loess. *Journal of Official Statistics* 1990;6:3-73 [[FREE Full text](#)]
20. Pfaff B. *Analysis of Integrated and Cointegrated Time Series With R*, 2nd Ed. New York, NY: Springer; 2008.
21. Akaike H. A new look at the statistical model identification. *IEEE Trans Automat Contr* 1974 Dec;19(6):716-723. [doi: [10.1109/TAC.1974.1100705](#)]
22. Johansen S. Estimation and Hypothesis Testing of Cointegration Vectors in Gaussian Vector Autoregressive Models. *Econometrica* 1991 Nov;59(6):1551. [doi: [10.2307/2938278](#)]
23. Johansen S. *Likelihood-Based Inference in Cointegrated Vector Autoregressive Models*. Oxford; New York: Oxford University Press; 1995.
24. Johansen S, Juselius K. Maximum Likelihood Estimation and Inference on Cointegration - With Applications to the Demand for Money. *Oxford Bulletin of Economics and Statistics* 2009;52:169-210. [doi: [10.1111/j.1468-0084.1990.mp52002003.x](#)]
25. Hyndman R, Athanasopoulos G. *Forecasting: Principles and Practice*, 3rd Edition. Melbourne, Australia: OTexts; 2021. URL: <https://otexts.com/fpp3/> [accessed 2021-07-27]
26. Incidence-Census-Model. GitHub. URL: <https://github.com/hmnguye/Incidence-Census-Model> [accessed 2021-07-28]
27. Jarque CM, Bera AK. A Test for Normality of Observations and Regression Residuals. *International Statistical Review / Revue Internationale de Statistique* 1987 Aug;55(2):163. [doi: [10.2307/1403192](#)]
28. Said SE, Dickey DA. Testing for unit roots in autoregressive-moving average models of unknown order. *Biometrika* 1984;71(3):599-607. [doi: [10.1093/biomet/71.3.599](#)]
29. Kwiatkowski D, Phillips PC, Schmidt P, Shin Y. Testing the null hypothesis of stationarity against the alternative of a unit root. *Journal of Econometrics* 1992 Oct;54(1-3):159-178. [doi: [10.1016/0304-4076\(92\)90104-y](#)]
30. Hamilton J. *Time Series Analysis*. Princeton, NJ: Princeton University Press; 1994.
31. Wang L, Zhou Y, He J, Zhu B, Wang F, Tang L. An epidemiological forecast model and software assessing interventions on COVID-19 epidemic in China. *Journal of Data Science* 2020;18(3):409-432. [doi: [10.6339/jds.202007_18\(3\).0003](#)]
32. Gardner ES, McKenzie E. Forecasting Trends in Time Series. *Management Science* 1985 Oct;31(10):1237-1246. [doi: [10.1287/mnsc.31.10.1237](#)]
33. Johansen S. Cointegration: Overview and Development. In: Mikosch T, Kreiß JP, Davis RA, Andersen TG, editors. *Handbook of Financial Time Series*. Berlin, Heidelberg: Springer; 2009:671-693.
34. Weissman GE, Crane-Droesch A, Chivers C, Luong T, Hanish A, Levy MZ, et al. Locally Informed Simulation to Predict Hospital Capacity Needs During the COVID-19 Pandemic. *Annals of Internal Medicine* 2020 Jul 07;173(1):21-28. [doi: [10.7326/m20-1260](#)]
35. Saikkonen P, Lütkepohl H, Lutkepohl H. Testing for the Cointegrating Rank of a VAR Process with Structural Shifts. *Journal of Business & Economic Statistics* 2000 Oct;18(4):451. [doi: [10.2307/1392226](#)]
36. Lutkepohl H, Saikkonen P, Trenkler C. Testing for the Cointegrating Rank of a VAR Process with Level Shift at Unknown Time. *Econometrica* 2004 Mar;72(2):647-662. [doi: [10.1111/j.1468-0262.2004.00505.x](#)]

Abbreviations

ADF: augmented Dickey-Fuller
AIC: Akaike information criterion
ARIMA: autoregressive integrated moving average
CRI: Cities Readiness Initiative
HIPAA: Health Insurance Portability and Accountability Act
IRB: Institutional Review Board
KPSS: Kwiatkowski-Phillips-Schmidt-Shin
MAPE: mean absolute percentage error
SARIMA: seasonal autoregressive integrated moving average

SARS: severe acute respiratory syndrome

STL: seasonal and trend decomposition using Loess

VAR: vector autoregressive

VECM: vector error correction model

Edited by T Sanchez; submitted 24.02.21; peer-reviewed by E Mahmoudi, R Gore; comments to author 02.06.21; revised version received 22.06.21; accepted 29.06.21; published 04.08.21.

Please cite as:

Nguyen HM, Turk PJ, McWilliams AD

Forecasting COVID-19 Hospital Census: A Multivariate Time-Series Model Based on Local Infection Incidence

JMIR Public Health Surveill 2021;7(8):e28195

URL: <https://publichealth.jmir.org/2021/8/e28195>

doi: [10.2196/28195](https://doi.org/10.2196/28195)

PMID: [34346897](https://pubmed.ncbi.nlm.nih.gov/34346897/)

©Hieu M Nguyen, Philip J Turk, Andrew D McWilliams. Originally published in JMIR Public Health and Surveillance (<https://publichealth.jmir.org>), 04.08.2021. This is an open-access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR Public Health and Surveillance, is properly cited. The complete bibliographic information, a link to the original publication on <https://publichealth.jmir.org>, as well as this copyright and license information must be included.

Original Paper

Census Tract Patterns and Contextual Social Determinants of Health Associated With COVID-19 in a Hispanic Population From South Texas: A Spatiotemporal Perspective

Cici Bauer¹, PhD; Kehe Zhang¹, MS; Miryoung Lee², PhD; Susan Fisher-Hoch², MD; Esmeralda Guajardo³, MA; Joseph McCormick², MD; Isela de la Cerda², MS; Maria E Fernandez⁴, PhD; Belinda Reininger⁵, DrPH

¹Department of Biostatistics and Data Science, School of Public Health, The University of Texas Health Science Center at Houston, Houston, TX, United States

²Department of Epidemiology, Human Genetics and Environmental Science, School of Public Health, The University of Texas Health Science Center at Houston, Brownsville, TX, United States

³Cameron County Public Health, San Benito, TX, United States

⁴Department of Health Promotion and Behavior Sciences, School of Public Health, The University of Texas Health Science Center at Houston, Houston, TX, United States

⁵Department of Health Promotion and Behavior Sciences, School of Public Health, The University of Texas Health Science Center at Houston, Brownsville, TX, United States

Corresponding Author:

Cici Bauer, PhD

Department of Biostatistics and Data Science

School of Public Health

The University of Texas Health Science Center at Houston

1200 Pressler Street

Houston, TX, 77030

United States

Phone: 1 713 500 9581

Email: cici.x.bauer@uth.tmc.edu

Related Article:

This is a corrected version. See correction statement: <https://publichealth.jmir.org/2021/8/e32870>

Abstract

Background: Previous studies have shown that various social determinants of health (SDOH) may have contributed to the disparities in COVID-19 incidence and mortality among minorities and underserved populations at the county or zip code level.

Objective: This analysis was carried out at a granular spatial resolution of census tracts to explore the spatial patterns and contextual SDOH associated with COVID-19 incidence from a Hispanic population mostly consisting of a Mexican American population living in Cameron County, Texas on the border of the United States and Mexico. We performed age-stratified analysis to identify different contributing SDOH and quantify their effects by age groups.

Methods: We included all reported COVID-19-positive cases confirmed by reverse transcription-polymerase chain reaction testing between March 18 (first case reported) and December 16, 2020, in Cameron County, Texas. Confirmed COVID-19 cases were aggregated to weekly counts by census tracts. We adopted a Bayesian spatiotemporal negative binomial model to investigate the COVID-19 incidence rate in relation to census tract demographics and SDOH obtained from the American Community Survey. Moreover, we investigated the impact of local mitigation policy on COVID-19 by creating the binary variable “shelter-in-place.” The analysis was performed on all COVID-19-confirmed cases and age-stratified subgroups.

Results: Our analysis revealed that the relative incidence risk (RR) of COVID-19 was higher among census tracts with a higher percentage of single-parent households (RR=1.016, 95% posterior credible intervals [CIs] 1.005, 1.027) and a higher percentage of the population with limited English proficiency (RR=1.015, 95% CI 1.003, 1.028). Lower RR was associated with lower income (RR=0.972, 95% CI 0.953, 0.993) and the percentage of the population younger than 18 years (RR=0.976, 95% CI 0.959, 0.993). The most significant association was related to the “shelter-in-place” variable, where the incidence risk of COVID-19

was reduced by over 50%, comparing the time periods when the policy was present versus absent (RR=0.506, 95% CI 0.454, 0.563). Moreover, age-stratified analyses identified different significant contributing factors and a varying magnitude of the “shelter-in-place” effect.

Conclusions: In our study, SDOH including social environment and local emergency measures were identified in relation to COVID-19 incidence risk at the census tract level in a highly disadvantaged population with limited health care access and a high prevalence of chronic conditions. Results from our analysis provide key knowledge to design efficient testing strategies and assist local public health departments in COVID-19 control, mitigation, and implementation of vaccine strategies.

(*JMIR Public Health Surveill* 2021;7(8):e29205) doi:[10.2196/29205](https://doi.org/10.2196/29205)

KEYWORDS

COVID-19; spatial pattern; social determinants of health; Bayesian; underserved population; health inequity

Introduction

COVID-19, which comes from SARS-CoV-2, has caused death, health care system stress, and global economic instability. In the United States, it also has disproportionately affected minority and underserved populations, where COVID-19 infection and fatality rates are significantly higher among African American and Hispanic populations [1-3]. Previous studies have shown various social determinants of health (SDOH) that may explain the disparity in COVID-19 incidence and mortality in ethnic and racial minorities [1,4,5].

The differential impact of COVID-19 on minorities and other groups facing health inequities has been described and underscores a critical need to target these underserved groups. However, the majority of these studies in the United States used aggregated county-level data from the COVID Tracking Project [6]. The geographical scale of the US county often lacks granularity to reveal the local spatial pattern and detect local hot spots (ie, areas with excessive infection rates). Moreover, the high variability of SDOH within a county population was not able to accurately examine the impact of SDOH on COVID-19 disparities in populations [5]. Studies that investigate the SDOH and COVID-19 incidence and mortality at a geographical scale smaller than the US county are limited [7,8]. The lack of studies on a granular spatial scale is largely due to insufficiently detailed COVID-19 surveillance data, particularly data that are publicly available.

In this study, we investigated the contextual SDOH and their potential association with COVID-19 incidence at the census tract level. The study population consists of a Hispanic population with mostly Mexican American people living in South Texas on the US-Mexico border. The Mexican American

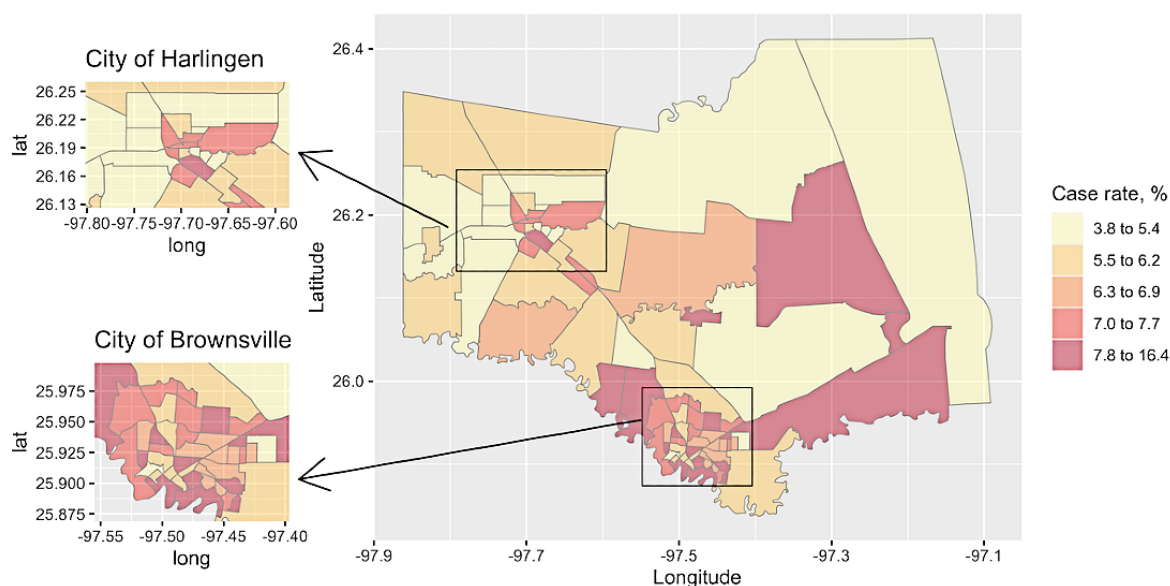
population are the largest and fastest-growing Hispanic subgroup in the United States and among those with low socioeconomic status compared to other ethnic groups in the nation [9]. The population in our study has high prevalence of obesity and diabetes [10]; both pre-existing conditions increase the risk of severe COVID-19 outcomes [11]. Our analysis provided a look at the SDOH at sufficient spatial granularity to detect local trends and hot spots for COVID-19 monitoring and control. Results from our study have informed the intervention strategies to increase COVID-19 testing uptake in underserved populations and the design of interventions and targeted vaccination programs.

Methods

Study Population

Our study population is from Cameron County, Texas with a current population of 423,163 and over 90% Hispanics [12], where the vast majority were Mexican-Hispanic [13]. Most Cameron County residents are uninsured (~29%) and live below the poverty line (~33%) [12]; additional research based on a well-documented cohort from this region estimated that around 52% of the population does not have any private or public health insurance coverage [10,13]. This population, similar to many others living in the South Texas region, also has a high prevalence of type 2 diabetes (over 27%) and obesity (over 50%) [10,13,14]. In our analysis, we included a total of 84 census tracts within Cameron County, as shown in Figure 1. The two largest cities in Cameron County are the City of Brownsville (population 183,677) on the US-Mexico border and the City of Harlingen (population 65,074) 20 miles north of Brownsville, together comprising 59% of the county population.

Figure 1. Choropleth map presenting the cumulative COVID-19 infection rate by census tract between March 18, 2020, and December 16, 2020, in Cameron County, TX. The two largest cities are the city of Brownsville on the border of the United States and Mexico (bottom left panel) and the city of Harlingen (top left panel).



COVID-19 Reported Cases

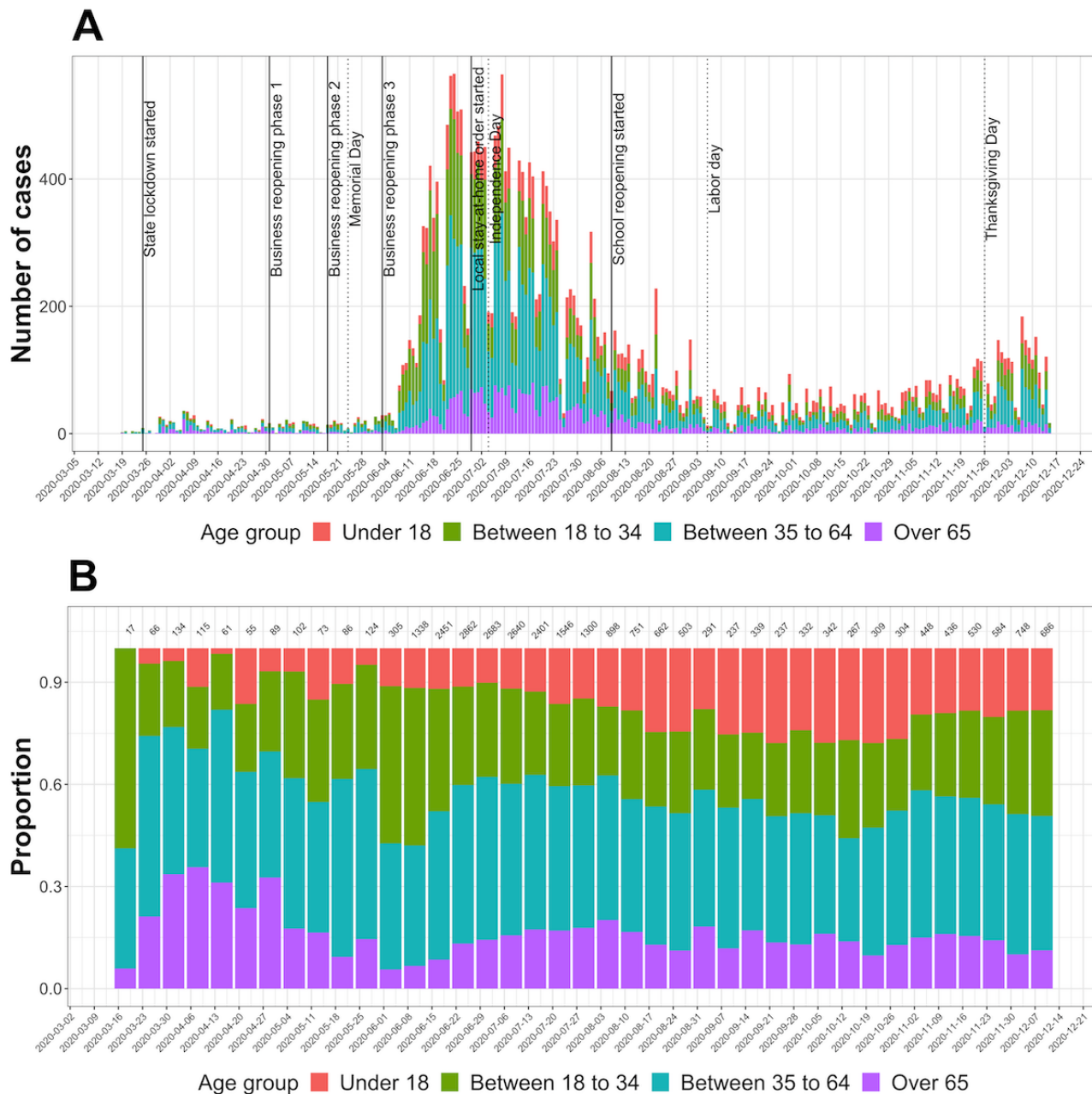
The first confirmed COVID-19 case in Cameron County, Texas was reported on March 18, 2020. By December 16, a total of 28,111 cases had been reported. The cumulative case rate, calculated as the number of positive cases per 100 people, was 1.93% in May 2020 and increased to 6.64% by December 2020, when the cumulative case rate of the general US population in December 2020 was 5.16% [6]. The case-fatality rate in Cameron County was 4% compared to 1.98% in Texas by December.

To facilitate local COVID-19 control and mitigation, Cameron County Public Health Department, the City of Brownsville, and the University of Texas School of Public Health (UTHealth) formed a collaborative group soon after the first COVID-19 case was reported in March 2020. Cameron County Public Health Department maintains a local database of reported and confirmed COVID-19 cases, which were concurrently reported to the Texas Department of State Health Services through the National Electronic Disease Surveillance System. Researchers from UTHealth were given access to the local database and provided data management and analytical support to investigate the trends and risk factors associated with COVID-19 spread. This study was approved by the UTHealth Committee for the Protection of Human Subjects (HSC-SPH-20-1372) and the Data User Agreement between the UTHealth and Cameron County.

During the initial stage of the COVID-19 pandemic, individual-level information associated with each reported case

included age, gender, ethnicity, race, residential address, and specimen collection date. Over time, additional variables were added to the database such as testing type, testing facility, case definition (ie, confirmed or probable), type of exposure, date of recovery, date of death, comorbidities, hospitalization, symptoms, and symptoms onset date. In this analysis, we included all reported COVID-19-positive cases confirmed by reverse transcription-polymerase chain reaction testing based on a sample collection date between March 18 and December 16, 2020, and with a reported residential address within Cameron County. We developed a geocoding algorithm that extracted residential address information and then obtained the corresponding census tract information using the Google application programming interface (API) and the Census Bureau API in R (R Foundation for Statistical Computing) [15,16]. Among the total of 28,111 cases, we were able to geocode 27,733 cases and obtained their census tract information. Of these, 27,731 cases had information on sex, with 14,903 (53.8%) females and 12,824 (46.2%) males. Of the 27,726 cases (missing 1.37%) with age information, 15% (n=4148) were younger than 18 years, 28% (n=7770) were between age 18 to 34 years, 42.7% (n=11,843) were between age 35 to 64 years, and 14.2% (n=3965) were 65 years and older. The age strata range was chosen based on the US Centers for Disease Control and Prevention (CDC) COVID-19 case reporting [6], with some age groups collapsed due to small case numbers. Figure 2 presented weekly confirmed cases stratified by these age groups during the study time.

Figure 2. Temporal pattern of COVID-19–confirmed cases by age groups in Cameron County, Texas between March 18 and December 16, 2020. Panel (A) presents the weekly counts by age groups, along with the event timeline of the state or local COVID-19 mitigation and control policies (solid line) and holidays (dashed line). Panel (B) presents the relative proportions of the weekly cases by age groups, where total weekly counts are shown at the top margin.



Demographic and Social Determinants of Health

Census tract demographic and SDOH variables for Cameron County were obtained from 2013 to 2018 US Census Bureau American Community Survey (ACS) 5-year estimates. These variables included total population, unemployment (%), racial minority (%), poverty level (% living under poverty), education level (% with no high school diploma), income (per capita income in dollars), insurance (% of population uninsured), living conditions (% renters and % living in crowded housing), and transportation (% without vehicles). We also created a population density variable for the census tracts, calculated as the population size per kilometer squared (km^2), ranging from 17 to 1360 per km^2 . We observed substantial spatial variation

of these demographic and SDOH within the Cameron County ([Multimedia Appendix 1](#)).

Shelter-in-Place Indicator Variable

To evaluate the impact of local policy on COVID-19, we created a binary indicator variable with value 1 for time periods when a state or local stay-at-home order was in place, and value 0 otherwise. Mandatory policies of facial coverings, curfew, limitations on gatherings, or beach access closure were present during the shelter-in-place periods [17]. The time period between March 26 and May 1, 2020, corresponded to the presence of the state-level lockdown, at the end of which the phased business reopening began. The local stay-at-home order started from July 1, 2020, and became less restrictive after schools reopened in

mid-August. The event timeline of the policy and holidays is shown in [Figure 2](#) (panel A).

Statistical Analysis

Due to potential reporting lag, we aggregated the number of COVID-19–confirmed cases to weekly counts by census tract. We considered the following Bayesian spatiotemporal model [18,19]. Let Y_{it} denote the number of confirmed cases from census tract i and week t ; we assumed a negative binomial distribution with incidence risk μ_{it} (ie, $Y_{it}|\mu_{it} \sim NB(N_{it}\mu_{it})$, with N_{it} the population size as the offset. The incidence risk was μ_{it} then modeled as follows:

$$\log(\mu_{it}) = \alpha + x_i'\beta + s_t\gamma + \phi_i + \delta_{it}$$

where α was the overall intercept, x_i was the vector of census tract covariates (eg, unemployment and crowded housing) with the associated coefficient vector β . Covariate s_t was the binary policy-in-place indicator previously described. To account for the tract-level spatial dependency, we included a spatial random effect ϕ_i using the intrinsic conditional autoregressive model [20]. The spatiotemporal interaction term δ_{it} captured the unexplained residuals and was assumed an independent and identically distributed normal distribution with variance σ^2 . We reported the relative risk (RR) associated with each covariate, which was calculated as the exponentiated coefficient, along with its 95% posterior credible intervals (CIs). We performed this model on the total COVID-19 cases and then on age subgroups of younger than 18 years, between ages 18 and 34 years, between ages 35 and 64 years, and older than 65 years. All analyses were performed in R [21] and R package INLA [22].

Results

Compared to the US general population, Cameron County has a higher proportion of people who are uninsured (29.1% vs 9.4%), living under poverty (29.6% vs 11.5%), less educated (36.2% with no high school diploma vs 13%), and with worse living conditions (11.8% with crowded housing vs 3.4%). Cameron County is also 90.6% Hispanic, in contrast to 38.3% nationally, and 75% of the population with Spanish as the primary language and 28% having limited English proficiency ([Table 1](#)).

[Figure 2](#) presents the temporal patterns of COVID-19–confirmed cases, in total numbers and by proportion, between March and December 2020. We observed a clear increase in new cases starting in June, that gradually decreased through the end of August. At the beginning of the pandemic in March and April, most cases were from the older population; more cases emerged from the younger population as the pandemic progressed to the

summer. Cases among those 18 years or younger substantially increased from June and peaked in September. Unlike the three waves seen in the US general population, we only observed one prominent wave during the summer, with a smaller second wave after the Thanksgiving holiday.

We fit the Bayesian spatiotemporal negative binomial model previously described to all COVID-19–confirmed cases and then to four age-stratified subgroups (age younger than 18 years, 19–35 years, 36–64 years, 65 years and older), and the results are presented in [Figures 3](#) and [4](#). Of the various demographic and SDOH variables included, the RR of COVID-19 incidence was higher among census tracts with a higher percentage of single-parent households (RR=1.016, 95% CI 1.005, 1.027) and a higher percentage of the population with limited English proficiency (RR=1.015, 95% CI 1.003, 1.028). Lower income was associated with a reduced risk of COVID-19 (RR=0.972, 95% CI 0.953, 0.993) as was the percentage of the population younger than 18 years (RR=0.976, 95% CI 0.959, 0.993). The most striking association was the *shelter-in-place* variable, where the RR of COVID-19 incidence was 0.506 (95% CI 0.454, 0.563) when comparing policy present versus policy absent. This suggests the risk of COVID-19 was reduced by almost 50% when the *shelter-in-place* policy was present.

Age-stratified analyses identified different significant SDOH for each group, and results are presented in [Figure 4](#). For the age group 19 to 34 years, the estimated RR associated with higher percentage of limited English proficiency was 1.025 (95% CI 1.010, 1.040), a higher risk compared to that of the overall population (RR=1.015, 95% CI 1.003, 1.028). Reduced COVID-19 risk was associated with census tracts with higher percentage of no high school education (RR=0.987, 95% CI 0.976, 0.998). For the age group 65 years and older, the percentages of renters and racial minority (ie, percentage of non-Hispanic White) were additional SDOH significantly associated with increased risk of COVID-19 (RR=1.014, 95% CI 1.008, 1.020 and RR=1.018, 95% CI 1.005, 1.032, respectively). The complete results are presented in [Multimedia Appendix 2](#).

The COVID-19 incidence risk was consistently and substantially lower during the time when the “shelter-in-place” policy was present. The effect was the most remarkable for the age group 19 to 35 years, where the risk was reduced by almost 60% when the policy was in place (RR=0.378, 95% CI 0.335, 0.425). For the age group 35 to 65 years, the risk was reduced by almost 50% (RR=0.475, 95% CI 0.424, 0.532). COVID-19 risk reduction was attenuated for the age group 65 years and older (RR=0.690, 95% CI 0.599, 0.793) and the smallest for the age group 18 years or younger (RR=0.767, 95% CI 0.667, 0.881).

Table 1. Summary statistics of the census tract demographics and social determinants of health in Cameron County, Texas and the whole United States. Data were obtained from American Community Survey 2013-2018 5-year estimates.

Variable	Cameron (n=84)	US (n=73,056)
Younger than 18 years (%)		
Mean (CV ^a %)	30.2 (18.9)	22.1 (30.1)
Median (Q1, Q3 ^b)	31.2 (27.5, 33.7)	22.2 (18.5, 26.0)
Older than 65 years (%)		
Mean (CV %)	14.0 (38.3)	16.0 (50.2)
Median (Q1, Q3)	13.2 (10.4, 16.9)	15.2 (11.0, 19.6)
Racial minority (%)		
Mean (CV %)	90.6 (12.1)	38.3 (78.3)
Median (Q1, Q3)	94.7 (86.8, 97.1)	29.7 (12.5, 60.8)
Single-parent household (%)		
Mean (CV %)	14.6 (39.0)	9.3 (69.2)
Median (Q1, Q3)	14.1 (10.1, 18.8)	7.9 (4.8, 12.2)
Disability (%)		
Mean (CV %)	13.5 (30.1)	13.4 (44.0)
Median (Q1, Q3)	13.7 (10.4, 16.2)	12.5 (9.2, 16.6)
Limited English (%)		
Mean (CV %)	27.7 (35.4)	8.0 (135.4)
Median (Q1, Q3)	27.6 (20.4, 35.7)	3.5 (1.1, 10.1)
Unemployed (%)		
Mean (CV %)	3.9 (51.5)	3.9 (68.1)
Median (Q1, Q3)	3.5 (2.4, 5.1)	3.3 (2.1, 4.9)
No high school diploma (%)		
Mean (CV %)	36.2 (39.1)	13.0 (81.2)
Median (Q1, Q3)	35.7 (24.3, 48.6)	10.1 (5.4, 17.6)
Per capita income (US \$)		
Mean (CV %)	16,100 (42.5)	32,300 (52.1)
Median (Q1, Q3)	14,000 (11,300, 19,500)	28,600 (21,700, 38,200)
Living poverty (%)		
Mean (CV %)	29.6 (38.6)	11.5 (93.1)
Median (Q1, Q3)	28.9 (21.0, 37.1)	8.2 (3.8, 15.9)
Uninsured (%)		
Mean (CV %)	29.1 (29.1)	9.4 (75.7)
Median (Q1, Q3)	29.1 (23.3, 34.3)	7.6 (4.2, 12.6)
Crowded housing (%)		
Mean (CV %)	11.8 (52.9)	3.6 (146.1)
Median (Q1, Q3)	11.2 (7.5, 15.0)	1.9 (0.6, 4.4)
Renters (%)		
Mean (CV %)	36.9 (46.5)	36.8 (62.3)
Median (Q1, Q3)	34.5 (24.9, 44.9)	31.7 (18.6, 51.4)
Rent burden (%)		
Mean (CV %)	55.8 (26.4)	48.2 (33.4)

Variable	Cameron (n=84)	US (n=73,056)
Median (Q1, Q3)	55.0 (48.9, 65.0)	48.6 (38.0, 59.0)
No vehicle (%)		
Mean (CV %)	8.6 (81.1)	9.4 (130.4)
Median (Q1, Q3)	7.40 (3.5, 11.0)	5.3 (2.5, 11.0)

^aCV: coefficient of variation.

^bQ1, Q3: first quartile, third quartile.

Figure 3. Estimated RRs and posterior 95% credible intervals associated with census tract social determinants of health. Estimates are obtained from fitting a Bayesian spatiotemporal negative binomial on all COVID-19 confirmed cases from Cameron County, TX between March 18, 2020, and December 16, 2020. RR: relative risk.

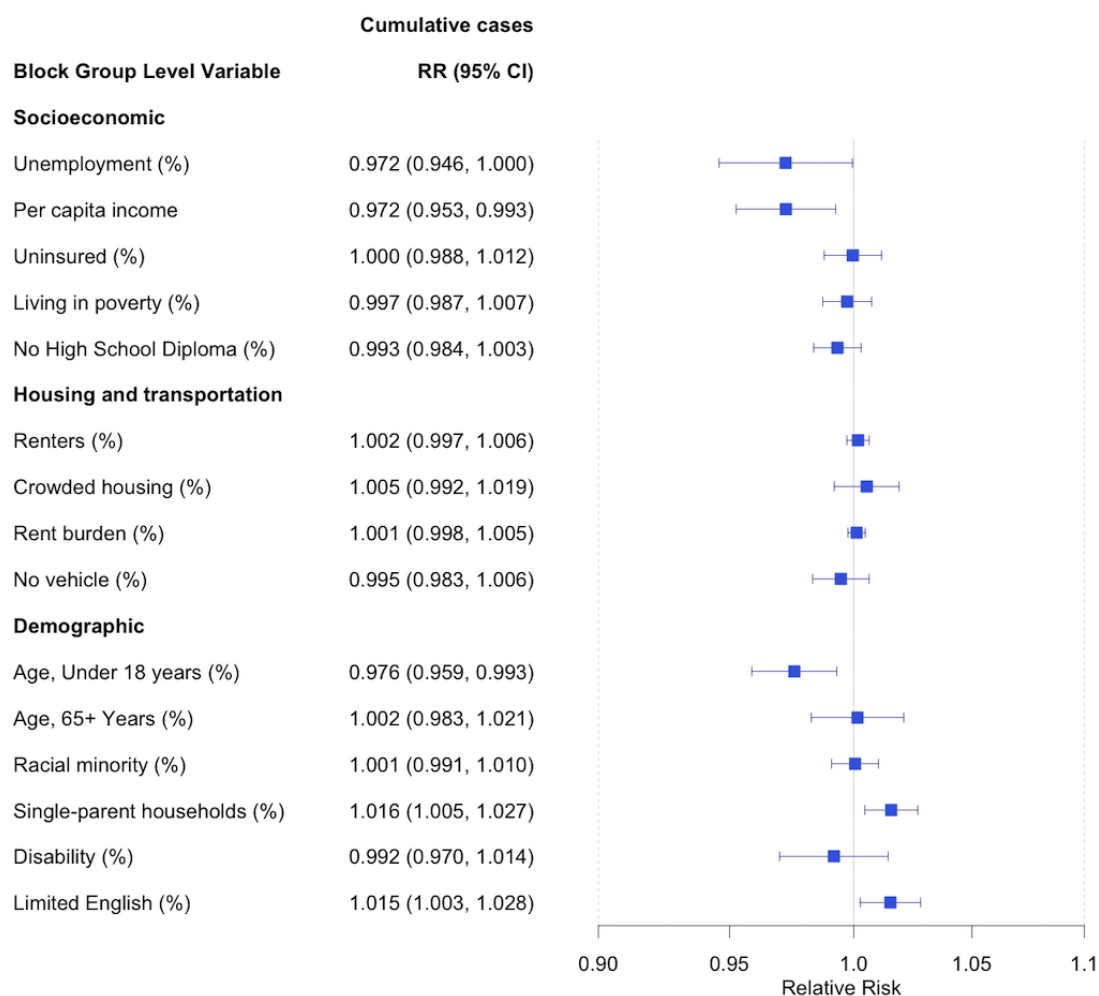
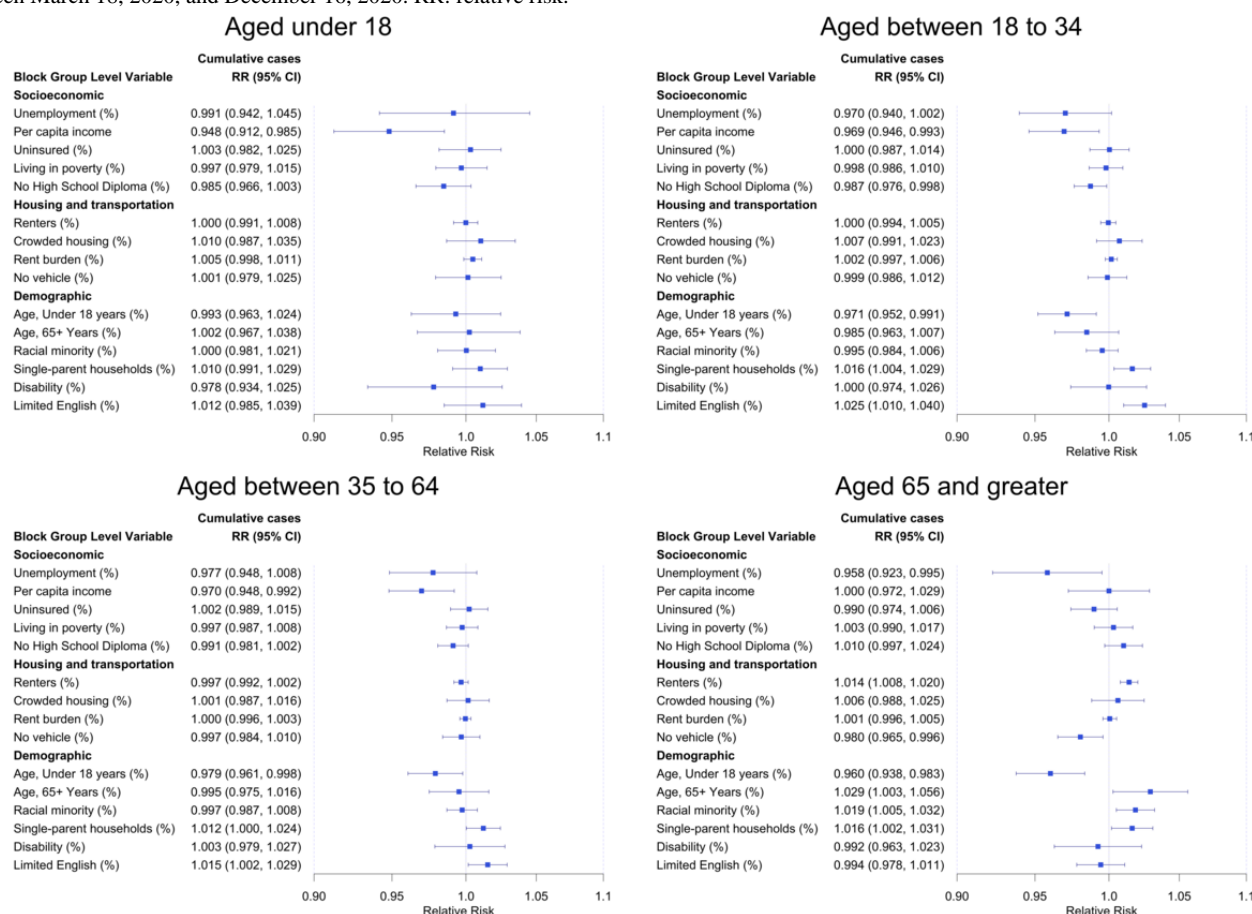


Figure 4. Estimated RRs and posterior 95% credible intervals associated with census tract social determinants of health. Estimates are obtained from fitting a Bayesian spatiotemporal negative binomial model and stratified by age groups, using COVID-19–confirmed cases from Cameron County, TX between March 18, 2020, and December 16, 2020. RR: relative risk.



Discussion

Using reported and confirmed COVID-19 cases from the Cameron County Public Health Department, we identified SDOH that were associated with COVID-19 incidence risk at the census tract level for the overall population and age subgroups. Risk of COVID-19 incidence was statistically significantly higher among areas with higher percentages of single-parent households and limited English-speaking proficiency but lower among areas with younger populations and lower income. The protective effects of lower income (for all cases) and lower education (for the age group 19-34 years) were difficult to decipher. On one hand, people living in low employment areas during the pandemic may have reduced contact with those infected and hence less likely to get infected. On the other hand, people with *essential jobs* (eg, food services) also tend to live in low income and low education areas. They may not be able to shelter at home like those in other jobs and hence have a higher risk of getting infected. For example, a previous study from Orange County, California showed an increase in COVID-19 cases in Hispanic and Latinx populations who lived in low-income census tracts and had low education attainment [23]. We were not able to further investigate the association with the census tract unemployment rate due to the lack of employment data at the census tract level during the pandemic. Other SDOH variables and social vulnerability indices such as those provided by the CDC [24] were not

included in this analysis since they are typically constructed using the ACS variables we included in this analysis or tend to be highly associated with those included. Our result on the *shelter-in-place* policy agreed with previous studies where stay-at-home orders were effective in decreasing the confirmed case growth rate [25], and cumulative COVID-19 cases fell by about 50% following 3 weeks of a *shelter-in-place* order [26], but the effects vary in magnitude by age subgroups.

Our study has some limitations. First, our analysis only included the reported and confirmed cases, and hence missed those that were unreported or undiagnosed. Second, we were not able to evaluate the individual contribution of each different mitigation plan on reducing COVID-19 incidence risk. Third, we could not include the pre-existing conditions such as diabetes and obesity prevalence in our analysis, which were shown to impact COVID-19 severity but were unclear on infection. Finally, and probably the most important one, is that we were not able to include the overall testing data due to the lack of complete and accurate testing data by census tract level in the study region. Accurately capturing the COVID-19 pandemic requires an enhanced surveillance database, where ideally testing and infection data can be linked at the individual level. We hope in our future endeavor to assist the county and city public health departments to construct a comprehensive surveillance database as such to provide real-time monitoring and early detection of future COVID-19 outbreaks.

The population we focus on in this analysis is one of the poorest in the United States, frequently uninsured, and with limited access to COVID-19 testing throughout the pandemic. Using a Bayesian spatiotemporal binomial model, we investigated the association of SDOH and COVID-19 *shelter-in-place* policies with confirmed COVID-19 cases. Though there has been a surge of studies investigating the association of SDOH and COVID-19-related health outcomes since the pandemic started, most of them focused on the county-level analysis [27-30]. This spatial unit may lack the granularity to detect local hotspots and, subsequently, is inadequate to inform the local public health officials for mitigation control and planning. To our knowledge, our study is the first conducted at a granular spatial scale of

census tracts and on a highly disadvantaged Hispanic population with limited health care access and a high chronic health risk including diabetes and obesity. The analysis also provided key information in guiding the intervention strategies to increase the testing uptake in the underserved population. For example, we are currently using this methodology as part of the Rapid Access to Diagnostics for Underserved Populations program that aims to increase knowledge about and access to testing in high-risk communities. The information generated from this study and the application of this methodology is informing both the development of targeted intervention strategies and the deployment of services to these areas.

Acknowledgments

We acknowledge the important contribution of data managers and staff in the extensive and tedious job of data entry and cleaning from Cameron County Public Health. These staff include Gabriela Saucedo, Raquel Castillo, Caludia Soto, and Saul Ruvalcaba. We also thank all the officials of the cities and counties for their tireless efforts in setting up and operating testing facilities, software for COVID-19 testing appointments, and all the other work necessary in epidemic control, in particular Art Rodriguez, Michelle Jones, and Alvaro Silva from the City of Brownsville.

This study was partially supported by National Institutes of Health funding 3UL1TR003167-02S1.

Authors' Contributions

CB conceived and designed the analysis. EG and BR contributed to the acquisition of data. IC and KZ contributed to data processing and data curation. CB and KZ conducted the data analysis. CB, ML, and BR contributed to the interpretation of the results. CB, KZ, and ML drafted the initial manuscript. SFH and JM supervised the findings of the project. CB, KZ, SFH, JM, MEF, and BR contributed to critical revision of the article. MEF contributed to the acquisition of the financial support for the project leading to this publication.

Conflicts of Interest

None declared.

Multimedia Appendix 1

Maps of demographic and social determinants of health variables at census tract level in Cameron County, TX. Data were obtained from American Community Survey 2013-2018 5-year estimates. At the tract level, the average percentage of Hispanics is 89.4% (SD 13%), much higher compared to the national average of 16.4%.

[PNG File, 2012 KB - [publichealth_v7i8e29205_app1.png](#)]

Multimedia Appendix 2

Census tract level estimated relative risks associated with social determinants of health variables, with posterior 95% credible intervals in parentheses. The relative risk (RR) estimates were obtained from fitting a Bayesian negative binomial regression model, with spatial and spatiotemporal random effects. RRs with statistically significant results are shown in bold. COVID-19 case data between March 19, 2020, and December 16, 2020, from Cameron County, TX was used in reporting the results.

[DOCX File, 16 KB - [publichealth_v7i8e29205_app2.docx](#)]

References

1. Abrams EM, Szeffler SJ. COVID-19 and the impact of social determinants of health. *Lancet Respir Med* 2020 Jul;8(7):659-661 [FREE Full text] [doi: [10.1016/S2213-2600\(20\)30234-4](#)] [Medline: [32437646](#)]
2. Golestaneh L, Neugarten J, Fisher M, Billett HH, Gil MR, Johns T, et al. The association of race and COVID-19 mortality. *EClinicalMedicine* 2020 Aug;25:100455 [FREE Full text] [doi: [10.1016/j.eclinm.2020.100455](#)] [Medline: [32838233](#)]
3. Yancy CW. COVID-19 and African Americans. *JAMA* 2020 May 19;323(19):1891-1892. [doi: [10.1001/jama.2020.6548](#)] [Medline: [32293639](#)]
4. Boserup B, McKenney M, Elkbuli A. Disproportionate impact of COVID-19 pandemic on racial and ethnic minorities. *Am Surg* 2020 Dec;86(12):1615-1622 [FREE Full text] [doi: [10.1177/0003134820973356](#)] [Medline: [33231496](#)]

5. Li D, Gaynor SM, Quick C, Chen JT, Stephenson BJK, Coull BA, et al. Unraveling US National COVID-19 racial/ethnic disparities using county level data among 328 million Americans. medRxiv. Preprint posted online on January 12, 2021. [doi: [10.1101/2020.12.02.20234989](https://doi.org/10.1101/2020.12.02.20234989)] [Medline: [33300014](#)]
6. COVID Data Tracker. Centers for Disease Control and Prevention. 2020 Mar 28. URL: <https://covid.cdc.gov/covid-data-tracker> [accessed 2021-04-07]
7. Cordes J, Castro MC. Spatial analysis of COVID-19 clusters and contextual factors in New York City. *Spat Spatiotemporal Epidemiol* 2020 Aug;34:100355 [FREE Full text] [doi: [10.1016/j.sste.2020.100355](https://doi.org/10.1016/j.sste.2020.100355)] [Medline: [32807400](#)]
8. Oluyomi AO, Gunter SM, Leining LM, Murray KO, Amos C. COVID-19 community incidence and associated neighborhood-level characteristics in Houston, Texas, USA. *Int J Environ Res Public Health* 2021 Feb 04;18(4):1495 [FREE Full text] [doi: [10.3390/ijerph18041495](https://doi.org/10.3390/ijerph18041495)] [Medline: [33557439](#)]
9. Krogstad JM. Hispanics have accounted for more than half of total U.S. population growth since 2010. Pew Research Center. 2020. URL: <https://www.pewresearch.org/fact-tank/2020/07/10/hispanics-have-accounted-for-more-than-half-of-total-u-s-population-growth-since-2010/> [accessed 2020-09-23]
10. Fisher-Hoch SP, Vatcheva KP, Rahbar MH, McCormick JB. Undiagnosed diabetes and pre-diabetes in health disparities. *PLoS One* 2015;10(7):e0133135 [FREE Full text] [doi: [10.1371/journal.pone.0133135](https://doi.org/10.1371/journal.pone.0133135)] [Medline: [26186342](#)]
11. Cole SA, Laviada-Molina HA, Serres-Perales JM, Rodriguez-Ayala E, Bastarrachea RA. The COVID-19 pandemic during the time of the diabetes pandemic: likely fraternal twins? *Pathogens* 2020 May 19;9(5):389 [FREE Full text] [doi: [10.3390/pathogens9050389](https://doi.org/10.3390/pathogens9050389)] [Medline: [32438687](#)]
12. QuickFacts Cameron County, Texas; Texas; United States. United States Census Bureau. URL: <https://www.census.gov/quickfacts/fact/table/cameroncountytexas,tx,us/PST045219> [accessed 2021-01-04]
13. Fisher-Hoch S, Rentfro AR, Salinas JJ, Pérez A, Brown HS, Reininger BM, et al. Socioeconomic status and prevalence of obesity and diabetes in a Mexican American community, Cameron County, Texas, 2004-2007. *Prev Chronic Dis* 2010 May;7(3):A53 [FREE Full text] [Medline: [20394692](#)]
14. Watt GP, Fisher-Hoch SP, Rahbar MH, McCormick JB, Lee M, Choh AC, et al. Mexican American and South Asian population-based cohorts reveal high prevalence of type 2 diabetes and crucial differences in metabolic phenotypes. *BMJ Open Diabetes Res Care* 2018;6(1):e000436 [FREE Full text] [doi: [10.1136/bmjdr-2017-000436](https://doi.org/10.1136/bmjdr-2017-000436)] [Medline: [29607048](#)]
15. Kahle D, Wickham H. ggmap: spatial visualization with ggplot2. *R J* 2013;5(1):144. [doi: [10.32614/rj-2013-014](https://doi.org/10.32614/rj-2013-014)]
16. Welcome to geocoder. United States Census Bureau. URL: <https://geocoding.geo.census.gov/> [accessed 2020-08-03]
17. Cameron County Public Health. URL: <https://www.cameroncounty.us/covid-19/> [accessed 2021-02-02]
18. Bauer C, Wakefield J, Rue H, Self S, Feng Z, Wang Y. Bayesian penalized spline models for the analysis of spatio-temporal count data. *Stat Med* 2016 May 20;35(11):1848-1865 [FREE Full text] [doi: [10.1002/sim.6785](https://doi.org/10.1002/sim.6785)] [Medline: [26530705](#)]
19. Knorr-Held L, Besag J. Modelling risk from a disease in time and space. *Stat Med* 1998 Sep 30;17(18):2045-2060. [doi: [10.1002/\(sici\)1097-0258\(19980930\)17:18<2045::aid-sim943>3.0.co;2-p](https://doi.org/10.1002/(sici)1097-0258(19980930)17:18<2045::aid-sim943>3.0.co;2-p)] [Medline: [9789913](#)]
20. Besag J, York J, Mollié A. Bayesian image restoration, with two applications in spatial statistics. *Ann Inst Stat Math* 1991 Mar;43(1):1-20. [doi: [10.1007/bf00116466](https://doi.org/10.1007/bf00116466)]
21. RStudio. 2020. URL: <http://www.rstudio.com/> [accessed 2020-09-22]
22. Lindgren F, Rue H. Bayesian spatial modelling with R-INLA. *J Stat Software* 2015;63(19):1-25. [doi: [10.18637/jss.v063.i19](https://doi.org/10.18637/jss.v063.i19)]
23. Chow DS, Soun JE, Glavis-Bloom J, Weinberg B, Chang PD, Mutasa S, et al. The Disproportionate Rise in COVID-19 Cases among Hispanic/Latinx in Disadvantaged Communities of Orange County, California: A Socioeconomic Case-Series. *MedRxiv* 2020 May 04.
24. CDC's Social Vulnerability Index (SVI). Centers for Disease Control and Prevention. 2021 Jan 19. URL: <https://www.atsdr.cdc.gov/placeandhealth/svi/index.html> [accessed 2021-04-01]
25. Li Y, Li M, Rice M, Zhang H, Sha D, Li M, et al. The Impact of Policy Measures on Human Mobility, COVID-19 Cases, and Mortality in the US: A Spatiotemporal Perspective. *Int J Environ Res Public Health* 2021 Jan 23;18(3):996 [FREE Full text] [doi: [10.3390/ijerph18030996](https://doi.org/10.3390/ijerph18030996)] [Medline: [33498647](#)]
26. Dave D, Friedson AI, Matsuzawa K, Sabia JJ. When Do Shelter-in-Place Orders Fight COVID-19 Best? Policy Heterogeneity Across States and Adoption Time. *Econ Inq* 2020 Aug 03;59(1):29-52 [FREE Full text] [doi: [10.1111/ecin.12944](https://doi.org/10.1111/ecin.12944)] [Medline: [32836519](#)]
27. Dasgupta S, Bowen VB, Leidner A, Fletcher K, Musial T, Rose C, et al. Association between social vulnerability and a county's risk for becoming a COVID-19 hotspot - United States, June 1-July 25, 2020. *MMWR Morb Mortal Wkly Rep* 2020 Oct 23;69(42):1535-1541. [doi: [10.15585/mmwr.mm6942a3](https://doi.org/10.15585/mmwr.mm6942a3)] [Medline: [33090977](#)]
28. Fielding-Miller RK, Sundaram ME, Brouwer K. Social determinants of COVID-19 mortality at the county level. *PLoS One* 2020;15(10):e0240151 [FREE Full text] [doi: [10.1371/journal.pone.0240151](https://doi.org/10.1371/journal.pone.0240151)] [Medline: [33052932](#)]
29. Paul R, Arif AA, Adeyemi O, Ghosh S, Han D. Progression of COVID-19 From Urban to Rural Areas in the United States: A Spatiotemporal Analysis of Prevalence Rates. *J Rural Health* 2020 Sep 30;36(4):591-601 [FREE Full text] [doi: [10.1111/jrh.12486](https://doi.org/10.1111/jrh.12486)] [Medline: [32602983](#)]
30. Rao JS, Zhang H, Mantero A. Contextualizing COVID-19 spread: a county level analysis, urban versus rural, and implications for preparing for the next wave. *F1000Res* 2020 May 21;9:418. [doi: [10.12688/f1000research.23903.1](https://doi.org/10.12688/f1000research.23903.1)]

Abbreviations

ACS: American Community Survey

API: application programming interface

CDC: Centers for Disease Control and Prevention

CI: credible interval

RR: relative risk

SDOH: social determinants of health

UTHealth: University of Texas School of Public Health

Edited by T Sanchez; submitted 29.03.21; peer-reviewed by V Verma, A Sheon; comments to author 18.05.21; revised version received 26.05.21; accepted 02.06.21; published 05.08.21.

Please cite as:

*Bauer C, Zhang K, Lee M, Fisher-Hoch S, Guajardo E, McCormick J, de la Cerda I, Fernandez ME, Reininger B
Census Tract Patterns and Contextual Social Determinants of Health Associated With COVID-19 in a Hispanic Population From
South Texas: A Spatiotemporal Perspective
JMIR Public Health Surveill 2021;7(8):e29205
URL: <https://publichealth.jmir.org/2021/8/e29205>
doi: [10.2196/29205](https://doi.org/10.2196/29205)
PMID: [34081608](https://pubmed.ncbi.nlm.nih.gov/34081608/)*

©Cici Bauer, Kehe Zhang, Miryoung Lee, Susan Fisher-Hoch, Esmeralda Guajardo, Joseph McCormick, Isela de la Cerda, Maria E Fernandez, Belinda Reininger. Originally published in JMIR Public Health and Surveillance (<https://publichealth.jmir.org>), 05.08.2021. This is an open-access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR Public Health and Surveillance, is properly cited. The complete bibliographic information, a link to the original publication on <https://publichealth.jmir.org>, as well as this copyright and license information must be included.

Original Paper

Natural Language Processing Insight into LGBTQ+ Youth Mental Health During the COVID-19 Pandemic: Longitudinal Content Analysis of Anxiety-Provoking Topics and Trends in Emotion in LGBTeens Microcommunity Subreddit

Hannah R Stevens¹, BA; Irena Acic¹, MA; Sofia Rhea¹, BA

Department of Communication, University of California, Davis, Davis, CA, United States

Corresponding Author:

Hannah R Stevens, BA
Department of Communication
University of California, Davis
1 Shields Ave
Davis, CA, 95616
United States
Phone: 1 530 752 0966
Email: hrstevens@ucdavis.edu

Abstract

Background: Widespread fear surrounding COVID-19, coupled with physical and social distancing orders, has caused severe adverse mental health outcomes. Little is known, however, about how the COVID-19 crisis has impacted LGBTQ+ youth, who disproportionately experienced a high rate of adverse mental health outcomes before the COVID-19 pandemic.

Objective: We aimed to address this knowledge gap by harnessing natural language processing methodologies to investigate the evolution of conversation topics in the most popular subreddit for LGBTQ+ youth.

Methods: We generated a data set of all r/LGBTeens subreddit posts (n=39,389) between January 1, 2020 and February 1, 2021 and analyzed meaningful trends in anxiety, anger, and sadness in the posts. Because the distribution of anxiety before widespread social distancing orders was meaningfully different from the distribution after ($P<.001$), we employed latent Dirichlet allocation to examine topics that provoked this shift in anxiety.

Results: We did not find any differences in LGBTQ+ youth anger and sadness before and after government-mandated social distancing; however, anxiety increased significantly ($P<.001$). Further analysis revealed a list of 10 anxiety-provoking topics discussed during the pandemic: attraction to a friend, coming out, coming out to family, discrimination, education, exploring sexuality, gender pronouns, love and relationship advice, starting a new relationship, and struggling with mental health.

Conclusions: During the COVID-19 pandemic, LGBTQ+ teens increased their reliance on anonymous discussion forums when discussing anxiety-provoking topics. LGBTQ+ teens likely perceived anonymous forums as safe spaces for discussing lifestyle stressors during COVID-19 disruptions (eg, school closures). The list of prevalent anxiety-provoking topics in LGBTQ+ teens' anonymous discussions can inform future mental health interventions in LGBTQ+ youth.

(*JMIR Public Health Surveill* 2021;7(8):e29029) doi:[10.2196/29029](https://doi.org/10.2196/29029)

KEYWORDS

COVID-19; natural language processing; LGBTQ+; mental health; anxiety; emotion; coronavirus; outbreak

Introduction

The COVID-19 pandemic has dramatically affected both physical and mental health worldwide. As of February 1, 2021, the novel coronavirus infected over 100 million people in the United States and has killed over 2.5 million people globally [1]. Widespread fear about COVID-19, coupled with physical

and social distancing orders, has caused severe adverse mental health outcomes and interpersonal relationship turmoil [2,3]. Before the pandemic, 1 in 10 US adults had symptoms of anxiety or depressive disorder. By January 2021, this figure had increased to 4 in 10 adults [4].

This sharp mental health decline may be different for LGBTQ+ youth, who disproportionately experienced a high rate of adverse

mental health outcomes before the COVID-19 pandemic due to prejudice, victimization, and unaccepting communities [5-7]. As of 2017, rates of suicidal ideation were 4 times greater in LGBTQ+ youth than those for their heterosexual, cisgender peers [8]. Amid stay-at-home orders, school closures, rollbacks of LGBTQ+ nondiscrimination protections, and the stresses of being home in potentially unsupportive environments, LGBTQ+ youth are even more vulnerable to mental health struggles during the COVID-19 crisis [5,9].

LGBTQ+ youth report that cost and parental consent are barriers to accessing mental health resources, and the inability to access confidential school counseling during COVID-19 school closures magnifies these obstacles [10,11]. Anonymous, confidential, and free online support groups are safe resources for LGBTQ+ youth during widespread school closures. LGBTQ+ youth's use of web-based platforms for support may be reinforced by the popularity of web-based platforms among younger generations as a means to create and maintain connections [5]. Decreased access to mental health resources and counseling, paired with a lack of family support, make anonymous discussion forums practical outlets for LGBTQ+ youth.

Questioning one's sexuality is a normal developmental aspect of adolescence [12]. Pubescent adolescents may experience same-sex attraction that causes them to question their sexual orientation [12]. Over time, adolescents become more certain of their sexuality and develop different sexual orientations and gender identities. As a result, individuals, who at one point disclose they are straight and cisgender, may later identify as LGBTQ+.

As individuals begin to realize their sexual orientation, they may choose to self-disclose their identity. Scholars conceptualize self-disclosure of sexual and gender identity as a dimension of the coming-out process that is closely linked to self-esteem, emotional distress, and well-being [13,14]. Yet LGBTQ+ youth may be hesitant to disclose their gender and sexual identities for fear of stigmatization, which spans across a variety of contexts, such as health care and education [15-17]. LGBTQ+ youth strategically tailor their identity disclosure to distinct social contexts to manage their stigmatized identities [18,19].

Research shows that LGBTQ+ youth resort to computer-mediated communication to explore their identities and find community [20]. For example, teens who do not disclose their sexuality to their classmates may feel comfortable revealing their sexuality to anonymous support forums for LGBTQ+ youth. Furthermore, the isolating nature of the COVID-19 pandemic makes anonymous discussion forums especially viable outlets for naturalistic studies of LGBTQ+ youth mental health. In online anonymous discussion forums, users can seek support and forge connections, which they might access otherwise, during isolation orders and school closures.

Additionally, because gender and sexual minorities are highly stigmatized, LGBTQ+ youth may not be comfortable disclosing their gender or sexual orientation to researchers as part of formal surveys and experiments [12]. Because it allows anonymity, the LGBTeens subreddit is a microcommunity that is well-suited to the investigation of an otherwise hard-to-reach population;

it is a popular microcommunity that focuses on LGBTQ+ issues and youth. While in some subreddit spaces, users' stigmatized identities might be faced with incivility [21], the norms of the LGBTQ+ microcommunity dictate that it is a safe space for LGBTQ+ youth to seek support in a space that validates their stigmatized identities.

In summary, the COVID-19 crisis has caused a concurrent mental health crisis. LGBTQ+ youth are especially vulnerable to adverse mental health outcomes, and online anonymous support forums are a uniquely accessible resource for LGBTQ+ youth to disclose their identities during the pandemic. LGBTQ+ self-disclosure is helpful for LGBTQ+ youth's mental health [13,14], yet research is needed to investigate how this vulnerable population manages their stigmatized identities while coping with the unique challenges of the pandemic and widespread social unrest [22].

However, at the time of this study, no longitudinal studies have investigated how discussions of the themes and sentiment of LGBTQ+ youth support forums unfold. We aimed to address this knowledge gap. We raised the following question: What patterns of emotions emerge from longitudinal analyses of LGBTQ+ youth conversation during the COVID-19 crisis?

Given that LGBTQ+ youth were disproportionately vulnerable to adverse mental health outcomes before the pandemic [5-7], we were interested in understanding how LGBTQ+ youth were impacted relative to the wider population of youth affected by lifestyle stressors during the COVID-19 crisis. Furthermore, research has revealed that individuals of all ages experienced interpersonal relationship turmoil during the COVID-19 crisis [3]; thus, we were interested in understanding whether this pattern was specific to teenagers or similar to the trajectory of all interpersonal relationships experiencing relationship turmoil related to the COVID-19 crisis—regardless of age. Accordingly, we posed the following question: How does the trajectory of patterns of emotions emerging from LGBTQ+ youth conversation compare to the patterns of emotions emerging from the wider population of youth as well as those emerging from any interpersonal relationships during the COVID-19 crisis?

In addition to being suitable for naturalistic investigations of emotion over time, online communities can illuminate which topics contribute to meaningful emotional trends. Knowing which topics are emotionally distressing to LGBTQ+ individuals is a requisite precursor to informing LGBTQ+ youth mental health interventions, yet at the time of this study, none had investigated themes related to LGBTQ+ youth online forums during the COVID-19 pandemic. To address this gap in the literature, we raised the following question: What conversation topics manifest from meaningful emotional trends?

Methods

Data Set

The pushshift (version 4.1) Python (version 3.9.0) package was used to extract all public posts made between January 1, 2020 and January 31, 2021 from the r/LGBTeens subreddit (n=38,389 posts). We chose this online community because of its popularity

as a community for LGBTQ+ youth and its specific focus on teens. Because we aimed to assess how users' textual expressions manifested amid global events, not in response to others' posts, comments were excluded from the data set. Although the anonymity of Reddit prevented us from accessing demographic information about the r/LGBTeens community, Reddit users live predominantly in the United States (49.3%) [23]. Notably, given the integral role of the United States in promoting LGBTQ+ rights in foreign policy, even LGBTQ+ youth outside the United States are impacted by the erosion of US LGBTQ+ advocacy [22].

To understand whether r/LGBTeens emotional patterns were specific to LGBTQ+ teenagers, we compared the trajectory of emotional tone in r/LGBTeens posts with those in 2 other subreddit microcommunities. After a review of relevant subreddits, we determined that r/Teenagers was the largest subreddit community, with $n=1,364,980$ posts, tailored toward a wide population of teens. To investigate LGBTeens post sentiment relative to widespread interpersonal relationship turmoil during the COVID-19 crisis [3], we compared the emotional tones of r/LGBTeens and r/Teenagers posts to that in the average of r/Relationships posts over time; r/Relationships was the largest subreddit dedicated to posts about interpersonal relationships ($n=193,282$).

This study only used information that could be accessed freely by the public. This study did not include any personally identifiable information. The institutional review board recognized that analysis of publicly available data does not constitute research on human participants. Thus, ethical review approval was not required for this study.

Trend Analysis

To track negative emotions over time, we analyzed aggregate post sentiment using the Linguistic Inquiry and Word Count program (LIWC) [24]. LIWC is a computerized coding tool that analyzes psychological processes (eg, positive and negative emotions) in texts by calculating the percentage of words in prevalidated lexicons relative to all words in a text. For example, we might find that 22 of 230 (9.56%) words in a post were words related to anxiety (eg, "scared" or "stressed"), and the program would assign that particular post an anxiety score of 9.56.

We focused on levels of anger, sadness, and anxiety present in posts because these psychological processes are symptoms of COVID-19-induced mental health challenges. For example, COVID-19 health threats and uncertainty may trigger feelings of anxiety [25-27]. Similarly, the COVID-19 pandemic has presented persistent stressful stimuli that some individuals may react to with anger [28-30]. Current research has demonstrated a positive correlation between the perceived threat of COVID-19 and moods of anger and hostility [31].

Likewise, the loss of loved ones, feelings of isolation, and routine disruptions associated with the rapidly changing COVID-19 pandemic may trigger feelings of sadness. For example, a recent study [32] showed a link between increased perceived loneliness and depression symptoms. Similarly, a longitudinal study [33], conducted with children in the United

Kingdom, demonstrated a substantial increase in childhood depression symptoms compared to childhood depression symptoms before lockdowns began. As anonymous support forums serve as space for individuals to build community resilience in times of crisis [34], we analyzed the overall emotional tone of posts (positive vs negative) using LIWC to understand the potential positive influences of events (eg, the US presidential election of Joe Biden, a supporter of LGBTQ+ rights) on LGBTQ+ youth mental health [35]. Furthermore, we holistically investigated the relative emotional valence of texts over time. To investigate if patterns of emotion were specific to members of the LGBTQ+ youth community, or general to the population, we compared the trajectory of LIWC emotional tone scores of r/LGBTeens posts to those of 2 other subreddit spaces.

Lifestyle Stressors

We explored differences in the trajectory of emotions displayed in posts by visualizing trends. In addition, we recorded salient events during the crisis to examine how they may have affected the changing patterns of COVID-related user responses and associated emotions.

We marked 10 major events in the course of the COVID-19 crisis. Events were selected if they considerably disrupted LGBTQ+ youth lifestyles (eg, widespread school closures) [36,37] or if they affected the rights of the wider LGBTQ+ community [38-41]. According to the Centers for Disease Control and Prevention, the first US coronavirus cases emerged on January 21, 2020 [36]. In mid-February, many school campuses were closed temporarily to assess the severity of the virus [37]. Permanent campus closure then followed on March 16, 2020 [37]. By May 7, 2020, many schools announced that they were extending web-based learning for the remainder of 2020 [37]. School closures and shifts toward web-based learning represented a turbulent time for many adolescents during the pandemic. In addition, Black Lives Matter protests in the US peaked in June 2020, which garnered support from LGBTQ+ Pride celebrations around the world [38-40]. Black Lives Matter and Pride were followed by widespread political unrest for LGBTQ+ populations, including a change by the Trump administration that eliminated nondiscrimination protections for LGBTQ+ health care on August 8, 2020, as well as the Supreme Court nomination of Amy Coney Barrett [41]. We marked both events, as well as the election of President Biden, an advocate of LGBTQ+ populations when compared to President Trump [35,41]. Finally, we marked the January 6 Capitol riots, which threatened advocates of LGBTQ+ rights as well as President Biden's inauguration [35].

Topic Analysis

We extrapolated meaningful conversation topics related to trends in anxiety, anger, sadness, as well as, overall emotional tone. We conducted 2-tailed independent sample *t* tests to determine whether the mean levels of anxiety, sadness, and anger were significantly different before and after March 21, 2020, when COVID-19 social distancing orders were enacted. To check for homogeneity of variance, we ran Levene tests for each of the 3 variables. Although results revealed the variance of sadness ($F_{1,39352}=0.11$, $P=.74$) and anger ($t_{39352}=-0.57$, $P=.57$) met the

assumption, the variance of anxiety ($F_{1,39352}=37.99$, $P<.001$) violated the assumption. We ran a 2-tailed Mann-Whitney 2-sample rank-sum test to supplement the results.

The results of the 2-tailed independent t tests suggested differences in the mean of anger ($t_{39352}=-0.57$, $P=.57$) and sadness ($t_{39352}=-0.34$, $P=.74$) were not meaningfully different before and during the widespread social distancing orders. However, results revealed the mean of anxiety was significantly higher during social distancing measures ($t_{17672.02}=-7.94$, $P<.001$). The distribution of anxiety before widespread social distancing orders was meaningfully different from the distribution after ($U=125657033$, $z=-6.20$, $P<.001$). To supplement the conclusion that anxiety levels decreased as a function of lifestyle stressors related to the COVID-19 crisis, we plotted the anxiety levels of posts with anxiety-related topics from January 1, 2015 through January 31, 2021. We found an unprecedented upsurge between January 1, 2020 and January 1, 2021 (Multimedia Appendix 1). Next, we employed latent Dirichlet allocation (LDA) to examine topics provoking this shift in anxiety.

LDA topic modeling is a bag-of-words machine learning algorithm that extrapolates meaningful topics from a large body of texts, in this case, subreddit posts [42]. Notably, we excluded posts with 0 LIWC anxiety from the LDA model. In line with the LDA procedures and guidelines outlined by prior researchers, we preprocessed the textual data and selected the optimal number of topics, and then human coders labeled those topics [43].

Text Processing

To generate a bag of words for the LDA model, we preprocessed the texts by tokenizing the text, removing stop words, lemmatizing the text, and generating a document term matrix. Tokenization separates sentences into bags of unordered words by removing all punctuation and making words lowercase. Given our relatively small sample size ($n=7882$ texts), lemmatization was necessary to reduce model noise. Lemmatizing texts removes prefixes and suffixes by transforming all words to their base lemma (eg, we converted “vaccinated” and “vaccinating” to “vaccine”). Only nouns were retained through the process of lemmatization because other parts of speech were not meaningful to our topics.

Model Selection

We set β to learn the asymmetric prior from the data [44]. Because multiple topics may have appeared in a single subreddit post, we set $\alpha=.9$ to reflect the nuanced distribution of topics

per text. To identify the optimal number (k) of topics for the model (Multimedia Appendix 2), we used the perplexity metric (the normalized log-likelihood of the model finding a previously unseen term), in which lower perplexity values suggest greater model accuracy. To keep the model parsimonious, we compared the perplexity values of models with $k=1$ to 30 topics and selected the 11 topic model because it yielded the lowest k value relative to the lowest perplexity value (-7.68).

Topic Labeling

Topics defined by the model require human labeling. LDA generates a list of the most relevant 30 terms, along with each term's β value (ie, their relative contribution to that topic) (Multimedia Appendix 3). The algorithm also identifies posts that belong primarily to a single topic. Two coders individually inspected the top 10 most relevant posts and the 30 terms with the highest β values, and then used these terms to label each topic. Human coders deemed 1 topic incoherent (Multimedia Appendix 4).

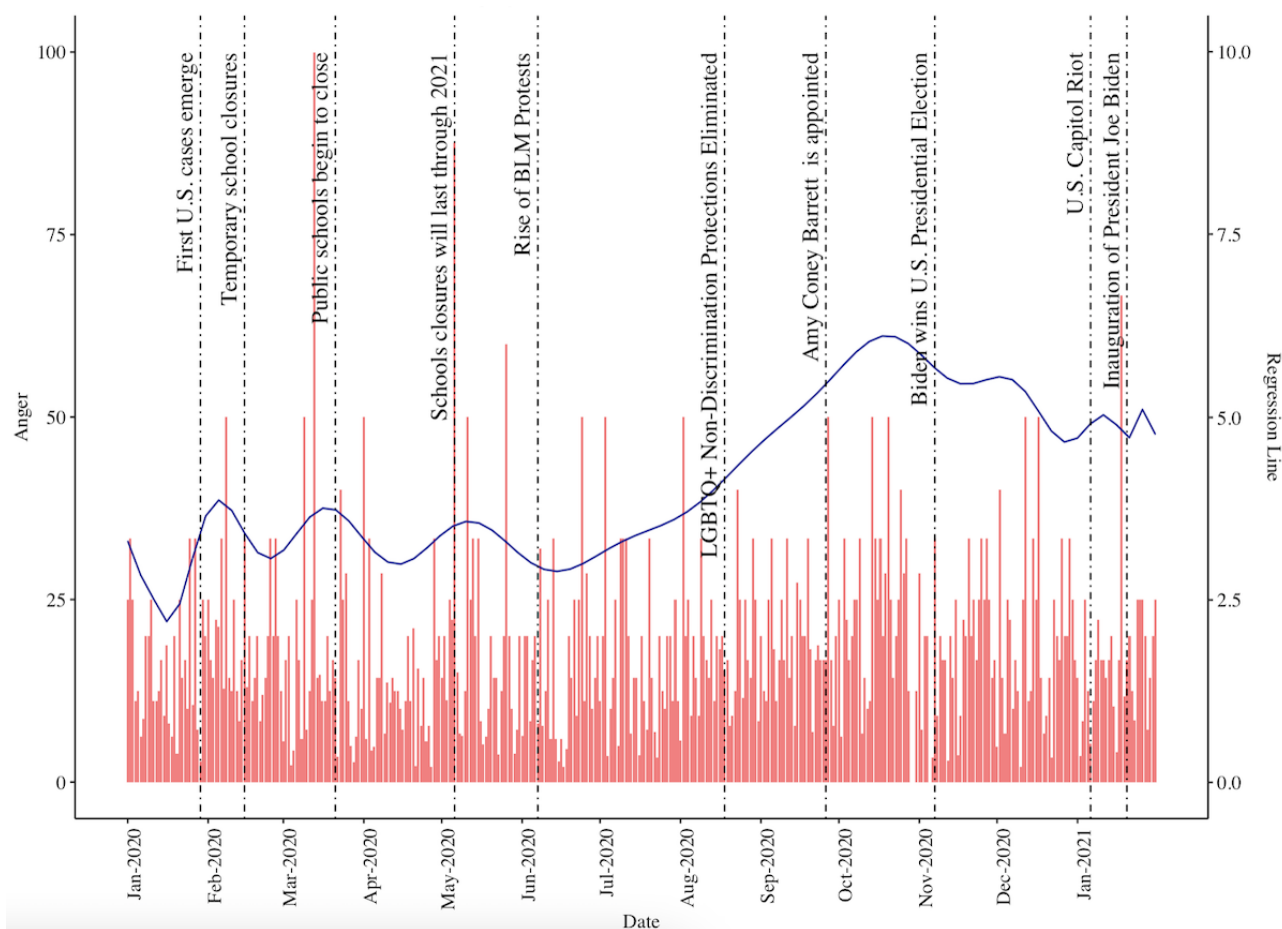
A third human coder validated the labels by reviewing the top 10 most relevant posts for each topic and confirming the human coder-assigned labels are present in those posts. The third human coder confirmed that the incoherent topic was indeed incoherent (Multimedia Appendix 5).

Results

Patterns of LGBTQ+ Teen Negative Emotion During the COVID-19 Crisis

We quantified the proportion of angry sentiment observed in r/LGBT Teens posts relative to each post's total number of words. Of the 39,389 posts in the data set, 17.67% (6961) were classified as containing anger. The mean percentage of words denoting anger relative to the total number of words in an anger-flagged post was 4.72% (SD 8.81%). There was an upsurge of anger following widespread school closures in the United States, and a second upsurge when many schools announced closures would last through the end of 2020 (Figure 1). Anger increased significantly again after the rise of the Black Lives Matter protests. More specifically, anger levels almost doubled between the Black Lives Matter protests on June 7, 2020 and when Biden won the US presidential election on November 7, 2020 [45,46]. The upsurge of anger may have been a result of frustration surrounding the Trump administration's stance on LGBTQ+ and racial injustice [47,48], which resulted in policies such as the elimination of nondiscrimination protections for LGBTQ+ health care on August 8, 2020 [41].

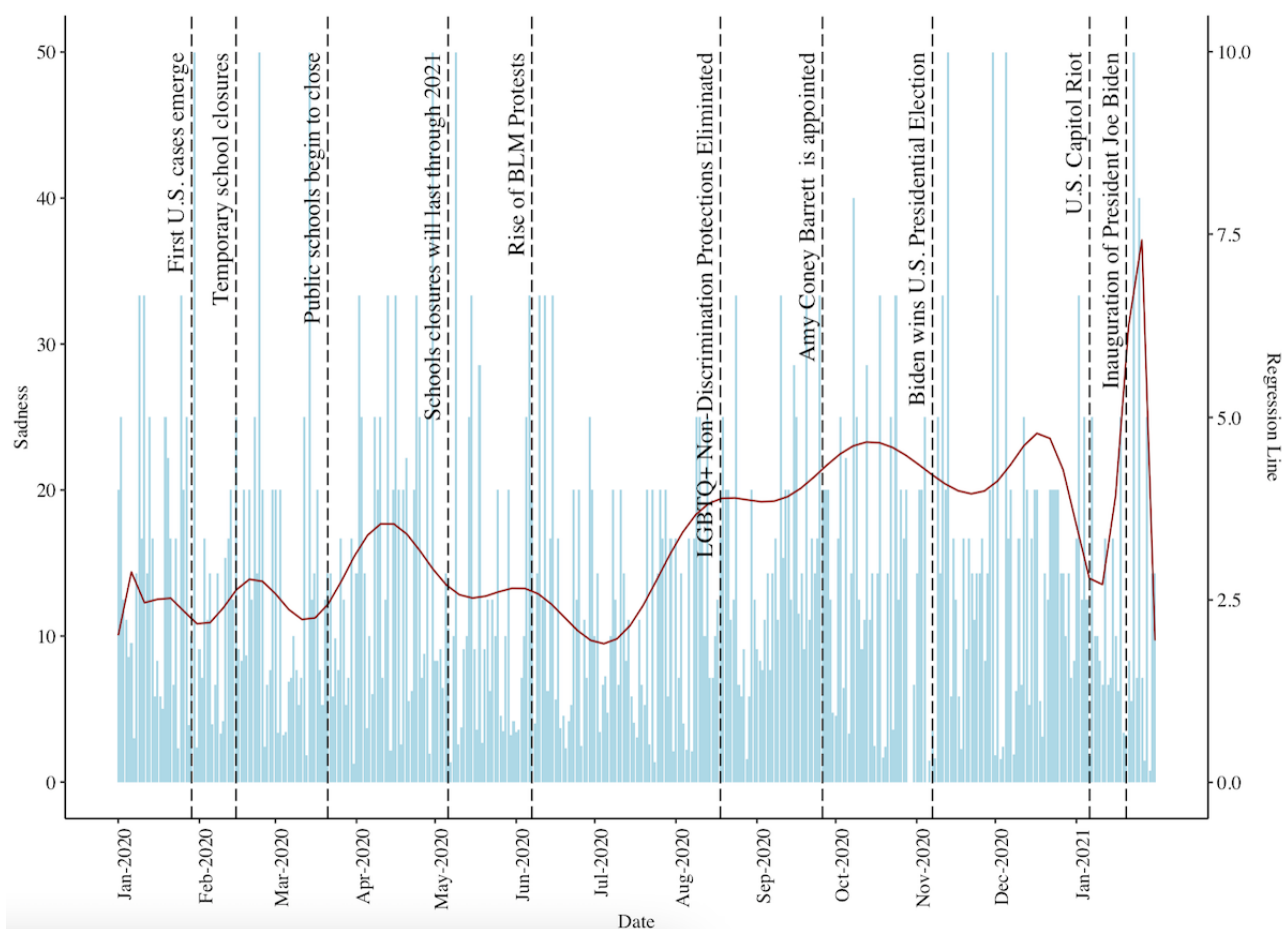
Figure 1. Histogram of the percentage of angry sentiment observed in r/LGBTeens posts over time, juxtaposed with a solid blue polynomial regression line showing average anger levels observed in angry r/LGBTeens posts over time. LGBTQ+: Lesbian, Gay, Bisexual, Transgender, Queer/Questioning, and Others.



We also quantified the proportion of sad sentiment observed in r/LGBTeens posts relative to the total number of words in each post (Figure 2). Sadness appeared in 14.4% (5687/39,389) of the posts in the data set (mean 3.82, SD 8.47). There was an upsurge of posts with high sad sentiment when the Trump administration enacted a change that eliminated nondiscrimination protections for LGBTQ+ health care on August 8, 2020 [41]. We found another increase in sadness

when Amy Coney Barrett was appointed to the US Supreme Court [41]. However, following the election of President Biden—an advocate of LGBTQ+ rights when compared to President Trump [35,41]—there was a decrease in sadness until the Capitol riots on January 6, when sadness levels sharply increased. Sadness levels decreased again after President Biden’s first week in office, perhaps because of positive feelings about President Biden’s pro-LGBTQ+ policies [35].

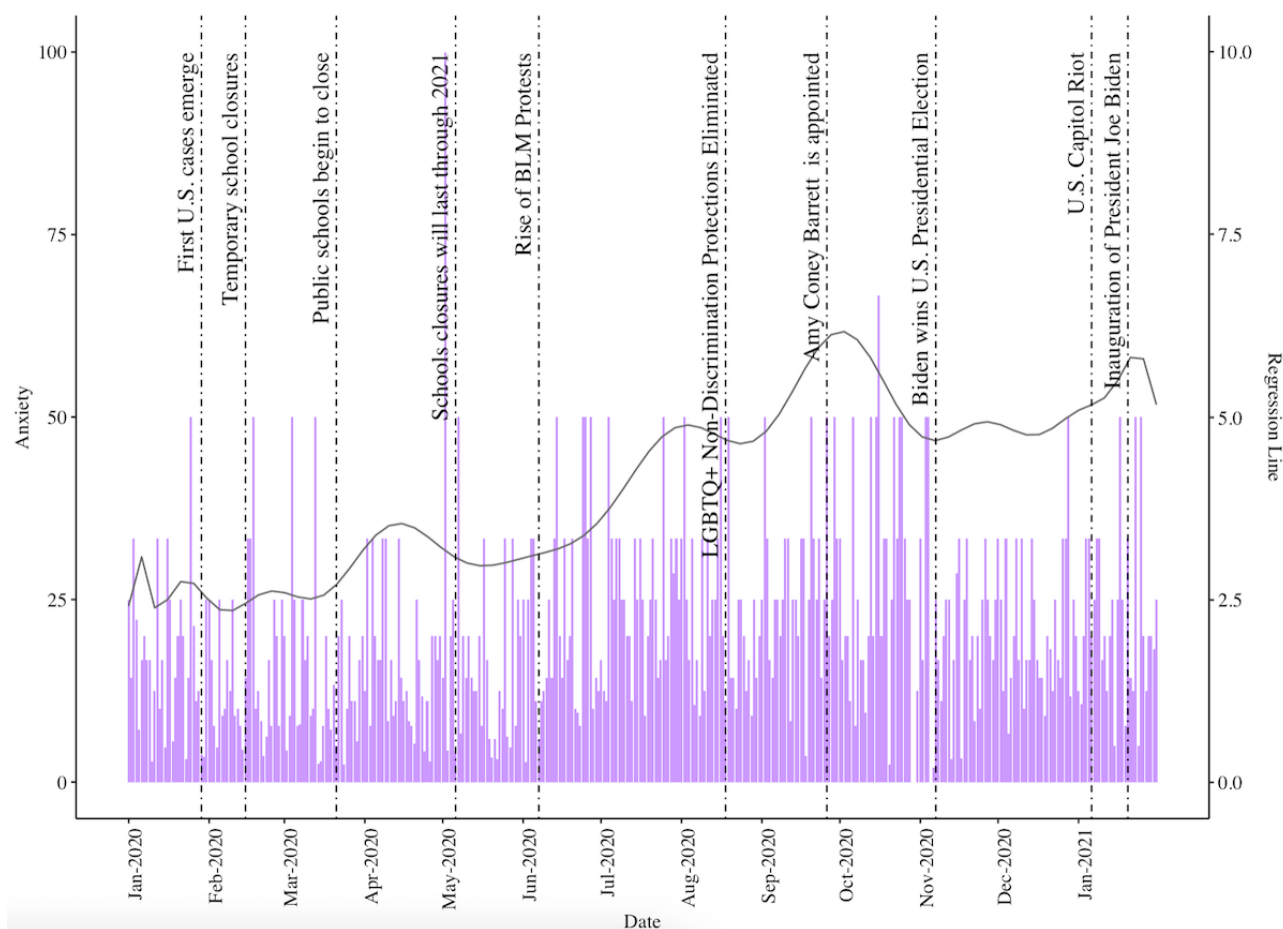
Figure 2. Histogram of the percentage of sad sentiment observed in r/LGBTeens posts over time, juxtaposed with a solid red polynomial regression line showing average sadness levels observed in sad r/LGBTeens posts over time. BLM: Black Lives Matter; LGBTQ+: Lesbian, Gay, Bisexual, Transgender, Queer/Questioning, and Others.



Of the 39,389 posts in the data set, 20.01% (7881) were classified as containing anxiety (anxiety level: mean 4.97, SD 10.43). Although all 3 negative emotions analyzed in r/LGBTeens posts increased over time, anxiety trended upward the most sharply. The histograms reveal a sharp spike in anxiety, sadness, and anger in the first week of May 2020, which may reflect emotional distress resulting from US schools closing for

the remainder of the school year (Figure 3). Another upsurge in anxiety occurred in late October, potentially a result of concern surrounding the future of LGBTQ+ constitutional rights following the Supreme Court appointment of Amy Coney Barrett [41]. This spike in anxiety flattened after President Biden won the US presidential election but spiked again on January 6, 2021, which was the day of the US Capitol riot [49].

Figure 3. Histogram of the percentage of anxiety sentiment observed in r/LGBTeens posts over time, juxtaposed with a solid grey polynomial regression line showing average anxiety levels observed in anxious r/LGBTeens posts. BLM: Black Lives Matter; LGBTQ+: Lesbian, Gay, Bisexual, Transgender, Queer/Questioning, and Others.

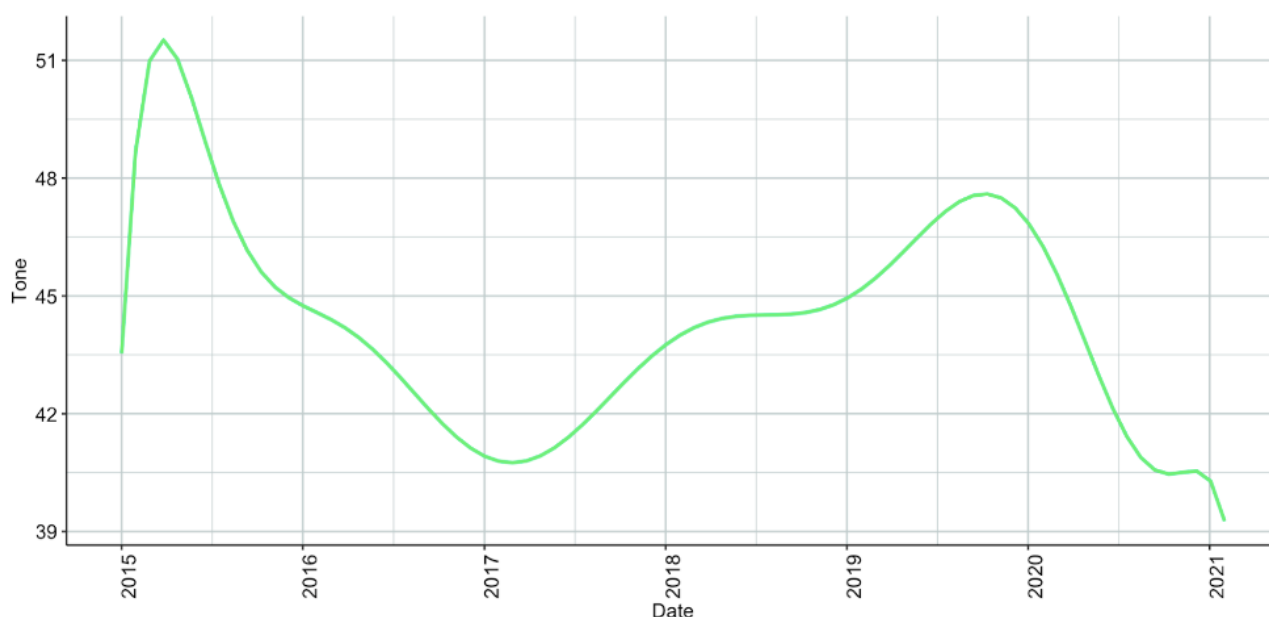


To supplement the conclusion that the mental health of r/LGBTeens community members had been adversely impacted by the COVID-19 crisis and concurrent social and political unrest, we measured the overall emotional tone of each post from the expanded timeline of January 1, 2015 through January 31, 2021. Emotional tone is measured by LIWC as the valence

of texts (ie, whether a text is positively valenced or negatively valenced) [24].

Of all posts encompassed in the 5-year period ($n=123,440$), the average post valence was negative (mean 44.33, SD 35.11, SEM 0.10, minimum 1.00, maximum 99.00, skewness -1.25). We found that posts became more negatively valenced throughout 2020 (Figure 4).

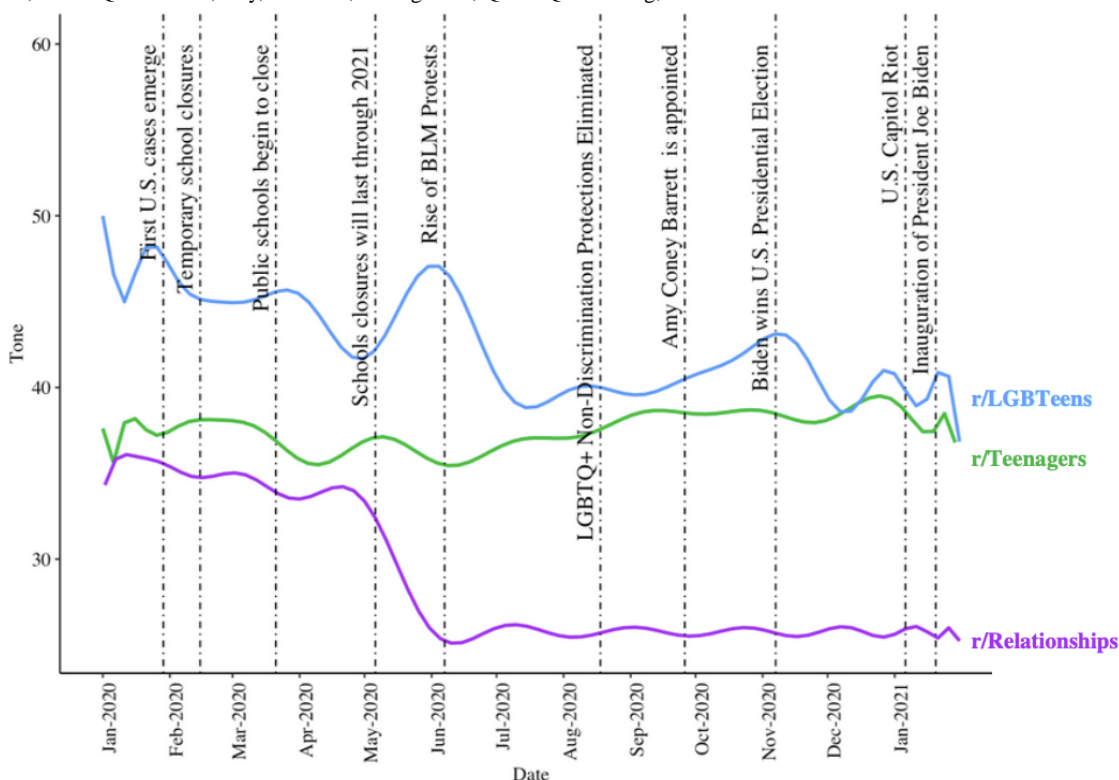
Figure 4. Green polynomial regression line representing mean emotional tone of r/LGBTeens posts from January 1, 2015 through January 31, 2021. We note a sharp decrease in the emotional tone of posts during 2020.



We compared the emotional tone of r/LGBTeens posts from January 1, 2020 to January 31, 2021 ($n=38,389$ posts) to the emotional tone of r/Teenagers ($n=1,364,980$) and r/Relationships ($n=193,282$) posts from the same time period (Figure 5). The mean emotional sentiment of r/LGBTeens posts was 42.34 (SD 34.68, SEM 0.18, minimum 1.00, maximum 99.00). Overall

post valence in the r/Teenagers subreddit ($n=1,364,980$ posts) was negative (mean 37.4, SD 32.95, SEM 0.03, minimum 1.00, maximum 99.00, skewness 0.94, kurtosis -0.52). Overall post valence in the r/Relationships subreddit ($n=193,282$) was negative (mean 28.38, SD 15.57, SEM 0.04, minimum 1.00, maximum 99.00, skewness 2.27, kurtosis 6.84).

Figure 5. Polynomial regression line representing the average post emotional tone (values of 100 represent maximally positive emotional tone; values below 50 represent more negatively valenced tone) of 3 subreddit communities (r/LGBTeens, r/Teenagers, and r/Relationships). BLM: Black Lives Matter; LGBTQ+: Lesbian, Gay, Bisexual, Transgender, Queer/Questioning, and Others.



A 2-tailed independent samples t test was conducted to examine whether the difference in mean emotional sentiment between the r/LGBTeens and r/Relationships posts were statistically

significant. Prior to the main analyses, the assumption of homogeneity of variance was checked; the results of a Levene test for r/Relationships and r/LGBTeens emotional tone was

significant based on $\alpha=.05$ ($F_{1,1403967}=37821.73$, $P<.001$), indicating the assumption was violated. The result of the 2-tailed independent samples t -test suggested the mean emotional sentiment of $r/LGBTeens$ and $r/Relationships$ posts were significantly different based on $\alpha=.05$ ($t_{42204.64}=77.99$, $P<.001$). Since the assumption of assumption of homogeneity of variance was violated, a Mann-Whitney U Test, which does not have any distributional assumptions, was conducted to supplement the t test results. The result of the 2-tailed Mann-Whitney U test examining differences between $r/LGBTeens$ and $r/Relationships$ posts was significant based on $\alpha=.05$ ($U=4301049386.5$, $z=-\text{Inf}$, $P<.001$), suggesting that the distribution of emotional tone in $r/LGBTeens$ posts was significantly different than that in $r/Relationships$.

Additionally, a 2-tailed independent samples t test was conducted to compare the mean emotional sentiment of $r/LGBTeens$ and $r/Teenagers$ posts. The results of the Levene test comparing $r/teenagers$ and $r/LGBTeens$ posts also indicated the variances of tone for $r/LGBTeens$ and $r/Relationships$ posts were unlikely to be equal based on $\alpha=.05$ ($F_{1,232265}=37821.73$, $P<.001$). The results of the 2-tailed independent sample t test suggested the mean emotional sentiment of $r/LGBTeens$ and $r/teenagers$ posts were also significantly different, based on $\alpha=.05$ ($t_{41019.96}=27.86$, $P<.001$). The homogeneity of variance violation necessitated a Mann-Whitney U Test to supplement the t test results; results were significant based on $\alpha=.05$ ($U=28602717011.5$, $z=-27.23$, $P<.001$), suggesting their distributions were significantly different.

Findings revealed that the emotional sentiment of $r/LGBTeens$ posts (mean 42.34, SD 34.68) was significantly greater than those of $r/Teenagers$ posts (mean 37.4, SD 32.95) and $r/Relationships$ posts (mean 28.38, SD 15.57).

We found a decrease in emotional tone in late January in both $r/Teenagers$ and $r/LGBTeens$ posts when many schools closed temporarily (Figure 5). This downturn was followed by a second decrease in $r/LGBTeens$ and $r/Teenagers$ tone when public schools announced that closures would last through 2020. Then, there was a sharp rise in emotional tone during the Black Lives Matter protests and simultaneous Pride celebrations for $r/LGBTeens$ posts but not for $r/Teenagers$ posts. After Pride celebrations and Black Lives Matter protests in June, the trend for $r/Teenagers$ posts seemed to stabilize. However, we saw a

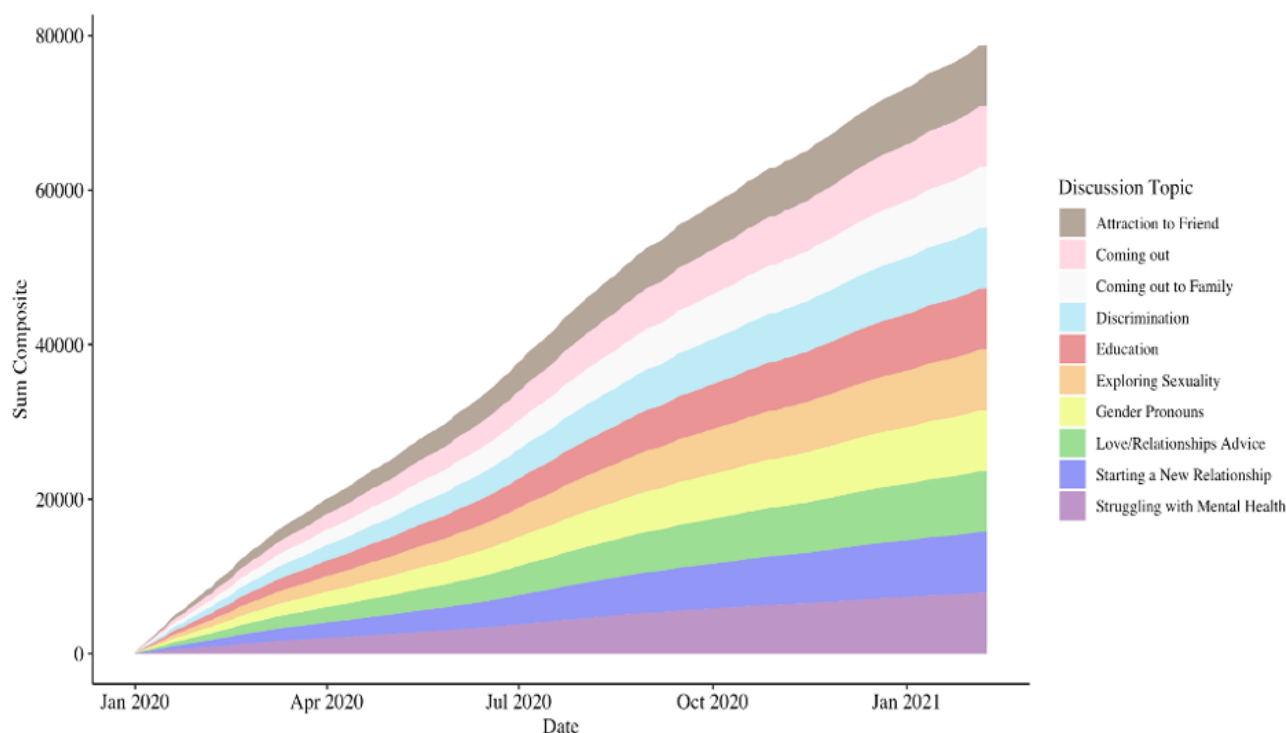
tone decrease in $r/LGBTeens$ posts, representing more negatively valenced posts, which persisted until a second prominent increase in early November 2020. This November upsurge may be explained by the election of President Biden, who has a pro-LGBTQ+ agenda in contrast to his predecessor [35,37,47].

Furthermore, we ran a point biserial correlation analysis to assess whether negative emotion changed as US social distancing orders were relaxed and vaccine distribution increased in January 1, 2021 [50]. LIWC positive emotion calculates the frequency of positive emotion lexicon words (eg, “happy” or “thankful”) relative to the total number of words in a post. To supplement our results, we likewise assessed whether positive emotion changed as US social distancing orders were gradually lifted. While there were no significant differences in anger and sadness, results showed a significant increase in anxiety after social distance orders were reduced in January of 2021 ($r=0.01$, $P=.04$, 95% CI 0 to 0.02). Similarly, we found a decrease in positive emotion after social distancing orders began to be lifted ($r=-0.01$, $P=.04$, 95% CI -0.02 to 0). This decrease in positive emotion and increase in anxiety may be because—regardless of the vaccine becoming available on December 10, 2020 and reduced social distancing orders—youth were not immediately impacted (eg, the vaccine was not approved for youth ages 12 to 16 years until May of 2021 [51]). We would expect that this trend might flatten out in Fall 2021—when students begin to return to normalcy (Multimedia Appendix 6).

Conversation Patterns That Manifest From Heightened Anxiety

Because the distribution of anxiety before widespread social distancing orders was meaningfully different from the distribution after lockdown ($U=125657033$, $z=-6.20$, $P<.001$), we used a topic model to extrapolate the specific conversation topics co-occurring with LGBTQ+ teen anxiety. Findings revealed the following conversation topics among the LGBTQ+ youth during the pandemic accompany sharp increases in anxiety: attraction to a friend, coming out, coming out to family, discrimination, education, exploring sexuality, gender pronouns, love/relationship advice, starting a new relationship, and struggling with mental health. There was no meaningful change in the anxiety-provoking topics discussed over time (Multimedia Appendix 7); however, the frequency of such discussions increased (Figure 6).

Figure 6. Streamgraph of topic percent contribution to the corpus of anxious r/LGBT teens posts over time. Each color represents a discussion topic across the corpus of anxious r/LGBT teens posts. This streamgraph shows the composite of topics overall (graph envelope) and the relative importance of the topic to posts over time (topic color stream width).



Discussion

We employed natural language processing to investigate the emotional trends in LGBTQ+ teens' anonymous online conversations during the COVID-19 pandemic. Results revealed that the overall emotional tone of posts sharply decreased during the 2020-2021 COVID-19 crisis, relative to prior years—revealing this emotional trend was specific to the COVID-19 crisis. Findings reveal that the emotional trajectory of LGBTQ+ youth fluctuated more drastically in response to impactful events during the COVID-19 crisis (eg, widespread school closures and Black Lives Matter protests) compared to the emotional trajectory of more neutral subreddit spaces [37,39,40].

Findings revealed that the trajectory of LGBTQ+ teens' overall emotional tone (positive vs negative) to be more affected by lifestyle stressors during the COVID-19 crisis than the general population of r/Teenage users. Results are consistent with those from previous research indicating that LGBTQ+ youth are disproportionately vulnerable to adverse mental health outcomes relative to their straight, cisgender peers [4-6].

While this study did not find pre and postlockdown differences in LGBTQ+ youth anger and sadness, results revealed that anxiety increased after government-mandated social distancing measures. In addition, further analysis revealed a list of 10 anxiety-provoking topics discussed during the pandemic: attraction to a friend, coming out, coming out to family, discrimination, education, exploring sexuality, gender pronouns, love/relationship advice, starting a new relationship, and struggling with mental health. These conversation topics were anxiety-provoking for LGBTQ+ youth both before and during

the pandemic. However, the increase in the frequency of these conversations coincided with the emergence of lifestyle disruptors related to the pandemic, reflecting LGBTQ+ teens' increased reliance on anonymous discussion forums as outlets for discussing lifestyle stressors during COVID-19 lifestyle disruptions (eg, school closures).

Findings revealing LGBTQ+ teens' increased reliance on an anonymous forum as a discussion outlet during the COVID-19 crisis were consistent with those from previous studies showing that individuals are likely to turn to social media in times of crisis to seek psychological support and build community resilience [34]. Furthermore, the results of this study are in line with those of existing studies demonstrating the importance of online support to LGBTQ+ youth while coping with the challenges of a global pandemic [5].

This study also shed light on the specific sources of anxiety for LGBTQ+ youth during the COVID-19 pandemic. Research has revealed links between LGBTQ+ youth anxiety disorders and self-harm and suicidal behavior—in part due to stigma and discrimination [6]. Additionally, research has shown that elevated anxiety levels can weaken individuals' immune systems, make them more vulnerable to certain illnesses, and elevate their risk of death from cardiovascular complications [52]. Clinical practices show that identifying sources of anxiety is an essential step in helping teenagers develop coping mechanisms [53]. By identifying sources of anxiety in LGBTQ+ youth, this study may help mental health professionals design effective strategies to address anxiety in that population.

Additionally, this study's findings suggest that mental health professionals should consider anonymous online supplements or alternatives to in-person treatment of LGBTQ+ youth anxiety,

especially during school closures. Despite mental health professionals' adaptation to web-based counseling, LGBTQ+ youth report that treatment cost and parental consent are barriers to accessing mental health resources outside of school [10,11]. Findings accentuate the need to develop web-based programming to address the needs of LGBTQ+ youth in times of crisis and were consistent with those of Pacey et al [54]. Research shows that youth whose social environments support their LGBTQ+ identities are less likely to have suicidal thoughts or attempt suicide, compared to those who live in nonsupportive social environments [55]. Furthermore, being connected to the LGBTQ+ community may provide a buffer to suicidal thoughts [56]. Amid widespread social distancing orders, LGBTQ+ youth may have turned to web-based resources as a way of connecting with the LGBTQ+ community [57]. For instance, there was an upsurge of r/LGBTeens positive emotion during the rise of the Black Lives Matter protests, which may be explained by concurrent Gay Pride celebration unity, where many Pride celebrations stood in solidarity with Black Lives Matter to support marginalized populations [38-40]. The upsurge of positive emotion during Gay Pride was followed by another upsurge when President Biden won the US presidential election (Multimedia Appendix 8). Future work should consider how LGBTQ+ youth well-being may benefit from identity-affirming online spaces to celebrate their individual and collective achievements [34].

Although this study provides valuable insight into LGBTQ+ youth mental health during the COVID-19 pandemic, the study had some limitations. First, using computerized coding tools such as LIWC does not allow for sophisticated coding that could

be achieved with human coders. Previous studies using LIWC have found that LIWC may overidentify emotional expression [58]; thus, LIWC may have captured extraneous sentiment. Second, research shows that observational studies may lead to faulty findings due to confounding factors. Although scholars admit that it is not possible to identify all confounders in practice [59], we expect that online expressions will vary between individuals according to their level of exposure to mental health risk factors, such as institutionalized and interpersonal discrimination and the lack of parental support [6,52]. Although these factors are consequential, we cannot test them using an observational study. To reduce bias and limitations through methodological triangulation, future research should marry observational work on LGBTQ+ youth with surveys and experimental data. Surveying LGBTQ+ youth may also help to identify the specific sources of emotional distress (eg, negative emotion resulting from the elimination of LGBTQ+ nondiscrimination protections vs negative emotion related to COVID-19 isolation).

The COVID-19 crisis has caused a concurrent mental health pandemic, and LGBTQ+ youth are especially vulnerable to adverse mental health outcomes [4-7]. Results reveal that online microcommunities serve as a viable space for LGBTQ+ youth to express their emotions about lifestyle stressors and navigate their stigmatized identities while simultaneously building community resilience [5,34]. Because LGBTQ+ self-disclosure is helpful for mental health [13,14], future work should explore the potential of web-based platforms in identifying the sources of LGBTQ+ youth emotional distress and developing safe online spaces to offer support to this vulnerable population.

Conflicts of Interest

None declared.

Multimedia Appendix 1

Polynomial regression of Linguistic Inquiry Word Count anxiety levels in r/LGBTeens posts from January 1, 2015 through January 31, 2021.

[DOCX File, 104 KB - [publichealth_v7i8e29029_app1.docx](#)]

Multimedia Appendix 2

Perplexity metric.

[DOCX File, 52 KB - [publichealth_v7i8e29029_app2.docx](#)]

Multimedia Appendix 3

Latent Dirichlet allocation results.

[DOCX File, 14 KB - [publichealth_v7i8e29029_app3.docx](#)]

Multimedia Appendix 4

Intertopic distance.

[DOCX File, 123 KB - [publichealth_v7i8e29029_app4.docx](#)]

Multimedia Appendix 5

Topic model data processing.

[DOCX File, 97 KB - [publichealth_v7i8e29029_app5.docx](#)]

Multimedia Appendix 6

Point biserial correlation analysis.

[\[DOCX File, 14 KB - publichealth_v7i8e29029_app6.docx\]](#)

Multimedia Appendix 7

Anxiety-related topics over time.

[\[DOCX File, 209 KB - publichealth_v7i8e29029_app7.docx\]](#)

Multimedia Appendix 8

Histogram of the percentage of positive sentiment in r/LGBTeens posts over time and average level of positive emotion in positive r/LGBTeens posts over time.

[\[PNG File, 417 KB - publichealth_v7i8e29029_app8.png\]](#)

References

1. Dong E, Du H, Gardner L. An interactive web-based dashboard to track COVID-19 in real time. *Lancet Infect Dis* 2020 May;20(5):533-534 [FREE Full text] [doi: [10.1016/S1473-3099\(20\)30120-1](#)] [Medline: [32087114](#)]
2. McElroy E, Patalay P, Moltrecht B, Shevlin M, Shum A, Creswell C, et al. Demographic and health factors associated with pandemic anxiety in the context of COVID-19. *Br J Health Psychol* 2020 Nov;25(4):934-944. [doi: [10.1111/bjhp.12470](#)] [Medline: [32860334](#)]
3. Goodwin R, Hou WK, Sun S, Ben-Ezra M. Quarantine, distress and interpersonal relationships during COVID-19. *Gen Psychiatr* 2020;33(6):e100385 [FREE Full text] [doi: [10.1136/gpsych-2020-100385](#)] [Medline: [33163857](#)]
4. Adults reporting symptoms of anxiety or depressive disorder during COVID-19 pandemic. Kaiser Family Foundation. 2021. URL: <https://www.kff.org/other/state-indicator/adults-reporting-symptoms-of-anxiety-or-depressive-disorder-during-covid-19-pandemic/?currentTimeframe=0&sortModel=%7B%22colId%22%3A%22Location%22%2C%22sort%22%3A%22asc%22%7D> [accessed 2021-07-09]
5. Fish JN, McInroy LB, Pacey MS, Williams ND, Henderson S, Levine DS, et al. "I'm kinda stuck at home with unsupportive parents right now": LGBTQ youths' experiences with COVID-19 and the importance of online support. *J Adolesc Health* 2020 Sep;67(3):450-452 [FREE Full text] [doi: [10.1016/j.jadohealth.2020.06.002](#)] [Medline: [32591304](#)]
6. Russell ST, Fish JN. Mental health in lesbian, gay, bisexual, and transgender (LGBT) youth. *Annu Rev Clin Psychol* 2016;12:465-487 [FREE Full text] [doi: [10.1146/annurev-clinpsy-021815-093153](#)] [Medline: [26772206](#)]
7. Shearer A, Herres J, Kodish T, Squitieri H, James K, Russon J, et al. Differences in mental health symptoms across lesbian, gay, bisexual, and questioning youth in primary care settings. *J Adolesc Health* 2016 Jul;59(1):38-43. [doi: [10.1016/j.jadohealth.2016.02.005](#)] [Medline: [27053400](#)]
8. Know the warning signs. National Alliance on Mental Illness. 2021. URL: <https://www.nami.org/Your-Journey/Identity-and-Cultural-Dimensions/LGBTQI> [accessed 2021-06-03]
9. Dole T. How LGBTQ youth can cope with anxiety and stress during COVID-19. The Trevor Project. 2020. URL: <https://www.thetrevorproject.org/2020/03/26/how-lgbtq-youth-can-cope-with-anxiety-and-stress-during-covid-19/> [accessed 2021-06-03]
10. Isselbacher J. LGBTQ youth say cost, parent consent pose barriers to mental health care. *Stat News*. 2020. URL: <https://www.statnews.com/2020/08/18/lgbtq-youth-mental-health-care/> [accessed 2020-09-07]
11. Valencia M. The challenges of the pandemic for queer youth. *The New York Times*. 2020 Jun 29. URL: <https://www.nytimes.com/2020/06/29/well/family/LGBTQ-youth-teenagers-pandemic-coronavirus.html> [accessed 2021-08-01]
12. Meyer IH, Wilson PA. Sampling lesbian, gay, and bisexual populations. *J Couns Psychol* 2009;56(1):23-31. [doi: [10.1037/a0014587](#)]
13. Kosciw JG, Palmer NA, Kull RM. Reflecting resiliency: openness about sexual orientation and/or gender identity and its relationship to well-being and educational outcomes for LGBT students. *Am J Community Psychol* 2015 Mar;55(1-2):167-178. [doi: [10.1007/s10464-014-9642-6](#)] [Medline: [24691967](#)]
14. Rosario M, Hunter J, Maguen S, Gwadz M, Smith R. The coming-out process and its adaptational and health-related associations among gay, lesbian, and bisexual youths: stipulation and exploration of a model. *Am J Community Psychol* 2001 Feb;29(1):133-160. [doi: [10.1023/A:1005205630978](#)] [Medline: [11439825](#)]
15. Payne EC, Smith MJ. Refusing relevance: school administrator resistance to offering professional development addressing LGBTQ issues in schools. *Educ Adm Q* 2017 Aug 03;54(2):183-215. [doi: [10.1177/0013161x17723426](#)]
16. Wickens CM, Sandlin JA. Homophobia and heterosexism in a college of education: a culture of fear, a culture of silence. *Int J Qual Stud Educ* 2010 Nov;23(6):651-670. [doi: [10.1080/09518390903551035](#)]
17. Whitehead J, Shaver J, Stephenson R. Outness, stigma, and primary health care utilization among rural LGBT populations. *PLoS One* 2016;11(1):e0146139 [FREE Full text] [doi: [10.1371/journal.pone.0146139](#)] [Medline: [26731405](#)]
18. Orne J. 'You will always have to "out" yourself': reconsidering coming out through strategic outness. *Sexualities* 2012 Mar 05;14(6):681-703. [doi: [10.1177/1363460711420462](#)]

19. Schmitz RM, Tyler KA. Contextual constraints and choices: strategic identity management among LGBTQ youth. *J LGBT Youth* 2018 May 31;15(3):212-226. [doi: [10.1080/19361653.2018.1466754](https://doi.org/10.1080/19361653.2018.1466754)]
20. McKenna KYA, Bargh JA. Coming out in the age of the internet: identity "demarginalization" through virtual group participation. *J Pers Soc Psychol* 1998;75(3):681-694. [doi: [10.1037/0022-3514.75.3.681](https://doi.org/10.1037/0022-3514.75.3.681)]
21. Stevens HR, Acic I, Taylor LD. Uncivil reactions to sexual assault online: linguistic features of news reports predict discourse incivility. *Cyberpsychol Behav Soc Netw* 2021 Oct 16:e26876 (forthcoming).
22. Angelo P, Bocci D. The changing landscape of global LGBTQ+ rights Internet. Council on Foreign Relations. 2021 Jan 29. URL: <https://www.cfr.org/article/changing-landscape-global-lgbtq-rights> [accessed 2021-06-06]
23. Clement J. Statista. Regional distribution of desktop traffic to Reddit.com as of June 2021, by country. 2020 Jul 29. URL: <https://www.statista.com/statistics/325144/reddit-global-active-user-distribution/> [accessed 2021-05-01]
24. Pennebaker J, Boyd R, Jordan K, Blackburn K. The development and psychometric properties of LIWC2015. University of Texas at Austin. 2015. URL: https://repositories.lib.utexas.edu/bitstream/handle/2152/31333/LIWC2015_LanguageManual.pdf [accessed 2021-06-02]
25. Gallagher MW, Zvolensky MJ, Long LJ, Rogers AH, Garey L. The impact of COVID-19 experiences and associated stress on anxiety, depression, and functional impairment in American adults. *Cognit Ther Res* 2020 Aug 29:1-9 [FREE Full text] [doi: [10.1007/s10608-020-10143-y](https://doi.org/10.1007/s10608-020-10143-y)] [Medline: [32904454](https://pubmed.ncbi.nlm.nih.gov/32904454/)]
26. Grupe DW, Nitschke JB. Uncertainty and anticipation in anxiety: an integrated neurobiological and psychological perspective. *Nat Rev Neurosci* 2013 Jul;14(7):488-501 [FREE Full text] [doi: [10.1038/nrn3524](https://doi.org/10.1038/nrn3524)] [Medline: [23783199](https://pubmed.ncbi.nlm.nih.gov/23783199/)]
27. Stevens HR, Oh Y, Taylor LD. Desensitization to fear-inducing COVID-19 health news on Twitter: observational study. *JMIR Infodemiology* 2021 Jul 16;1(1):e26876. [doi: [10.2196/26876](https://doi.org/10.2196/26876)]
28. Westbrook T. Why is COVID-19 making me so angry? Ohio State University. 2020 Apr 19. URL: <https://wexnermedical.osu.edu/blog/why-so-angry-covid> [accessed 2021-06-09]
29. Lobbestael J, Arntz A, Wiers RW. How to push someone's buttons: a comparison of four anger-induction methods. *Cogn Emot* 2008 Feb;22(2):353-373. [doi: [10.1080/02699930701438285](https://doi.org/10.1080/02699930701438285)]
30. Clore GL, Centerbar DB. Analyzing anger: how to make people mad. *Emotion* 2004 Jun;4(2):139-44; discussion 151. [doi: [10.1037/1528-3542.4.2.139](https://doi.org/10.1037/1528-3542.4.2.139)] [Medline: [15222850](https://pubmed.ncbi.nlm.nih.gov/15222850/)]
31. Wang C, Pan R, Wan X, Tan Y, Xu L, Ho CS, et al. Immediate psychological responses and associated factors during the initial stage of the 2019 coronavirus disease (COVID-19) epidemic among the general population in China. *Int J Environ Res Public Health* 2020 Mar 06;17(5):1705-1729. [doi: [10.3390/ijerph17051729](https://doi.org/10.3390/ijerph17051729)] [Medline: [32155789](https://pubmed.ncbi.nlm.nih.gov/32155789/)]
32. Palgi Y, Shrira A, Ring L, Bodner E, Avidor S, Bergman Y, et al. The loneliness pandemic: loneliness and other concomitants of depression, anxiety and their comorbidity during the COVID-19 outbreak. *J Affect Disord* 2020 Oct 01;275:109-111 [FREE Full text] [doi: [10.1016/j.jad.2020.06.036](https://doi.org/10.1016/j.jad.2020.06.036)] [Medline: [32658811](https://pubmed.ncbi.nlm.nih.gov/32658811/)]
33. Bignardi G, Dalmaijer ES, Anwyl-Irvine AL, Smith TA, Siugzdaite R, Uh S, et al. Longitudinal increases in childhood depression symptoms during the COVID-19 lockdown. *Arch Dis Child* 2020 Dec 09;320-372 [FREE Full text] [doi: [10.1136/archdischild-2020-320372](https://doi.org/10.1136/archdischild-2020-320372)] [Medline: [33298552](https://pubmed.ncbi.nlm.nih.gov/33298552/)]
34. Taylor M, Wells G, Howell G, Raphael B. The role of social media as psychological first aid as a support to community resilience building. *Australian J Emerg Manag* 2012 Feb 02;27(1):20-27.
35. President Biden's pro-LGBT timeline. Human Rights Campaign. 2020 Mar 03. URL: <https://www.hrc.org/resources/president-bidens-pro-lgbtq-timeline> [accessed 2021-07-08]
36. First travel-related case of 2019 novel coronavirus detected in United States. Center for Disease Control and Prevention. 2020 Jan 21. URL: <https://www.cdc.gov/media/releases/2020/p0121-novel-coronavirus-travel-case.html> [accessed 2021-03-23]
37. The coronavirus spring: the historic closing of U.S. schools (a timeline). Education Week. 2020 Jul 01. URL: <https://www.edweek.org/leadership/the-coronavirus-spring-the-historic-closing-of-u-s-schools-a-timeline/2020/07> [accessed 2021-03-21]
38. Salam M. Your 2020 virtual pride guide. New York Times. 2020 Jun 24. URL: <https://www.nytimes.com/article/gay-pride-2020-events-online.html> [accessed 2021-03-21]
39. LaCrosse M. Boston Pride shows solidarity with Black Lives Matter by shifting focus this June. CBS Boston. 2020 Jun 12. URL: <https://boston.cbslocal.com/2020/06/12/boston-pride-solidarity-black-lives-matter-unity-flag/> [accessed 2021-03-21]
40. Rad A. Why Pride 2020 is all about Black Lives Matter. Advocate. 2020 Jun 23. URL: <https://www.advocate.com/commentary/2020/6/23/why-pride-2020-all-about-black-lives-matter> [accessed 2021-03-21]
41. Trump administration change will eliminate nondiscrimination protections for LGBTs in health care. Out in Jersey. 2020 Aug 03. URL: <https://outinjersey.net/trump-administration-change-will-eliminate-non-discrimination-protections-for-lgbts-in-health-care/> [accessed 2021-03-21]
42. Mohr JW, Bogdanov P. Introduction—topic models: what they are and why they matter. *Poetics (Amst)* 2013 Dec;41(6):545-569. [doi: [10.1016/j.poetic.2013.10.001](https://doi.org/10.1016/j.poetic.2013.10.001)]
43. Maier D, Waldherr A, Miltner P, Wiedemann G, Niekler A, Keinert A, et al. Applying LDA topic modeling in communication research: toward a valid and reliable methodology. *Commun Methods Meas* 2018 Feb 16;12(2-3):93-118. [doi: [10.1080/19312458.2018.1430754](https://doi.org/10.1080/19312458.2018.1430754)]

44. Huang J. Maximum likelihood estimation of Dirichlet distribution parameters. CMU Technique Report 2005:1-9 [[FREE Full text](#)]
45. Buchanan L, Bui Q, Patel J. Black Lives Matter may be the largest movement in US history. The New York Times. 2020 Jul 03. URL: <https://www.nytimes.com/interactive/2020/07/03/us/george-floyd-protests-crowd-size.html> [accessed 2021-03-21]
46. Slodysko B. Why AP called the 2020 election for Joe Biden. The Associated Press. 2020 Nov 07. URL: <https://apnews.com/article/why-did-AP-call-election-for-Biden-fe79276cd9175fffc7cf4fb58045fcf9> [accessed 2021-03-22]
47. Berg K, Syed M. Under Trump, LGBTQ progress is being reversed in plain sight. Pro Publica. 2019 Nov 22. URL: <https://projects.propublica.org/graphics/lgbtq-rights-rollback> [accessed 2021-03-23]
48. Branson-Poots H, Stiles M. All Black Lives Matter march calls for LGBTQ rights and racial justice. Los Angeles Times; 2020 Jun 15. 2020 Jun 15. URL: <https://www.latimes.com/california/story/2020-06-15/lgbtq-pride-black-lives-controversy> [accessed 2021-03-21]
49. Riley J. LGBTQ advocates call for Trump's removal from office following unrest at US Capitol. Metro Weekly. 2021 Jan 07. URL: <https://www.metroweekly.com/2021/01/lgbtq-advocates-call-for-trumps-removal-from-office-following-unrest-at-u-s-capitol/> [accessed 2021-03-24]
50. Map of COVID-19 case trends, restrictions and mobility. USA Today. 2020 Apr 04. URL: <https://www.usatoday.com/storytelling/coronavirus-reopening-america-map/> [accessed 2021-03-22]
51. Coronavirus (COVID-19) update: FDA authorizes Pfizer-BioNTech COVID-19 Vaccine for emergency use in adolescents in another important action in fight against pandemic. U.S. Food & Drug Administration. 2021 Apr 10. URL: <https://www.fda.gov/news-events/press-announcements/coronavirus-covid-19-update-fda-authorizes-pfizer-biontech-covid-19-vaccine-emergency-use> [accessed 2021-06-11]
52. Suinn RM. The terrible twos--anger and anxiety. hazardous to your health. Am Psychol 2001 Jan;56(1):27-36. [doi: [10.1037/0003-066x.56.1.27](https://doi.org/10.1037/0003-066x.56.1.27)] [Medline: [11242985](https://pubmed.ncbi.nlm.nih.gov/11242985/)]
53. Weisz J, Kazdin E. Evidence-Based Psychotherapies for Children and Adolescents. New York City: The Guilford Press; Jun 13, 2017.
54. Paceley MS, Okrey-Anderson S, Fish JN, McInroy L, Lin M. Beyond a shared experience: queer and trans youth navigating COVID-19. Qual Soc Work 2021 Mar 16;20(1-2):97-104 [[FREE Full text](#)] [doi: [10.1177/1473325020973329](https://doi.org/10.1177/1473325020973329)] [Medline: [34025216](https://pubmed.ncbi.nlm.nih.gov/34025216/)]
55. Hatzembuehler ML. The social environment and suicide attempts in lesbian, gay, and bisexual youth. Pediatrics 2011 May;127(5):896-903 [[FREE Full text](#)] [doi: [10.1542/peds.2010-3020](https://doi.org/10.1542/peds.2010-3020)] [Medline: [21502225](https://pubmed.ncbi.nlm.nih.gov/21502225/)]
56. Kaniuka A, Pugh KC, Jordan M, Brooks B, Dodd J, Mann AK, et al. Stigma and suicide risk among the LGBTQ population: are anxiety and depression to blame and can connectedness to the LGBTQ community help? J Gay Lesbian Ment Health 2019 Mar 08;23(2):205-220. [doi: [10.1080/19359705.2018.1560385](https://doi.org/10.1080/19359705.2018.1560385)]
57. Implications of COVID-19 for LGBTQ youth mental health and suicide prevention. The Trevor Project. 2020 Apr 03. URL: <https://www.thetrevorproject.org/2020/04/03/implications-of-covid-19-for-lgbtq-youth-mental-health-and-suicide-prevention/> [accessed 2021-03-03]
58. Bantum EO, Owen JE. Evaluating the validity of computerized content analysis programs for identification of emotional expression in cancer narratives. Psychol Assess 2009 Mar;21(1):79-88. [doi: [10.1037/a0014643](https://doi.org/10.1037/a0014643)] [Medline: [19290768](https://pubmed.ncbi.nlm.nih.gov/19290768/)]
59. Gueyffier F, Cucherat M. The limitations of observation studies for decision making regarding drugs efficacy and safety. Therapie 2019 Apr;74(2):181-185 [[FREE Full text](#)] [doi: [10.1016/j.therap.2018.11.001](https://doi.org/10.1016/j.therap.2018.11.001)] [Medline: [30514576](https://pubmed.ncbi.nlm.nih.gov/30514576/)]

Abbreviations

LDA: latent Dirichlet allocation

LGBTQ+: Lesbian, Gay, Bisexual, Transgender, Queer/Questioning, and Others

LIWC: Linguistic Inquiry and Word Count program

Edited by T Sanchez; submitted 23.03.21; peer-reviewed by M Paceley, R Gregson, K Jacques; comments to author 02.06.21; revised version received 22.06.21; accepted 15.07.21; published 17.08.21.

Please cite as:

Stevens HR, Acic I, Rhea S

Natural Language Processing Insight into LGBTQ+ Youth Mental Health During the COVID-19 Pandemic: Longitudinal Content Analysis of Anxiety-Provoking Topics and Trends in Emotion in LGBTeens Microcommunity Subreddit

JMIR Public Health Surveill 2021;7(8):e29029

URL: <https://publichealth.jmir.org/2021/8/e29029>

doi:[10.2196/29029](https://doi.org/10.2196/29029)

PMID:[34402803](https://pubmed.ncbi.nlm.nih.gov/34402803/)

©Hannah R Stevens, Irena Acic, Sofia Rhea. Originally published in JMIR Public Health and Surveillance (<https://publichealth.jmir.org>), 17.08.2021. This is an open-access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR Public Health and Surveillance, is properly cited. The complete bibliographic information, a link to the original publication on <https://publichealth.jmir.org>, as well as this copyright and license information must be included.

Original Paper

The Roles of General Health and COVID-19 Proximity in Contact Tracing App Usage: Cross-sectional Survey Study

Dirk Witteveen^{1*}, PhD; Pablo de Pedraza^{2*}, PhD

¹Nuffield College, University of Oxford, Oxford, United Kingdom

²European Commission, DG Joint Research Centre, Directorate I – Competences, Unit I.1 - Monitoring, Indicators and Impact Evaluation, Ispra (VA), Italy

* all authors contributed equally

Corresponding Author:

Pablo de Pedraza, PhD

European Commission, DG Joint Research Centre

Directorate I – Competences

Unit I.1 - Monitoring, Indicators and Impact Evaluation

Via E. Fermi 2749, TP 361

Ispra (VA), I-21027

Italy

Phone: 39 033278380

Email: pablo.depedraza@ec.europa.eu

Abstract

Background: Contact tracing apps are considered useful means to monitor SARS-CoV-2 infections during the off-peak stages of the COVID-19 pandemic. Their effectiveness is, however, dependent on the uptake of such COVID-19 apps.

Objective: We examined the role of individuals' general health status in their willingness to use a COVID-19 tracing app as well as the roles of socioeconomic characteristics and COVID-19 proximity.

Methods: We drew data from the WageIndicator Foundation Living and Working in Coronavirus Times survey. The survey collected data on labor market status as well as the potential confounders of the relationship between general health and COVID-19 tracing app usage, such as sociodemographics and regular smartphone usage data. The survey also contained information that allowed us to examine the role of COVID-19 proximity, such as whether an individual has contracted SARS-CoV-2, whether an individual has family members and colleagues with COVID-19, and whether an individual exhibits COVID-19 pandemic-induced depressive and anxiety symptoms. We selected data that were collected in Spain, Italy, Germany, and the Netherlands from individuals aged between 18 and 70 years (N=4504). Logistic regressions were used to measure individuals' willingness to use a COVID-19 tracing app.

Results: We found that the influence that socioeconomic factors have on COVID-19 tracing app usage varied dramatically between the four countries, although individuals experiencing forms of not being employed (ie, recent job loss and inactivity) consistently had a lower willingness to use a contact tracing app (effect size: 24.6%) compared to that of employees (effect size: 33.4%; $P<.001$). Among the selected COVID-19 proximity indicators, having a close family member with SARS-CoV-2 infection was associated with higher contact tracing app usage (effect size: 36.3% vs 27.1%; $P<.001$). After accounting for these proximity factors and the country-based variations therein, we found that having a poorer general health status was significantly associated with a much higher likelihood of contact tracing app usage; compared to a self-reported "very good" health status (estimated probability of contact tracing app use: 29.6%), the "good" (estimated probability: +4.6%; 95% CI 1.2%-8.1%) and "fair or bad" (estimated probability: +6.3%; 95% CI 2.3%-10.3%) health statuses were associated with a markedly higher willingness to use a COVID-19 tracing app.

Conclusions: Current public health policies aim to promote the use of smartphone-based contact tracing apps during the off-peak periods of the COVID-19 pandemic. Campaigns that emphasize the health benefits of COVID-19 tracing apps may contribute the most to the uptake of such apps. Public health campaigns that rely on digital platforms would also benefit from seriously considering the country-specific distribution of privacy concerns.

(JMIR Public Health Surveill 2021;7(8):e27892) doi:[10.2196/27892](https://doi.org/10.2196/27892)

KEYWORDS

COVID-19; contact tracing; socioeconomic factors; labor market status; privacy; data sharing; pandemic; mobile health; public health; smartphone; mobile phone

Introduction

Over the course of 2020, governments have adopted a range of strategies to reduce the spread of COVID-19—an infectious disease that resulted in a pandemic—while trying to keep their economies afloat. Since mobility restrictions were gradually lifted during the last phase of the first pandemic wave in Europe, contact tracing has been considered to be an effective method for disease control, particularly for preventing disease transmission via contagious individuals who are not (yet) symptomatic [1-3]. Several governments have rolled out a version of a COVID-19 contact tracing app to help identify individuals who have been in close physical contact with an infected individual. However, as contact tracing apps inevitably rely on the collection of personal health data and mobility data, privacy concerns have been raised among the public [4].

In Europe, where participation in contact tracing via smartphone apps is voluntary, the effectiveness of contact tracing is dependent on the uptake of such apps. One study showed that in order to successfully suppress virus transmission during the peak of an outbreak in a hypothetical city with 1 million inhabitants, about 80% of all smartphone users or 56% of the population aged under 70 years would have to install the contact tracer [5]. Further, by modeling data from Washington State, researchers found that over the course of 300 days in 2020, infections and deaths could be reduced by 8% and 6%, respectively, if only 15% of the population were to participate in digital contact tracing [6]. Other researchers have also found that app-based tracing remains a more effective system than conventional contact tracing if coverage exceeds 20% [7]. Thus, although COVID-19 tracing apps are relatively ineffective during pandemic spikes, they can still help to slow the spread of SARS-COV-2 in subsequent periods, even though these apps have relatively low coverage.

This study concentrated on the association between individuals' general health status and their willingness to use a COVID-19 tracing app across several European countries and focused on Spain, Italy, Germany, and the Netherlands. These attitudes were measured in the fall of 2020, which was when the daily number of new cases was increasing (ie, the “second wave”). We ask the following question: are poorer health statuses associated with a higher willingness to share personal information in COVID-19 tracing apps? This dynamic could occur if individuals prioritize personal or public health concerns over possible data privacy concerns. Such health risk calculations tend to operate differently for individuals who perceive themselves to be more vulnerable. The higher sense of danger among at-risk groups is often found to positively correlate with distress [8]. Furthermore, research has shown that individuals' level of engagement with disease prevention behavior increases as soon as they are able to translate an abstract societal risk into a likelihood of experiencing a disease's most severe consequences [9,10].

We also concentrated on the moderating role of COVID-19 proximity in the relationship between general health status and the willingness to use a contact tracing app. This is because risk perceptions are known to be partially influenced by the experiences of other individuals in one's social circle, such as family, friends, and colleagues [10,11]. Recent studies have also suggested that individuals' risk behaviors are rather susceptible to information treatments about COVID-19 during the pandemic. For instance, learning about the severe symptoms of COVID-19 positively influences a range of protective behaviors [12] and results in individuals being less accepting of the incautious behavior of others [13]. In other words, first-hand physical and psychological experiences and observations of nearby people being affected by COVID-19 are likely to impact individuals' attitudes and risk behaviors. Hence, we examined the relationship between a set of COVID-19 proximity indicators and individuals' willingness to install and use a contact tracing app. We used indicators such as being tested for COVID-19, having a close family member or colleague with COVID-19, and self-reporting depression and anxiety symptoms resulting from the pandemic.

We also addressed the role that individuals' socioeconomic characteristics have as covariates of COVID-19 tracing app support. Such characteristics included gender, migration status, age, household status, and labor market status. It is important to account for these factors because of their expected relationship with the dependent variable (COVID-19 tracing app usage). As contact tracing systems are being rapidly rolled out by current administrations, skepticism toward COVID-19 tracing apps may be rooted in general distrust toward the government, which is why the selected sociodemographics served as necessary control variables [14]. Sociodemographic factors are also predictive of general smartphone app usage and COVID-19 tracing app installation, as shown in recent studies [15]. Furthermore, it is important to account for possible general health effect heterogeneity across the aforementioned individual-level socioeconomic characteristics [16], which could also vary across European countries [15].

In sum, we expected to find significant self-reported general health status gradients in individuals' willingness to use a COVID-19 tracing app that are dependent on a range of socioeconomic characteristics. This relationship could be mediated by country-specific associations between socioeconomic attributes and COVID-19 tracing app support. We also hypothesized that as observable pandemic-related health risks increase for individuals, their willingness to use a COVID-19 tracing app also increases.

Methods

Data

Observational data were drawn from the WageIndicator Living and Working in Coronavirus Times (LWCV) survey, which was filled out by web respondents between week 42 and week

49 of 2020 [17]. Respondents provided consent for their data be used in scientific research and did not receive financial compensation for their participation in the survey. All individual-level data were anonymized by WageIndicator prior to their use by academic researchers. The data set used and the analyses conducted did not contain identifiable information.

We selected respondents aged between 18 and 70 years ($N=4504$) from Spain ($n=1936$), Italy ($n=562$), Germany ($n=1294$), and the Netherlands ($n=712$). This was because adults aged up to 70 years have relatively large social networks and stronger connections to the labor market (eg, coworkers). Contact tracing is also believed to be the most effective when it is performed with this population [5]. The LWCV survey collects data about family structure, COVID-19 testing, self-perceived health status, and depressive and anxiety symptoms. It also contains a series of questions about individuals' willingness to use a COVID-19 tracing app as well as data on relevant confounders, such as general smartphone and app usage. [Multimedia Appendix 1](#) contains the sections of the questionnaire that were used for this study.

Voluntary web surveys have become common data collection tools during the pandemic. A range of policy-relevant studies that documented the initial impact that COVID-19 has on health, work, personal, and family situations relied on data from voluntary web surveys [18–21]. Two important advantages of this data collection method are that sampling is continuous and that questionnaires can be adjusted to rapidly changing situations, such as the 2020 COVID-19 pandemic. A significant drawback of voluntary web surveys is that the samples are not representative of the full population (ie, individuals who use and do not use web-based platforms). The results of such surveys therefore have to be interpreted with caution. The application of poststratification techniques can help to partly correct the bias resulting from self-selection and underrepresentation [22].

The WageIndicator Foundation is a global research organization that relies on a long-standing survey of workforces across 150 countries. The WageIndicator Foundation website receives millions of visitors annually. The WageIndicator Foundation has produced reliable estimates of mental health, data on subjective feelings such as well-being and insecurity, and web survey weighting techniques for balancing selectivity bias [23–25]. During the COVID-19 pandemic, it has enabled the exploration of mental health, anxiety, and life satisfaction determinants [26,27]. In [Multimedia Appendix 2](#), we benchmark the LWCV study samples against those of the European Social Survey based on key sociodemographics; relatively comparable sample distributions across age and respondents' highest education level are displayed. However, the LWCV study samples contained slightly more individuals from the 30- to 54-year age groups than those in the general population, with the exception of Spain's population (more individuals from the ≥ 55 -year age group). [Multimedia Appendix 3](#) indicates that including European Social Survey-based weights led to same substantive conclusions. In accordance with recent studies that used LWCV survey data [26,27], we report unweighted estimates for the main findings. [Multimedia Appendix 4](#) documents the model statistics and model specification checks.

Measures

Data on COVID-19 tracing app support were derived by asking whether a respondent was willing to share both their health status and geographical location on a COVID-19 tracing app (yes vs no or do not know). The key independent variable was self-reported general health status, which was based on the following question: "How would you rate your overall health?" Respondents answered with "very good" (832/4504, 18.5%), "good" (2409/4504, 53.5%), "fair" (1082/4504, 24%), "bad" (153/4504, 3.4%), and "very bad" (28/4504, 0.6%). We merged the smaller categories—the "bad" and "very bad" categories—with the "fair" health status to aid with interpretation and used the "very good" category as the reference. Indicators of COVID-19 proximity were measured with questions on whether a close colleague or a family member has ever contracted COVID-19, one's own COVID-19 test-taking status and their results (none, positive, negative, and awaiting result), and self-reported COVID-19 pandemic-induced depression symptoms (5-point Likert scale) and COVID-19 pandemic-induced anxiety symptoms (5-point Likert scale). Socioeconomic variables included gender, age (age group), migration background (dichotomous), partnership status (whether partners are present in the household), whether children aged under 18 years were present in the household, the highest education level obtained (low, medium, and high), urbanicity (3 categories), labor market position (employee, freelance, self-employed, inactive, and other), and how labor market position has been affected by the COVID-19 pandemic in terms of job loss and income reduction. All models accounted for the timing of the survey (week number). [Multimedia Appendix 5](#) depicts the correlation matrix.

Estimation

In order to gain a thorough understanding of the critical structural pathways for contact tracing app uptake and its potential country-based variation, we first estimated the marginal effects of the socioeconomic factors and COVID-19 proximity indicators by using 2 separate series of logit models. We present the results of the bivariate models (independent variables and outcome only) and multivariate models (all independent variables combined). These analyses also included an overall model with country-fixed effects, which allowed us to account for dynamics that are altogether country specific (eg, debates on general data privacy and its consequences for people's trust in governments). Aside from learning about the relevance of these covariates for social and health policies, they also informed us about how the relationship between general health status and contact tracing app usage should be modeled. Two-sided significance tests ($\alpha=.05$) were performed for all analyses.

We also estimated COVID-19 tracing app support (Y_{prob}) based on the general health status indicator (H) in nested models; socioeconomic variables and COVID-19 proximity variables were added in separate steps (equation 1). The use of nested models allowed us to examine the mechanism for explaining how general health status is related to the level of contact tracing app support. The baseline model only contained country-fixed effects (F) and a control for survey week (W). In a second series of models, we added the socioeconomic matrix (D) and a

variable matrix for respondents' regular smartphone usage (P), that is, the ownership of other apps that collect health and geographic location data (dichotomous) and the total number of phone apps. The third series of models were further adjusted for the COVID-19 proximity indicators (C). We calculated average marginal effects to aid our interpretation of the coefficients, as per social science conventions [28]. In [Multimedia Appendix 6](#), we replicate the key results by using country random intercepts, which present the same quantitative results as those of the reported country-fixed effects models. Equation 1 is as follows:

$$Y_{\text{prob}} = \beta_0 + H_i\beta_1 + F_i\gamma + D_i\omega + C_i\phi + P_i\delta + \varepsilon_i \quad (1)$$

Results

Descriptive Statistics

[Table 1](#) presents the number and proportion of respondents and the average proportion of respondents who support a COVID-19 tracing app for each of the independent variables. It should be noted that both of the Southern European countries have a much higher average proportion of respondents who support contact tracing apps (Italy: 282/562, 50.2%; Spain: 716/1936, 37%). Both Germany (209/1294, 16.2%) and the Netherlands (127/712, 17.8%) display distinctly lower levels of support for a COVID-19 tracing app.

Table 1. Descriptive statistics for COVID-19 tracing app usage in Spain, Italy, Germany, and the Netherlands. Data are from weeks 42 through 49 (year: 2020; N=4504).

Variables	Respondents, n (proportion)	Proportion of respondents who support COVID-19 tracing apps
Key independent variable		
Health status		
Strong	832 (0.185)	0.209
Good	2409 (0.535)	0.310
Fair or bad	1263 (0.280)	0.327
Independent variables		
Gender		
Woman	1373 (0.305)	0.268
Man	3131 (0.695)	0.309
Migration background		
Native-born	4323 (0.960)	0.298
Foreign-born	181 (0.040)	0.249
Age group (years)		
18-29	667 (0.148)	0.237
30-44	1489 (0.331)	0.320
45-54	1286 (0.286)	0.297
55-70	1062 (0.236)	0.298
Partnership status		
No partner	1617 (0.359)	0.276
Partner in household	2887 (0.641)	0.307
Children (in the household)		
No children	2679 (0.595)	0.299
Children	1825 (0.405)	0.292
Highest education level		
Low	933 (0.207)	0.251
Medium	1625 (0.361)	0.260
High	1946 (0.432)	0.348
Labor market position		
Employee	2419 (0.537)	0.324
Freelance	189 (0.042)	0.307
Self-employed with employees	59 (0.013)	0.203
Other employment	151 (0.034)	0.311
Job loss and income reduction due to the COVID-19 pandemic	861 (0.191)	0.253
Inactive	825 (0.183)	0.261
Urbanicity		
City or metropole	2454 (0.545)	0.313
Small city or town	1297 (0.288)	0.282
Village or rural	753 (0.167)	0.264
COVID-19 pandemic-induced depression symptoms		
Disagree	1773 (0.394)	0.267
Neutral	1084 (0.241)	0.318

Variables	Respondents, n (proportion)	Proportion of respondents who support COVID-19 tracing apps
Agree	1647 (0.366)	0.313
COVID-19 pandemic-induced anxiety symptoms		
Disagree	1415 (0.314)	0.246
Neutral	1090 (0.242)	0.255
Agree	1999 (0.444)	0.354
COVID-19 test		
No	2976 (0.661)	0.269
Yes, positive	139 (0.031)	0.381
Yes, awaiting result	20 (0.004)	0.400
Yes, negative	1369 (0.304)	0.344
Close colleague with COVID-19		
No	2060 (0.457)	0.265
Yes	1242 (0.276)	0.368
Do not know or N/A ^a	1202 (0.267)	0.285
Family member with COVID-19		
No	3200 (0.710)	0.263
Yes	1170 (0.260)	0.387
Do not know or N/A	134 (0.030)	0.291

^aN/A: not applicable.

Socioeconomic Factors

Table 2 presents the marginal effects that socioeconomic factors had on the willingness to use a COVID-19 tracing app among the full sample and the four countries separately. The bivariate associations (marginal effect sizes) in the *Country-fixed effects* column suggest that older individuals (aged 45-54 years: 30.9%; aged 55-70 years: 31.5%; $P<.001$) are significantly more willing to use a COVID-19 tracing app than young adults (about 22.5%). Partnered individuals who also live in the same household are also more likely to use a contact tracing app than nonpartnered individuals, as indicated by the 4% marginal effects gap. Individuals with medium (effect size: 30.1%) and high (effect size: 31.6%) levels of education had a significantly higher willingness to use a COVID-19 tracing app compared to that of individuals with low levels of education (effect size: 24.7%; $P<.001$). Furthermore, compared to employees (effect size: 33.4%), individuals who are not active in the labor force (effect size: 24.6%; $P<.001$) and those who lost their job or income during the COVID-19 pandemic (effect size: 25.2%; $P<.001$) are significantly less likely to use a contact tracing app. These

independent socioeconomic variables remained statistically significant in the multivariate model. Notably, gender, migration background, and urbanicity are not associated with the probability of using a COVID-19 tracing app.

The columns of Table 2 present the bivariate and multivariate model results for each of the four countries. The main finding from these models is the striking cross-national variation in the relationship between socioeconomic variables and COVID-19 tracing app usage. The impact that respondents from Spain had on the full-sample results appears to be substantial, as the country largely exhibits the same socioeconomic relationships in terms of the significance levels and magnitudes reported by the multivariate models. However, in Italy, the most important socioeconomic factors are labor market status and urbanicity. Relative to employees (effect size: 53.5%), respondents who experienced a recent job loss (effect size: 40.2%; $P=.03$) and inactive respondents (effect size: 32.2%; $P=.04$) had significantly lower probabilities of being willing to use a COVID-19 tracing app. These substantial differences were derived from multivariate models that accounted for important confounders, such as age and household status.

Table 2. The marginal effects that sociodemographic factors have on respondents' willingness to use a COVID-19 tracing app. Data are from weeks 42 through 49 (2020).

Sociodemographic factors	Country-fixed effects		Spain		Italy		Germany		The Netherlands	
	Bivariate model	Multivariate model	Bivariate model	Multivariate model	Bivariate model	Multivariate model	Bivariate model	Multivariate model	Bivariate model	Multivariate model
Gender										
Woman (referent)	0.302	0.304	0.376	0.379	0.509	0.512	0.176	0.175	0.170	0.173
Man	0.282	0.279	0.356	0.348	0.477	0.467	0.135	0.136	0.197	0.191
Migration background										
Native-born (referent)	0.298	0.298	0.372	0.373	0.504	0.503	0.165	0.165	0.178	0.179
Foreign-born	0.248	0.244	0.327	0.321	0.375	0.396	0.091	0.084	0.214	0.166
Age group (years)										
18-29 (referent)	0.225	0.227	0.231	0.235	0.466	0.494	0.213	0.197	0.165	0.158
30-44	0.306 ^a	0.300 ^a	0.381 ^a	0.374 ^a	0.536	0.527	0.184	0.189	0.122	0.115
45-54	0.309 ^a	0.314 ^a	0.426 ^a	0.428 ^a	0.439	0.448	0.133 ^a	0.145	0.225	0.232
55-70	0.315 ^a	0.316 ^a	0.412 ^a	0.418 ^a	0.537	0.528	0.152	0.137	0.184	0.191
Partnership										
No partner (referent)	0.271	0.277	0.333	0.357	0.462	0.474	0.153	0.133	0.167	0.173
Partner in household	0.311 ^a	0.307 ^a	0.392 ^a	0.377	0.525	0.518	0.166	0.179 ^a	0.184	0.181
Children (in household)										
No children (referent)	0.297	0.310	0.355	0.380	0.503	0.508	0.185	0.195	0.181	0.189
Children	0.295	0.277	0.390	0.357	0.500	0.489	0.128 ^a	0.118 ^a	0.175	0.164
Highest education level										
Low (referent)	0.247	0.246	0.330	0.326	— ^b	—	0.104	0.105	0.113	0.113
Medium	0.301 ^a	0.304 ^a	0.366	0.368	0.507	0.517	0.172 ^a	0.171 ^a	0.173	0.173
High	0.316 ^a	0.314 ^a	0.396 ^a	0.398 ^a	0.499	0.493	0.171 ^a	0.171 ^a	0.227 ^a	0.227 ^a
Labor market position										
Employee (referent)	0.334	0.327	0.434	0.411	0.536	0.535	0.181	0.181	0.186	0.183
Freelance	0.277	0.277	0.333	0.347	0.500	0.509	0.125	0.113	0.214	0.212
Self-employed with employees	0.246	0.237	0.316	0.305	0.333	0.316	0.094	0.096	0.400	0.354
Other employment	0.294	0.298	0.372	0.382	0.429	0.409	0.241	0.228	0.095	0.107
Job loss and income reduction due to the COVID-19 pandemic	0.252 ^a	0.249 ^a	0.336 ^a	0.327 ^a	0.397 ^a	0.402 ^a	0.100 ^a	0.101 ^a	0.229	0.232
Inactive	0.246 ^a	0.266 ^a	0.311 ^a	0.345 ^a	0.321 ^a	0.322 ^a	0.188	0.189	0.068 ^a	0.078 ^a
Urbanicity										
City or metropole (referent)	0.296	0.294	0.371	0.367	0.458	0.459	0.176	0.176	0.206	0.201
Small city or town	0.293	0.295	0.361	0.366	0.558 ^a	0.560 ^a	0.146	0.147	0.170	0.172
Village or rural	0.303	0.306	0.383	0.390	0.727 ^a	0.712 ^a	0.140	0.140	0.163	0.163

^aSignificant at the $P < .05$ level (two-tailed tests).^bNot available.

For Germany and the Netherlands, where COVID-19 tracing app support is, on average, much lower than that in the two Southern European countries, educational level and labor market position were the only statistically significant independent variables of the willingness to use a COVID-19 tracing app in the multivariate models. Holding higher education credentials in Germany and the Netherlands yielded a 6.6% higher marginal effect (Germany: $P=.04$) and an 11.4% higher marginal effect (the Netherlands: $P=.01$), respectively, on COVID-19 tracing app support compared to those for holding lower education credentials (reference group). In Germany, individuals who recently experienced a job loss or income loss are significantly less likely to use a contact tracing app (effect size: 10.1%; $P=.01$) than employees (effect size: 18.1%). In the Netherlands, the marginal effect for support for a COVID-19 tracing app is only 7.8% ($P=.04$) for those who remained inactive during the pandemic, regardless of other socioeconomic factors. Finally, in Germany, where age remained a nonsignificant independent variable, the presence of a partner ($P=.04$) and children ($P<.001$) in the household are positively associated with respondents' willingness to use a contact tracing app.

COVID-19 Proximity

We also examined key indicators of COVID-19 proximity. Table 3 presents the marginal effects that these indicators had on the willingness to use a COVID-19 tracing app. As shown in the *Country-fixed effects* column (the combined sample), all COVID-19 proximity factors had significant positive associations with respondents' support for a COVID-19 tracing app. Specifically, having (potentially) contracted COVID-19 was suggestive of a higher willingness to use a COVID-19

tracing app, as indicated by the substantially higher marginal effects of being tested (positive: 35.7%; negative: 33%; awaiting results: 41.2%). However, these marginal effects ceased to be statistically significantly different from those of the reference group (no test: 27.6%) in the multivariate models. The multivariate model further indicated that reporting anxiety symptoms (the 33.2% marginal effect of the "agree" response vis-à-vis the 27.1% marginal effect of the neutral response; $P<.001$) was significantly associated with contact tracing app support. Similarly, having a family member (effect size: 35.2%; $P<.001$) or colleague (effect size: 34.9%; $P<.001$) who has ever contracted COVID-19 was also associated with greater contact tracing app support compared to not having such a family member (effect size: 27.4%) or colleague (effect size: 28.4%).

Similar to socioeconomic factors, the relationship between COVID-19 proximity and the willingness to use a COVID-19 tracing app appears to vary across the four countries studied. For instance, having no anxiety symptoms had a positive effect on COVID-19 tracing app support in Germany. Furthermore, having a colleague who tested positive for COVID-19 was associated with a higher likelihood of using a contact tracing app overall, but this was not the case in the Netherlands. In fact, having a family member who has ever contracted COVID-19 was the only indicator that yielded a significant positive marginal effect on contact tracing support across all countries. This implies that, in addition to including country-fixed effects and controls for sociodemographics and COVID-19 proximity in models, modeling the association between general health and COVID-19 tracing app support required us to ensure that the controls interact with the country dummies to account for heterogeneity.

Table 3. Marginal effects that COVID-19 proximity indicators had on the willingness to use a COVID-19 tracing app. Data are from weeks 42 through 49 (2020).

COVID-19 proximity indicators	Country-fixed effects		Spain		Italy		Germany		The Netherlands	
	Bivariate model	Multivariate model	Bivariate model	Multivariate model	Bivariate model	Multivariate model	Bivariate model	Multivariate model	Bivariate model	Multivariate model
COVID-19 pandemic-induced depression symptoms										
Disagree	0.291	0.298	0.380	0.403	0.523	0.539	0.131	0.135	0.186	0.206
Neutral (referent)	0.325	0.321	0.430	0.417	0.543	0.520	0.166	0.165	0.166	0.156
Agree	0.283 ^a	0.280	0.335 ^a	0.329 ^a	0.456	0.453	0.232 ^a	0.218	0.178	0.165
COVID-19 pandemic-induced anxiety symptoms										
Disagree	0.256	0.261	0.267 ^a	0.263 ^a	0.425	0.428	0.250 ^a	0.230 ^a	0.117	0.113 ^a
Neutral (referent)	0.274	0.271	0.369	0.362	0.413	0.414	0.134	0.141	0.171	0.178
Agree	0.334 ^a	0.332 ^a	0.431 ^a	0.437 ^a	0.570 ^a	0.568 ^a	0.115	0.119	0.268 ^a	0.272 ^a
COVID-19 test										
No (referent)	0.276	0.286	0.343	0.356	0.515	0.524	0.127	0.134	0.184	0.197
Yes, positive	0.357 ^a	0.316	0.451	0.433	0.526	0.499	0.200	0.173	0.235	0.202
Yes, awaiting result	0.412	0.411	0.571	0.608	0.333	0.289	0.333	0.250	— ^b	—
Yes, negative	0.330 ^a	0.313	0.408 ^a	0.385	0.475	0.461	0.260 ^a	0.235 ^a	0.159	0.140
Close colleague with COVID-19										
No (referent)	0.277	0.284	0.363	0.374	0.454	0.464	0.132	0.144	0.192	0.209
Yes	0.364 ^a	0.349 ^a	0.448 ^a	0.420	0.550 ^a	0.543	0.262 ^a	0.219 ^a	0.208	0.197
Family member with COVID-19										
No (referent)	0.271	0.274	0.344	0.349	0.464	0.470	0.144	0.148	0.157	0.158
Yes	0.363 ^a	0.352 ^a	0.430 ^a	0.415 ^a	0.592 ^a	0.582 ^a	0.265 ^a	0.222 ^a	0.225 ^a	0.219 ^a

^aSignificant at the $P < .05$ level (two-tailed tests).^bNot available.

General Health Status

We illustrate the association between general health status and the willingness to use a COVID-19 tracing app in Figure 1. For the leftmost graph of Figure 1, we estimated a baseline model that adjusts for timing (survey week) and country-fixed effects. The graph suggests that both the “good” (estimated probability of contact tracing app use: +7.3%; $P < .001$) and “fair or bad” (estimated probability of contact tracing app use: +9.2%; $P < .001$) health statuses are positively associated with COVID-19 tracing app usage when compared to the reference category (the “strong” health status). The negative association between self-reported general health and COVID-19 tracing app usage persisted in terms of significance and magnitude when we also controlled for socioeconomic variables. However, as shown in the rightmost graph of Figure 1, the effect sizes of the poorer general health statuses were reduced by about one-third when adding the COVID-19 proximity variables; the estimated probability of contact tracing app use based on having a “good” health status and “fair or bad” health status increased by 4.6% ($P = .01$) and 6.3% ($P = .002$), respectively. This attenuation suggests that the relationship between general health status and

support for a contact tracing app partially operates through recent personal health scares or forms of stress that are related to direct experiences with the COVID-19 pandemic (ie, exhibiting depressive or anxiety symptoms and a diagnosis of COVID-19 among close contacts).

In order to correct the estimates for country heterogeneity in the relationship between the socioeconomic and COVID-19 proximity variables and the outcome variable, we analyzed the marginal effects that poorer general health has on COVID-19 tracing app support while adjusting for the interactions between all covariates and the country dummies (Figure 2). As the point estimates in the baseline and socioeconomic models presented in Figure 2 are nearly identical to those in Figure 1, we conclude that neither timing or socioeconomic heterogeneity influences the identified relationship between general health status and the willingness to use a COVID-19 tracing app across the four European countries. However, adjusting for the country heterogeneity in the COVID-19 proximity factors yielded evident null effects (rightmost graph in Figure 2). Hence, the association between (poorer) general health status and COVID-19 tracing app support varies across the four countries,

but this is likely due to the varying degree to which the tracing app support. COVID-19 proximity measures are associated with COVID-19

Figure 1. The marginal effects that poorer health statuses have on the willingness to use a COVID-19 tracing app. Country-fixed effects (Spain, Italy, Germany, and the Netherlands) are applied. Data are from weeks 42 through 49 (2020). The plots show the marginal effects of poorer health statuses, and the “strong” health status was used as the reference category. The baseline model only controlled for survey week. The error bars represent 95% CIs.

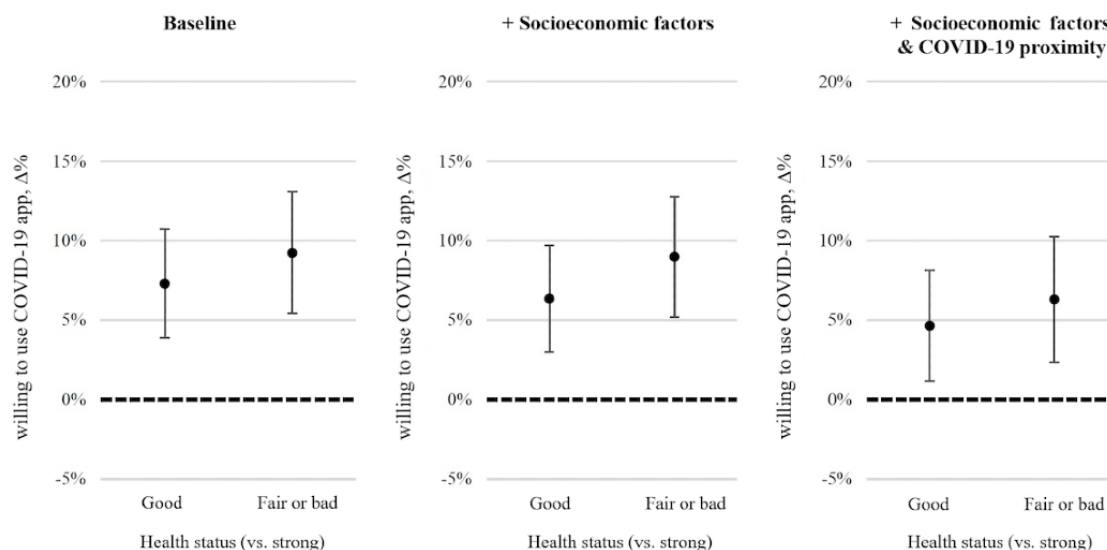
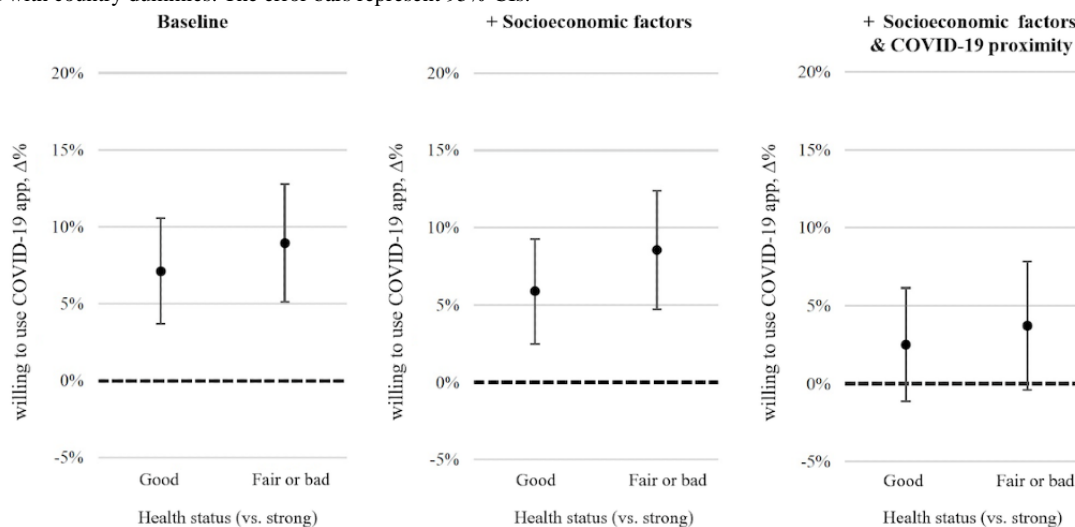


Figure 2. The marginal effects that poorer health statuses have on the willingness to use a COVID-19 tracing app. Country dummy interactions (Spain, Italy, Germany, and the Netherlands) with all covariates are applied. Data are from weeks 42 through 49 (2020). The plots show the marginal effects of poorer health statuses, and the “strong” health status was used as the reference category. The baseline model only controlled for survey week, which also interacted with country dummies. The error bars represent 95% CIs.



Discussion

Principal Findings

An important component of public health policy with regard to the spread of COVID-19 is the possibility of using mobile phone-based contact tracing in response to a positive COVID-19 test [2,3]. This is particularly relevant outside of major peaks in SARS-CoV-2 infection rates and lockdowns because it could help with avoiding rapid and uncontrollable disease transmission within communities [4]. COVID-19 tracing apps have been introduced in several countries. The extent to which these contact tracing apps can have a positive effect on public health (ie, reduce the chance of rapid outbreaks) is dependent on their uptake [5-7]. We found considerable country-based variation

in the willingness to use a COVID-19 tracing app, which ranged from 16.2% (209/1294) in Germany to 50.2% (282/562) in Italy. In addition, we found evidence indicating that several socioeconomic and demographic factors are associated with the willingness to use a COVID-19 contact tracing app. Such evidence was also found for the following COVID-19 proximity variables: exhibiting depression and anxiety symptoms, being tested, and having family members or colleagues who have ever contracted COVID-19. Importantly, based on these dynamics, we found that poorer health statuses are associated with significantly higher support for COVID-19 contact tracing apps.

We examined indicators of COVID-19 proximity because, based on exiting literature, we presumed that being confronted with the detrimental effects of the pandemic within one's social

network may trigger the motivation to install and use a COVID-19 tracing app [8-13]. Although we cannot observe changes in individuals' attitudes or behaviors in response to COVID-19 proximity over time, our findings support this relationship. We observed some cross-national variation in the significance levels of the COVID-19 proximity indicators. Nonetheless, having a family member has ever tested positive for COVID-19 appears to be the strongest and most consistent independent variable of the increased intended usage of a contact tracing app. We argue that these results are suggestive of a relationship between the social context of the consequences of the COVID-19 pandemic and individuals' perceived risk.

Given the cross-national variation in the associations between our two sets of control variables (socioeconomic factors and COVID-19 proximity), we fitted a comprehensive model that adjusted for this heterogeneity (ie, the heterogeneity resulting from the covariates interacting with the country-fixed effects) [14-16]. The results from this model were straightforward. We found that having a poorer general health status (ie, the "fair or bad" or "good" health status vis-à-vis the "very good" health status) positively affects the willingness to use a COVID-19 tracing app, even after adjusting for socioeconomic factors, indicators of how close the pandemic came to an individual (COVID-19 proximity indicators), and regular levels of smartphone usage. Our results indicate that this relationship is moderated by COVID-19 proximity. In other words, it is plausible that the association between general health status and contact tracing app uptake is affected by how close the pandemic has come to an individual.

It is also important to note that the recent loss of one's job or main income source during the pandemic yielded a significantly lower marginal effect on the willingness to use a contact tracing app in Spain, Italy, and Germany (differences of almost 10%), and similar results were observed for long-term inactivity in the Netherlands. We suspect that some of these associations reflect the impacts that economic security and insecurity have on sentiments regarding the pandemic or even a broader (structural) rejection of government (social) policies. These dynamics require much more in-depth research in the social science field.

Limitations

Our analyses relied on cross-sectional data that were obtained during the COVID-19 pandemic. We conducted several robustness checks to avoid having strong selection bias in the reported marginal effects of general health status, socioeconomic

characteristics, and COVID-19 proximity. Such bias is likely to result from the underrepresentation of individuals with the poorest general health conditions—a demographic group that is difficult to include in all kinds of social surveys. Nationally representative panels would be the preferred data source for future research on the relationship between health and any kind of COVID-19-related measures. This is because longitudinal data are better equipped to measure (deteriorating) health-related attrition. Furthermore, panel data are also best suited for effectively measuring the degree to which attitudes and behaviors of individuals change over time as a function of their health status or in response to the acquisition of new information. Future research may also benefit from expanding the operationalization of health risk and risk perceptions, such as those related to wearing a mask in close proximity to an individual.

Conclusions

This paper builds upon existing evidence indicating that contact tracing apps are an important element of public health [2,3] and that their positive effect is dependent on their uptake [5-7]. We studied whether general health status and COVID-19 proximity can be linked to contact tracing app uptake. This research question was motivated by a discussion in public health literature about the necessity of effective contact tracing in combating the COVID-19 pandemic as well as research in the social science field regarding the individual-level drivers of attitudes toward contact tracing apps [8-13]. We conclude that poorer general health statuses are positively associated with the willingness to use a COVID-19 tracing app. Moreover, the extent to which one's general health status impacts their likelihood of using COVID-19 tracing apps partially operates through the pandemic-related experiences that occur in their social circle.

To date, public debates have mainly revolved around issues regarding apps' capacity to meet data privacy goals and legislation criteria. We suspect that the country-based variation we found in people's willingness to use a COVID-19 tracing app reflects path-dependent societal dimensions, such as large personal data leaks in recent history or underlying distrust in the government [14]. This implies that public policies that are intended to expand the usage of digital COVID-19 contact tracing apps always have to consider country-specific societal concerns. Our study suggests that once these conditions are met, public health policies that aim to increase contact tracing app uptake would benefit from campaigns that stress these apps' benefits for users (both physical and mental benefits), their family members, and the economy [1].

Conflicts of Interest

None declared.

Multimedia Appendix 1

Survey questions.

[DOCX File, 20 KB - [publichealth_v7i8e27892_app1.docx](#)]

Multimedia Appendix 2

Key demographic proportions: European Social Survey and WageIndicator LWCV study sample.

[DOCX File, 22 KB - [publichealth_v7i8e27892_app2.docx](#)]

Multimedia Appendix 3

Replication with European Social Survey-derived weights.

[DOCX File, 43 KB - [publichealth_v7i8e27892_app3.docx](#)]

Multimedia Appendix 4

Model statistics and specification checks.

[DOCX File, 20 KB - [publichealth_v7i8e27892_app4.docx](#)]

Multimedia Appendix 5

Correlation matrix of all predictor variables.

[DOCX File, 25 KB - [publichealth_v7i8e27892_app5.docx](#)]

Multimedia Appendix 6

Associations of socioeconomic factors and COVID-19 proximity with country — random intercepts.

[DOCX File, 25 KB - [publichealth_v7i8e27892_app6.docx](#)]

References

1. Acemoglu D, Chernozhukov V, Werning I, Whinston MD. Optimal targeted lockdowns in a multi-group SIR model. National Bureau of Economic Research. 2020. URL: <https://www.nber.org/papers/w27102> [accessed 2020-11-03]
2. Eames KTD, Keeling MJ. Contact tracing and disease control. *Proc Biol Sci* 2003 Dec 22;270(1533):2565-2571 [FREE Full text] [doi: [10.1098/rspb.2003.2554](#)] [Medline: [14728778](#)]
3. Ferretti L, Wymant C, Kendall M, Zhao L, Nurtay A, Abeler-Dörner L, et al. Quantifying SARS-CoV-2 transmission suggests epidemic control with digital contact tracing. *Science* 2020 May 08;368(6491):eabb6936 [FREE Full text] [doi: [10.1126/science.abb6936](#)] [Medline: [32234805](#)]
4. Zhang B, Kreps S, McMurphy N, McCain RM. Americans' perceptions of privacy and surveillance in the COVID-19 pandemic. *PLoS One* 2020 Dec 23;15(12):e0242652. [doi: [10.1371/journal.pone.0242652](#)] [Medline: [33362218](#)]
5. Hinch R, Probert WJM, Nurtay A, Kendall M, Wymatt C, Hall M, et al. OpenABM-Covid19-an agent-based model for non-pharmaceutical interventions against COVID-19 including contact tracing. *medRxiv*. Preprint posted online on September 22, 2020. [FREE Full text] [doi: [10.1101/2020.09.16.20195925](#)]
6. Abueg M, Hinch R, Wu N, Liu L, Probert W, Wu A, et al. Modeling the effect of exposure notification and non-pharmaceutical interventions on COVID-19 transmission in Washington state. *NPJ Digit Med* 2021 Mar 12;4(1):49 [FREE Full text] [doi: [10.1038/s41746-021-00422-7](#)] [Medline: [33712693](#)]
7. Kretzschmar ME, Rozhnova G, Bootsma MCJ, van Boven M, van de Wijert JHHM, Bonten MJM. Impact of delays on effectiveness of contact tracing strategies for COVID-19: a modelling study. *Lancet Public Health* 2020 Aug;5(8):e452-e459 [FREE Full text] [doi: [10.1016/S2468-2667\(20\)30157-2](#)] [Medline: [32682487](#)]
8. Braun-Lewensohn O, Al-Sayed K. Syrian adolescent refugees: How do they cope during their stay in refugee camps? *Front Psychol* 2018 Jul 20;9:1258 [FREE Full text] [doi: [10.3389/fpsyg.2018.01258](#)] [Medline: [30079046](#)]
9. Glanz K, Rimer B, Viswanath K. *Health Behavior and Health Education: Theory, Research, and Practice*. San Francisco, California, United States: Jossey-Bass; 2008.
10. Kraut A, Graff L, McLean D. Behavioral change with influenza vaccination: factors influencing increased uptake of the pandemic H1N1 versus seasonal influenza vaccine in health care personnel. *Vaccine* 2011 Oct 26;29(46):8357-8363. [doi: [10.1016/j.vaccine.2011.08.084](#)] [Medline: [21888939](#)]
11. Kumar S, Quinn SC, Kim KH, Musa D, Hilyard KM, Freimuth VS. The social ecological model as a framework for determinants of 2009 H1N1 influenza vaccine uptake in the United States. *Health Educ Behav* 2012 Apr;39(2):229-243 [FREE Full text] [doi: [10.1177/1090198111415105](#)] [Medline: [21984692](#)]
12. de Bruin WB, Carman KG, Parker AM. Mental associations with COVID-19 and how they relate with self-reported protective behaviors: A national survey in the United States. *Soc Sci Med* 2021 Apr;275:113825 [FREE Full text] [doi: [10.1016/j.socscimed.2021.113825](#)] [Medline: [33735777](#)]
13. Lunn PD, Timmons S, Belton CA, Barjaková M, Julianne H, Lavin C. Motivating social distancing during the COVID-19 pandemic: An online experiment. *Soc Sci Med* 2020 Nov;265:113478 [FREE Full text] [doi: [10.1016/j.socscimed.2020.113478](#)] [Medline: [33162198](#)]
14. Citrin J, Stoker L. Political trust in a cynical age. *Annu Rev Polit Sci (Palo Alto)* 2018 May 11;21(1):49-70 [FREE Full text] [doi: [10.1146/annurev-polisci-050316-092550](#)]
15. Altmann S, Milsom L, Zillesen H, Blasone R, Gerdon F, Bach R, et al. Acceptability of app-based contact tracing for COVID-19: Cross-country survey study. *JMIR Mhealth Uhealth* 2020 Aug 28;8(8):e19857 [FREE Full text] [doi: [10.2196/19857](#)] [Medline: [32759102](#)]

16. Ivers LC, Weitzner DJ. Can digital contact tracing make up for lost time? *Lancet Public Health* 2020 Aug;5(8):e417-e418 [FREE Full text] [doi: [10.1016/S2468-2667\(20\)30160-2](https://doi.org/10.1016/S2468-2667(20)30160-2)] [Medline: [32682488](https://pubmed.ncbi.nlm.nih.gov/32682488/)]
17. WageIndicator Survey of Living and Working in Coronavirus Times 2020. IZA Institute of Labor Economics. URL: <https://datasets.iza.org/dataset/1388/living-and-working-in-coronavirus-times-survey> [accessed 2020-11-03]
18. Fetzer TR, Witte M, Hensel L, Jachimowicz J, Haushofer J, Ivchenko A, et al. Global behaviors and perceptions at the onset of the COVID-19 pandemic. National Bureau of Economic Research. 2020 May. URL: https://www.nber.org/system/files/working_papers/w27082/w27082.pdf [accessed 2021-04-20]
19. Lu H, Nie P, Qian L. Do quarantine experiences and attitudes towards COVID-19 affect the distribution of mental health in China? A quantile regression analysis. *Appl Res Qual Life* 2020 Jun 29:1-18 [FREE Full text] [doi: [10.1007/s11482-020-09851-0](https://doi.org/10.1007/s11482-020-09851-0)] [Medline: [32837605](https://pubmed.ncbi.nlm.nih.gov/32837605/)]
20. Baert S, Lippens L, Moens E, Sterkens P, Weytjens J. How do we think the COVID-19 crisis will affect our careers (if any remain)? IZA Institute of Labor Economics Discussion Paper Series. 2020 Apr. URL: <http://ftp.iza.org/dp13164.pdf> [accessed 2021-04-20]
21. Brodeur A, Clark AE, Fleche S, Powdthavee N. Assessing the impact of the coronavirus lockdown on unhappiness, loneliness, and boredom using Google Trends. arXiv. Preprint posted online on April 25, 2020. [FREE Full text]
22. Tourangeau R, Conrad FG, Couper MP. *The Science of Web Surveys*. Oxford, United Kingdom: Oxford University Press; 2013.
23. Guzi M, de Pedraza P. A web survey analysis of subjective well-being. *Int J Manpow* 2015;36(1):48-67. [doi: [10.1108/ijm-12-2014-0237](https://doi.org/10.1108/ijm-12-2014-0237)]
24. Pedraza PD, Tjidsen K, de Bustillo RM, Steinmetz S. A Spanish continuous volunteer web survey: Sample bias, weighting and efficiency. *Rev Esp Invest Sociol* 2010;131:109-130 [FREE Full text]
25. de Bustillo RM, de Pedraza P. Determinants of job insecurity in five European countries. *European Journal of Industrial Relations* 2010 Feb 19;16(1):5-20. [doi: [10.1177/0959680109355306](https://doi.org/10.1177/0959680109355306)]
26. de Pedraza P, Guzi M, Tjidsen K. Life dissatisfaction and anxiety in COVID-19 pandemic. European Commission Joint Research Centre Technical Report. 2020. URL: <https://core.ac.uk/download/pdf/326047012.pdf> [accessed 2021-01-19]
27. Witteveen D, Velthorst E. Economic hardship and mental health complaints during COVID-19. *Proc Natl Acad Sci U S A* 2020 Nov 03;117(44):27277-27284 [FREE Full text] [doi: [10.1073/pnas.2009609117](https://doi.org/10.1073/pnas.2009609117)] [Medline: [33046648](https://pubmed.ncbi.nlm.nih.gov/33046648/)]
28. Mood C. Logistic regression: Why we cannot do what we think we can do, and what we can do about it. *Eur Sociol Rev* 2009 Mar 09;26(1):67-82. [doi: [10.1093/esr/jcp006](https://doi.org/10.1093/esr/jcp006)]

Abbreviations

LWCV: Living and Working in Coronavirus Times

Edited by T Sanchez; submitted 11.02.21; peer-reviewed by F Kreute, F Wirth, A Khurshid; comments to author 04.03.21; revised version received 24.03.21; accepted 10.05.21; published 18.08.21.

Please cite as:

Witteveen D, de Pedraza P

The Roles of General Health and COVID-19 Proximity in Contact Tracing App Usage: Cross-sectional Survey Study

JMIR Public Health Surveill 2021;7(8):e27892

URL: <https://publichealth.jmir.org/2021/8/e27892>

doi: [10.2196/27892](https://doi.org/10.2196/27892)

PMID: [34081602](https://pubmed.ncbi.nlm.nih.gov/34081602/)

©Dirk Witteveen, Pablo de Pedraza. Originally published in *JMIR Public Health and Surveillance* (<https://publichealth.jmir.org>), 18.08.2021. This is an open-access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in *JMIR Public Health and Surveillance*, is properly cited. The complete bibliographic information, a link to the original publication on <https://publichealth.jmir.org>, as well as this copyright and license information must be included.

Viewpoint

With Great Hopes Come Great Expectations: Access and Adoption Issues Associated With COVID-19 Vaccines

Zhaohui Su¹, MA, PhD; Dean McDonnell², PhD; Ali Cheshmehzangi^{3,4}, PhD; Xiaoshan Li⁵, PhD; Daniel Maestro⁶, PhD; Sabina Šegalo⁷, PhD; Junaid Ahmad⁸, PhD; Xiaoning Hao⁹, PhD

¹Center on Smart and Connected Health Technologies, Mays Cancer Center, School of Nursing, UT Health San Antonio, San Antonio, TX, United States

²Department of Humanities, Institute of Technology Carlow, Carlow, Ireland

³Department of Architecture and Built Environment, University of Nottingham Ningbo China, Ningbo, China

⁴Hiroshima University, Hiroshima, Japan

⁵Program of Public Relations and Advertising, Beijing Normal University-Hong Kong Baptist University United International College, Zhuhai, China

⁶Department of Environmental Health, Institute for Public Health of Federation Bosnia and Herzegovina, Sarajevo, Bosnia and Herzegovina

⁷Department of Microbiology, Faculty of Medicine, University of Sarajevo, Sarajevo, Bosnia and Herzegovina

⁸Prime Institute of Public Health, Peshawar Medical College, Peshawar, Pakistan

⁹Division of Health Security Research, China National Health Development Research Center, National Health Commission, P.R. China, Beijing, China

Corresponding Author:

Xiaoning Hao, PhD

Division of Health Security Research, China National Health Development Research Center, National Health Commission, P.R. China

No. 9 Chegongzhuang Street, Xicheng District

Beijing, 100044

China

Phone: 86 010 88385748

Email: haoxn@nhei.cn

Abstract

Although COVID-19 vaccines are becoming increasingly available, their ability to effectively control and contain the spread of the COVID-19 pandemic is highly contingent on an array of factors. This paper discusses how limitations to vaccine accessibility, issues associated with vaccine side effects, concerns regarding vaccine efficacy, along with the persistent prevalence of vaccine hesitancy among the public, including health care professionals, might impact the potential of COVID-19 vaccines to curb the pandemic. We draw insights from the literature to identify practical solutions that could boost people's adoption of COVID-19 vaccines and their accessibility. We conclude with a discussion on health experts' and government officials' moral and ethical responsibilities to the public, even in light of the urgency to adopt and endorse "the greatest amount of good for the greatest number" utilitarian philosophy in controlling and managing the spread of COVID-19.

(*JMIR Public Health Surveill* 2021;7(8):e26111) doi:[10.2196/26111](https://doi.org/10.2196/26111)

KEYWORDS

COVID-19; coronavirus; COVID-19 vaccine; made in China; vaccine efficacy; vaccine safety; vaccine; China; expectation; safety; efficacy; infectious disease; public health; consequence; public health; standard

Background

How would we as a society remember COVID-19 after the postpandemic realities become the new normal? The destruction it caused, or the construction it motivated? The fears it inflamed, or the hopes it inspired? The transmission it kindled, or the determination it fanned? The sickness it roused, or the human solidarity it helped cement? As COVID-19 is still evolving, it might be difficult to determine how the pandemic might fold

or the contours of the finale. What is clear, though, is that COVID-19 vaccines have already forged and will continue to shape society's collective memories of the great pandemic of the 21st century—the almost-stopped-the-world-go-round global crisis that the COVID-19 pandemic is [1-3]. As of May 24, 2021, COVID-19 has already claimed 167 million infections and 3.46 million deaths [1]—numbers that might only represent a fraction of the true toll, as indicated by investigations led by the World Health Organization (WHO) and other organizations

(eg, the Economist) [2,3]. Although the pandemic has upended the lives and livelihoods of thousands of millions of families and has dragged the world economy into unknown terrain [4], COVID-19 vaccines offer rays of hope that continue to draw people closer to the end of the tunnel [5].

Determined to build some “normalcy,” a global race to develop vaccines that can halt the pandemic has elevated on decades of experience and knowledge on immunization, the most advanced establishment of infrastructure, and an unwavering talent and motivation united to curb the spread of the virus [6,7]. Starting from December 2020, nine months after the WHO first labeled COVID-19 a global pandemic [8], the United Kingdom became the first nation to roll out mass vaccination [9]. Owing to its success in administering shots at the arms, as of May 2021, after months of strict social distancing mandates and within the confines of certain rules, UK residents have once again been able to enjoy shots at pubs indoors [10], with the company of strangers, friends, or one’s inner peace. COVID-19 vaccines, essentially, are the shots of hope people have been anxiously waiting for; when human contact is no longer as contagious as it used to be, hugs, handshakes, and heart-shaped selfies will become possible again across the world. Not to mention the happiness experienced when reuniting with families through nursing home visits, rekindling friendships with face-to-face lectures, and the bittersweet dash to a closing gate for business and leisure travel.

However, it is important to note that COVID-19 vaccines are not equally distributed silver bullets. How well COVID-19 vaccines can help curb the pandemic is contingent upon factors ranging from vaccine accessibility and vaccine efficacy to vaccine hesitancy, particularly in light of uncertainties associated with COVID-19 mutations [11–15]. Not much is discussed about critical issues associated with COVID-19 vaccine access and adoption while sharing positive news on COVID-19 vaccines and during talks about recovery and normalcy. Therefore, in this paper, we examine key factors that shape people’s access to and adoption of COVID-19 vaccines. Furthermore, we draw insights from the literature and aim to identify strategies that could boost people’s adoption of and the availability of COVID-19 vaccines, and ethical considerations associated with these strategies.

Issues Associated With Vaccine Inequity and Accessibility

It is important to note that vaccine availability does not equate to vaccine accessibility [11]. As a result of limitations in vaccine production, although more COVID-19 vaccines will become available in the coming months, not all people will have the same level of access. In the United Kingdom, for instance, older

adults and frontline workers (eg, health care professionals) will be vaccinated first [16]. Simultaneously, in the United States, vaccine distribution policy will be heavily influenced by federal and state policies [17]. In addition to the prioritized distribution of vaccines, trial data availability also affects COVID-19 vaccine accessibility to individual end-users. For example, although expectant mothers are susceptible to COVID-19 [18], most vaccines were not tested on pregnant or lactating women; these individuals will not have access to COVID-19 vaccines until data become available [19]. In other words, although COVID-19 vaccines are available to use, they are not available to use for everyone [19]. This revelation speaks volumes—even though women have been historically ignored and underrepresented in clinical trials [20], it is difficult to contemplate that the same gender inequality could occur amid a pandemic of COVID-19’s scale.

Equally disturbing, evidence further suggests that 90% of people living in 70 poor-income countries across the world will not have access to COVID-19 mass immunization campaigns until 2022 or 2023 [21,22], with the worst estimate pointing to 2024 [23]. On the other hand, high-income countries are hoarding vaccines; by early December 2020, Canada, for instance, had ordered enough doses of COVID-19 vaccines to inoculate each Canadian five times [24]. Overall, as of May 21, 2021, wealthy countries such as the United States, the United Kingdom, Australia, and other nations within the European Union (EU), have collectively ordered approximately 7.8 billion doses of COVID-19 vaccines, whereas only 270,200,000 vaccines are available for low-income countries [23].

COVID-19 vaccines often require advanced infrastructure for storage and delivery (see Table 1) [25]. For instance, to safeguard their potency, Pfizer-BioNTech vaccines are required to be stored and transported between -112°F and 76°F (-80°C to -60°C) [26]; this condition can only be achieved with advanced cold chain systems that are difficult to build and navigate [27,28]. In the United States, due to a failure in storage, a company responsible for vials of the Moderna vaccine, which must be kept cold, spoiled 890 doses destined for older adults in eight nursing home residents in Ohio [29]. At least in the United States, even though several states are not sharing their data, available evidence already shows that vaccine waste is prevalent across states [15]. Considering how higher-income countries face logistical issues using state-of-the-art and high-capacity cold chain systems [28], it is difficult to imagine how low- and middle-income countries will gain access to these vaccines, and how will they deliver these vials to their citizens. Furthermore, pressing issues such as accessibility of glass vials, syringes, and needles may further worsen the COVID-19 accessibility conundrum [28,30].

Table 1. Details of leading COVID-19 vaccines with known efficacy (as of June 2, 2021).

Name	Developer	Country	Type	Efficacy (Dose)	Status	Storage
Convidecia (or Ad5-nCoV)	CanSino	China	Adenovirus	65.28% (single dose)	Approved in China, emergency use in Chile, Hungary, Pakistan, etc	Stable in regular refrigerator for at least 6 months
BBIBP-CorV	Sinopharm	China	Inactivated	86% (2 doses, 3 weeks apart)	Approved in UAE ^a and Bahrain; emergency use in Egypt and Jordan	Stable in regular refrigerator for at least 6 months
— ^b	Sinopharm-Wuhan	China	Inactivated	72.8%	Limited use in China and UAE	Stable in regular refrigerator for at least 6 months
CoronaVac (formerly PiCoVacc)	Sinovac	China	Inactivated	50.38%-78% (2 doses, 2 weeks apart)	Limited use in China, Brazil, etc	Stable in regular refrigerator for at least 6 months
Covaxin (or BBV152 A, B, C)	Bharat Biotech	India	Inactivated	78% (2 doses, 4 weeks apart)	Emergency use in India, Philippines, Zimbabwe, etc	At least a week at room temperature
Sputnik V	Gamaleya	Russia	Adenovirus	91.4% (2 doses, 3 weeks apart)	Early use in Russia	Freezer storage
EpiVacCorona	Vector Institute	Russia	Protein	— (2 doses, 3 weeks apart)	Limited use in Russia and Turkmenistan	Stable in refrigerator for up to 2 years
Vaxzevria (or AZD1222/Covishield)	Oxford-AstraZeneca	UK ^c and Sweden	Adenovirus	60%-90% (2 doses, 4 weeks apart)	Stopped use in Denmark and Norway; emergency use in UK, Lebanon, Canada, etc	Stable in regular refrigerator for at least 6 months
Ad26.COV2.S	Johnson & Johnson	US ^d	Adenovirus	57%-72% (1 dose)	Stopped use in Denmark and Finland; emergency use in US, the European Union, etc	Up to 2 years at -4°F (-20°C) or up to 3 months at 36-46°F (2-8°C)
mRNA-1273	Moderna	US	mRNA	94.5% (2 doses, 4 weeks apart)	Approved in Canada; emergency use in US, UK, etc	Stable in refrigerator for up to 30 days
NVX-CoV2373	Novavax	US	Protein	49.4%-89.3% (2 doses, 3 weeks apart)	—	Stable in regular refrigerator for at least 6 months
Tozinameran or Comirnaty or BNT162b2	Pfizer-BioNTech	US and Germany	mRNA	95% (2 doses, 3 weeks apart)	Approved in Canada, Saudi Arabia, UAE, Bahrain, and Kuwait; emergency use in UK, US, etc	Freezer storage only at -94°F (-70°C)

^aUAE: United Arab Emirates.^bNot available.^cUS: United States.^dUK: United Kingdom.

Issues Associated With COVID-19 Vaccine Safety and Vaccine Hesitancy

Assuming everything goes as planned, COVID-19 vaccine efficacy will still be contingent upon the abilities of individual health facilities to administer their doses. Emerging concerns point to the fact that these institutions often vary in terms of safety protocols, equipment maintenance, and staff training—critical competency criteria that could impact the vaccine administration process, and in turn, vaccine efficacy [17,31]. Competency of vaccine distribution centers also impacts end-user safety. For instance, in the state of West Virginia, 42 people who were scheduled to receive COVID-19 vaccines were

mistakenly injected with an experimental monoclonal antibody treatment that should be administered via an intravenous infusion [32]. In reality, hospitals and medical centers across the world are overstretched and are at a breaking point in addressing the skyrocketing COVID-19 cases [33-36]; many further compound the moral (eg, who should receive COVID-19 vaccines?) and logistical (eg, how to administer these vaccines effectively and safely?) issues associated with vaccine administration.

After severe allergic reaction cases were first reported in the United Kingdom, regulators warned that Pfizer-BioNTech vaccine administration should not be carried out on people with a history of serious allergies [37]. It is worth noting that these

reports occurred prior to the incidents of blood clots reported across the globe, especially in the EU nations [38]. The ever-emerging reports on COVID-19 vaccine side effects are alarming [39-41], as some individuals may not be aware of their allergies or underlying conditions that could expose them to severe vaccine side effects [11]. When they do, vaccine distribution facilities will have to face medical emergencies that they may or may not be capable of tackling. For the Pfizer-BioNTech vaccine trial alone, four volunteers developed Bell palsy or partial facial paralysis during the trial period [42]. For most established vaccines, such as seasonal influenza vaccines, allergic reactions often occur at a low rate estimated at one in a million people [43]; this number is substantially lower number compared to the current known allergic cases associated with COVID-19 vaccines, which is 11.1 per million people for the Pfizer-BioNTech COVID-19 vaccine [44].

In Norway, 23 older adults died shortly after COVID-19 vaccination [45]. Although the investigation is still underway, reports on vaccine side effects, especially if taken out of context, be it by legacy media outlets or conspiracy theory influencers, may further deepen the public's fear, uncertainty, and distrust over COVID-19 vaccines [46]. Not to mention the tsunami of fact-based reports or how fake news may further exacerbate the public's pandemic fatigue, along with potential mental health issues [47-49]. Inevitably, concerns associated with vaccine safety and reports on vaccine side effects may further hinder COVID-19 vaccine adoption [50-53], especially among those who spread unfounded vaccine rumors (eg, vaccine conspirators) or those who are already hesitant about COVID-19 vaccine uptake (eg, vaccine hesitants) [12]. Emerging reports on the impacts of COVID-19 mutations on vaccine efficacy may further compound the situation. Trial data on the Johnson & Johnson vaccines, for instance, show that although the vaccine efficacy is 72% in the United States, it dropped significantly in places where COVID-19 mutations are more prevalent—66% in Latin America and 57% in South Africa [54].

Strategies to Promote COVID-19 Vaccine Accessibility and Adoption

“Reimagining” COVID-19 Vaccine Doses to Improve Vaccine Accessibility

One way to increase vaccine accessibility that many governments are considering is by giving as many people as possible one dose instead of the original and approved two-dose vaccination regimen for fewer people [55]. Britain, for instance, along with other European countries [56], has already delayed administering the scheduled second doses of COVID-19 vaccines on the ground that “vaccinating a greater number of people with a single dose will prevent more deaths and hospitalizations than vaccinating a smaller number with two doses” [55]. In addition to delaying the administration of the second vaccine dose and capitalizing on vaccine overfill (ie, some Pfizer-BioNTech vaccine vials were found to contain a greater amount of the vaccine dose than expected), a surprise that many health care professionals are happy to unveil [57], epidemiologists are also weighing in the option of cutting COVID-19 vaccine doses in half (ie, from 100 mg to 50 mg),

hoping to double the available Moderna vaccine supply in a timely fashion [58].

“Extra” Doses or “Expected” Doses?

Although all the abovementioned measures could help health experts and government officials to capitalize on available vaccine doses, they each come with their own sets of caveats. Among all three measures, the least problematic approach is probably leveraging the vaccine overfill issue. However, even this approach has issues. The first problem lies in the knowledge and experience needed to extract extra doses from the vaccine vials. COVID-19 vaccines are fancy magic delicately packaged in tiny glass vials—they are exceedingly expensive in the way they are designed, developed, delivered, and deployed with care, or lack thereof [11]. The vaccine extraction procedures require medical expertise and special equipment to succeed, which could be an issue considering that hospitals in worst-hit places are often stretched thin. The particular syringe needed for the procedure is in short supply [59]. The second issue is rooted in Pfizer-BioNTech's very business-minded calculations. Not wishing to break its Big Pharma stereotypes, Pfizer will deliver fewer numbers of vaccine vials to account for the difficult-to-extract “extra” doses—Pfizer's contractual agreement with the US government counts doses, rather than vials [60].

In other words, health care professionals in the United States may soon have to extract the “expected” doses from each Pfizer-BioNTech vaccine vial. It is important to incentivize businesses, especially powerful Big Pharmas, amid COVID-19 to contribute to social goods. However, particularly in light of mechanisms such as the Defense Production Act of 1950 [61], it is questionable whether financial incentives are the only approaches governments can use. When all members of the public have to follow the COVID-19 safety measures, such as the United Kingdom's waves of lockdowns, or get fined or jailed, for the greater good, then why are Big Pharma companies such as Pfizer-BioNTech not expected to do the same? Perhaps rather than arguing with governments about wording, dosing, and business bottom lines, Big Pharma companies like Pfizer-BioNTech should focus on producing more COVID-19 vaccines. Overall, it is not a sustainable approach to allow Big Pharma to see lucrative financial benefits in pandemics; societies at large have too many of these already for them to secure their astronomical bonus payments, ranging from the obesity epidemic, HIV epidemic, cancer epidemic, to communicable disease epidemics such as the annual seasonal influenza epidemics.

Utilitarianism Without Consequentialism?

For the approaches that disregard the originally and only clinically tested and approved sets of dosing guidelines, both the problems and solutions may be substantially more challenging to obtain. Essentially, the splitting doses (ie, getting more people to receive one dose of COVID-19 vaccine) and halving doses (ie, getting more people to receive at least some dose of COVID-19 vaccine) methods are a manifestation of “the greatest amount of good for the greatest number” utilitarian philosophy developed by famed scholars such as John Stuart Mill [62]. These approaches have the potential to allow more

people to have access to COVID-19 vaccines without actually improving COVID-19 production rates; however, an important caveat is that there is a lack of data on what might be the health consequences of administering one or halved dose of COVID-19 vaccines, rather than the clinically validated dosing regimen. Data on Pfizer-BioNTech vaccines already shows that the high threshold efficacy for single dose of COVID-19 vaccine is 52%. It could only reach the much-lauded 95% efficacy after the second dose is administered successfully within the prescribed time frame [63]. Available evidence from real-world mass vaccination in Israel further suggests that the actual efficacy of a single dose Pfizer-BioNTech vaccine may have a more disappointing number [64].

It is important to note that the statistics above only address the vaccine efficacy issue rather than other looming issues such as side effects and the interaction between coronavirus and vaccination. Some epidemiologists have already aired their concerns about the potential impacts of inoculating a large portion of the society with the same vaccine in a short time. Collectively, we have yet to figure out how coronavirus might evolve in light of these triggers; will a more potent and powerful variant of SARS-CoV-2 develop that is even more worrisome than the B.1.1.7 mutation first identified in the United Kingdom? It is important to note that some governments have already voiced their concerns over splitting and halving dosing COVID-19 vaccines. Even before data from Israel become available, making it the first country that has managed to vaccinate over 20% of its population and en route to inoculate the entire nation [65], the US Food and Drug Administration, for instance, warned public health officials of the danger associated with tempering with vaccine doses, citing that the idea is not supported by scientific evidence and “may ultimately be counterproductive to public health” [66].

A group of international advisers to the WHO, on the other hand, have recommended public health officials to follow the Pfizer-BioNTech vaccine schedule (ie, two doses given 3–4 weeks apart) rigorously when possible, but they have also suggested that countries with limited supplies of vaccines can consider delaying the second dose for up to 6 weeks [67]. It is important to note that initial evidence on dose splitting and extending intervals between shots is available from the AstraZeneca-Oxford trial [68]. Researchers found comparable efficacies between the two different time frames but disparate efficacies between dosages. Although these insights cannot be directly applied to mRNA vaccines developed by Pfizer-BioNTech and Moderna, they provide preliminary data on the interaction between dosing the vaccine efficacy, which should be further validated or updated by the mentioned effort undertaken by Moderna. Moreover, the Strategic Advisory Group of Experts on Immunization (SAGE), the committee that is tasked to advise WHO when it comes to immunization research and development (eg, COVID-19 vaccine guidelines) [69,70], recommended WHO and all health officials to follow Pfizer-BioNTech dosing and timeframe scheme as the group was reporting the results of their discussion on the approval of the WHO’s emergency use listing of Pfizer-BioNTech vaccines [71]. In extrapolation, then, it can be argued that it is

recommended for officials to follow evidence-based schedules of the corresponding vaccines.

Moral and Ethical Obligations in Public Health Policy-Making

Overall, considering the tsunami of information—fact-based or not—on COVID-19 vaccines, data are urgently needed to shed light on the safety and practicality of changing previously agreed-upon vaccine dosing regimens. Promisingly, a group of scientists in the United States is currently collecting and analyzing data on Moderna vaccines to evaluate the possibility of halving COVID-19 vaccine doses to increase vaccine accessibility [72]. It is essential to digest the fact that devising a different dosing schedule is different from squeezing an additional dose from COVID-19 vaccine vials due to overfilling; the former changes the clinically tested and validated guidelines, whereas the latter simply capitalizes on the fact that some glass vials contain more amount of vaccine.

Although the exact impact of changing the COVID-19 vaccine dosing schedule on personal and public health amid the pandemic is still unclear, what is clear is that governments need to make sure they base their decisions on scientific evidence rather than hopeful assumptions [73]. Yet, baseless assumptions, let alone politics, influencing any decision about COVID-19 vaccines could potentially impact thousands of millions of people’s lives and livelihoods. What is also clear is that, for people who have already received their first dose of COVID-19 vaccines, denying their access to the second dose is a blatant violation of informed consent, the very foundation of medical ethics, a baseline that should not be violated even in a time like the COVID-19 pandemic, particularly in light of dark events ranging from the Nazi’s medical experiments [74], Unit 731 atrocities [75], and the Tuskegee scandal [76]. Obtaining informed consent from potential vaccine receivers has been a tricky task [11], and the violation of informed consent—a contractual trust between individuals and health organizations and governments—may further exacerbate vaccination hurdles for all other immunization efforts.

Conclusions

In this paper, we identified vaccine accessibility and adoption issues that can be collaboratively addressed by both private and public health sectors. Overall, more research is needed to shed light on these tasks, especially factoring in the ever-evolving nature of COVID-19 (eg, mutations) and phenomena such as “pandemic fatigue.” Great hopes have been invested in COVID-19 vaccines. However, it is important to understand that, for COVID-19 vaccines to effectively protect people from the pandemic, issues such as vaccine accessibility, vaccine efficacy, and vaccine hesitancy need to be solved first. In the context of COVID-19, great hopes will almost always mean great expectations—health experts and government officials have a fiduciary and an unwavering duty to the public to make sure they promise what can be delivered and they deliver what is promised. Although even rays of hope can light up the tunnel,

in an environment where distrust is rampant, hope could be easily lost and difficult to rebuild.

Acknowledgments

The authors wish to express their gratitude to the editor and reviewers for their constructive input and insightful feedback along with the kindness they showed to the team throughout the process. This work was supported by the Asia-Pacific Economic Cooperation (APEC) Funded Projects: Building the New Leadership of Infectious Disease Prevention and Control among APEC Economies and the United Nations Development Program (UNDP) South-South Cooperation: Learning from China's Experience to improve the Ability of Response to COVID-19 in Asia and the Pacific Region.

Authors' Contributions

ZS developed the research idea and drafted the manuscript. DMD, AC, XL, DM, SS, JA, and XH reviewed and revised the manuscript. All authors have read and approve the final manuscript.

Conflicts of Interest

None declared.

References

1. COVID-19 Dashboard by the Center for Systems Science and Engineering (CSSE) at Johns Hopkins University (JHU). Johns Hopkins Coronavirus Resource Center. URL: <https://coronavirus.jhu.edu/map.html> [accessed 2021-05-24]
2. The true death toll of COVID-19: Estimating global excess mortality. World Health Organization. URL: <https://www.who.int/data/stories/the-true-death-toll-of-covid-19-estimating-global-excess-mortality> [accessed 2021-05-24]
3. There have been 7m-13m excess deaths worldwide during the pandemic. The Economist. 2021 May 15. URL: <https://www.economist.com/briefing/2021/05/15/there-have-been-7m-13m-excess-deaths-worldwide-during-the-pandemic> [accessed 2021-05-24]
4. Burns D, John M. COVID-19 shook, rattled and rolled the global economy in 2020. Reuters. 2020 Dec 31. URL: <https://www.reuters.com/article/us-global-economy-year-end-graphic/covid-19-shook-rattled-and-rolled-the-global-economy-in-2020-idUSKBN2950GH> [accessed 2021-01-22]
5. Haas EJ, Angulo FJ, McLaughlin JM, Anis E, Singer SR, Khan F, et al. Impact and effectiveness of mRNA BNT162b2 vaccine against SARS-CoV-2 infections and COVID-19 cases, hospitalisations, and deaths following a nationwide vaccination campaign in Israel: an observational study using national surveillance data. The Lancet 2021 May;397(10287):1819-1829. [doi: [10.1016/s0140-6736\(21\)00947-8](https://doi.org/10.1016/s0140-6736(21)00947-8)]
6. Thanh Le T, Andreadakis Z, Kumar A, Gómez Román R, Tollefsen S, Saville M, et al. The COVID-19 vaccine development landscape. Nat Rev Drug Discov 2020 May;19(5):305-306. [doi: [10.1038/d41573-020-00073-5](https://doi.org/10.1038/d41573-020-00073-5)] [Medline: [32273591](https://pubmed.ncbi.nlm.nih.gov/32273591/)]
7. Lurie N, Saville M, Hatchett R, Halton J. Developing Covid-19 vaccines at pandemic speed. N Engl J Med 2020 May 21;382(21):1969-1973. [doi: [10.1056/nejmp2005630](https://doi.org/10.1056/nejmp2005630)]
8. Timeline: WHO's COVID-19 response. World Health Organization. Timeline URL: <https://www.who.int/emergencies/diseases/novel-coronavirus-2019/interactive-timeline> [accessed 2021-05-24]
9. Baraniuk C. Covid-19: How the UK vaccine rollout delivered success, so far. BMJ 2021 Feb 18;372:n421. [doi: [10.1136/bmj.n421](https://doi.org/10.1136/bmj.n421)] [Medline: [33602672](https://pubmed.ncbi.nlm.nih.gov/33602672/)]
10. Restaurants and pubs are reopening inside, but what are the rules? BBC News. 2021 May 21. URL: <https://www.bbc.com/news/business-52977388> [accessed 2021-05-24]
11. Su Z, Wen J, McDonnell D, Goh E, Li X, Šegalo S, et al. Vaccines are not yet a silver bullet: The imperative of continued communication about the importance of COVID-19 safety measures. Brain Behav Immun Health 2021 Mar;12:100204 [FREE Full text] [doi: [10.1016/j.bbih.2021.100204](https://doi.org/10.1016/j.bbih.2021.100204)] [Medline: [33495754](https://pubmed.ncbi.nlm.nih.gov/33495754/)]
12. Su Z, Wen J, Abbas J, McDonnell D, Cheshmehzangi A, Li X, et al. A race for a better understanding of COVID-19 vaccine non-adopters. Brain Behav Immun Health 2020 Dec;9:100159 [FREE Full text] [doi: [10.1016/j.bbih.2020.100159](https://doi.org/10.1016/j.bbih.2020.100159)] [Medline: [33052327](https://pubmed.ncbi.nlm.nih.gov/33052327/)]
13. Callaway E. Fast-spreading COVID variant can elude immune responses. Nature 2021 Jan;589(7843):500-501. [doi: [10.1038/d41586-021-00121-z](https://doi.org/10.1038/d41586-021-00121-z)] [Medline: [33479534](https://pubmed.ncbi.nlm.nih.gov/33479534/)]
14. Freeman D, Loe BS, Chadwick A, Vaccari C, Waite F, Rosebrock L, et al. COVID-19 vaccine hesitancy in the UK: the Oxford coronavirus explanations, attitudes, and narratives survey (Oceans) II. Psychol Med 2020 Dec 11:1-15 [FREE Full text] [doi: [10.1017/S0033291720005188](https://doi.org/10.1017/S0033291720005188)] [Medline: [33305716](https://pubmed.ncbi.nlm.nih.gov/33305716/)]
15. Gabrielson R, Chen C, Simon M. How Many Vaccine Shots Go to Waste? Several States Aren't Counting. ProPublica. 2021 Jan. URL: <https://www.propublica.org/article/covid-vaccine-wastage> [accessed 2021-01-22]
16. Smout A. In COVID-19 milestone for West, Britain starts mass vaccination. Reuters. 2020 Dec 7. URL: <https://www.reuters.com/article/us-health-coronavirus-britain-vaccine-idUSKBN28I01T> [accessed 2020-12-09]

17. COVID-19 Vaccination Program Jurisdiction Operations. Centers for Disease Control and Prevention. 2020 Oct 29. URL: https://www.cdc.gov/vaccines/imz-managers/downloads/COVID-19-Vaccination-Program-Interim_Playbook.pdf [accessed 2020-12-09]
18. Khalil A, Kalafat E, Benlioglu C, O'Brien P, Morris E, Draycott T, et al. SARS-CoV-2 infection in pregnancy: A systematic review and meta-analysis of clinical features and pregnancy outcomes. *EClinicalMedicine* 2020 Aug;25:100446 [FREE Full text] [doi: [10.1016/j.eclinm.2020.100446](https://doi.org/10.1016/j.eclinm.2020.100446)] [Medline: [32838230](https://pubmed.ncbi.nlm.nih.gov/32838230/)]
19. Honderich H. Will pregnant women receive the Covid-19 vaccine? It depends. BBC News. 2020 Dec 22. URL: <https://www.bbc.com/news/world-us-canada-55340244> [accessed 2020-12-24]
20. Vitale C, Fini M, Spoletini I, Lainscak M, Seferovic P, Rosano GM. Under-representation of elderly and women in clinical trials. *Int J Cardiol* 2017 May 01;232:216-221. [doi: [10.1016/j.ijcard.2017.01.018](https://doi.org/10.1016/j.ijcard.2017.01.018)] [Medline: [28111054](https://pubmed.ncbi.nlm.nih.gov/28111054/)]
21. Rich countries hoarding Covid vaccines, says People's Vaccine Alliance. BBC News. 2020 Dec 9. URL: <https://www.bbc.co.uk/news/health-55229894> [accessed 2020-12-09]
22. Dyer O. Covid-19: Many poor countries will see almost no vaccine next year, aid groups warn. *BMJ* 2020 Dec 11;371:m4809. [doi: [10.1136/bmj.m4809](https://doi.org/10.1136/bmj.m4809)] [Medline: [33310819](https://pubmed.ncbi.nlm.nih.gov/33310819/)]
23. Interactive COVAX Map. Launch & Scale Speedometer. Durham, NC: Duke Global Health Innovation Center, Duke University URL: <https://launchandscalefaster.org/COVID-19> [accessed 2021-01-06]
24. Doucleff M. How rich countries are 'hoarding' the world's vaccines, in charts. NPR. 2020 Dec 3. URL: <https://www.npr.org/sections/goatsandsoda/2020/12/03/942303736/how-rich-countries-are-hoarding-the-worlds-vaccines-in-charts?t=1607496526830> [accessed 2020-12-09]
25. COVID-19 vaccines. World Health Organization. URL: <https://www.who.int/emergencies/diseases/novel-coronavirus-2019/covid-19-vaccines> [accessed 2020-12-24]
26. Pfizer-BioNTech COVID-19 vaccine. Centers for Disease Control and Prevention – Vaccines and Immunization. URL: <https://www.cdc.gov/vaccines/covid-19/info-by-product/pfizer/index.html> [accessed 2020-12-24]
27. Checklist for storage, handling, and preparation of the Pfizer-BioNTech COVID-19 Vaccine. Pfizer-BioNTech. 2020. URL: <https://www.cvdvaccine-us.com/images/pdf/Checklist.pdf> [accessed 2020-12-24]
28. Developing and delivering covid-19 vaccines around the world. World Trade Organization. URL: https://www.wto.org/english/tratop_e/covid19_e/vaccine_report_e.pdf [accessed 2020-12-24]
29. Wright W. A company vaccinating Ohio nursing-home residents lets 890 doses go bad. *The New York Times*. 2021 Feb 1. URL: <https://www.nytimes.com/live/2021/01/20/world/covid-19-coronavirus/a-company-vaccinating-ohio-nursing-home-residents-lets-890-doses-go-bad> [accessed 2021-05-31]
30. UNICEF to stockpile over half a billion syringes by year end, as part of efforts to prepare for eventual COVID-19 vaccinations. United Nations International Children's Emergency Fund. 2020 Oct 19. URL: <https://www.unicef.org/press-releases/unicef-stockpile-over-half-billion-syringes-year-end-part-efforts-prepare-eventual> [accessed 2020-12-09]
31. COVID-19 vaccination: Governance, handling and preparation of vaccines in Hospital Hubs and Vaccination Centres. National Health Service. 2020 Dec 4. URL: <https://www.england.nhs.uk/coronavirus/wp-content/uploads/sites/52/2020/12/C0926-COVID-19-vaccination-Governance-handling-and-preparation-of-vaccines-in-Hospital-Hubs-and-Vaccination-Ce.pdf> [accessed 2020-12-09]
32. Wolfe L. 42 people in West Virginia are mistakenly given a virus treatment instead of the vaccine. *The New York Times*. 2020 Dec 31. URL: <https://www.nytimes.com/2020/12/31/us/west-virginia-covid-vaccine-regeneron.html> [accessed 2021-01-06]
33. Bernstein S. No intensive care beds for most Californians as COVID-19 surges. Reuters. 2020 Dec 21. URL: <https://www.reuters.com/article/us-health-coronavirus-usa-california/no-intensive-care-beds-for-most-californians-as-covid-19-surges-idUSKBN28V2R6> [accessed 2021-01-22]
34. Covid-19: Brazil hospitals 'run out of oxygen' for virus patients. BBC News. 2021 Jan 15. URL: <https://www.bbc.com/news/world-latin-america-55670318> [accessed 2021-01-22]
35. Ranzani OT, Bastos LSL, Gelli JGM, Marchesi JF, Baião F, Hamacher S, et al. Characterisation of the first 250 000 hospital admissions for COVID-19 in Brazil: a retrospective analysis of nationwide data. *The Lancet Respiratory Medicine* 2021 Apr;9(4):407-418. [doi: [10.1016/s2213-2600\(20\)30560-9](https://doi.org/10.1016/s2213-2600(20)30560-9)]
36. El-Naggar M, Al-Hlou Y. Egypt denied an oxygen failure killed COVID patients. We found that it did. *The New York Times*. 2021 Jan 18. URL: <https://www.nytimes.com/2021/01/18/world/middleeast/egypt-hospital-oxygen-covid.html> [accessed 2021-01-22]
37. Trigg N, Schraer R. Covid-19 vaccine: Allergy warning over new jab. BBC News. 2020 Dec 9. URL: <https://www.bbc.com/news/health-55244122> [accessed 2020-12-24]
38. Wise J. Covid-19: European countries suspend use of Oxford-AstraZeneca vaccine after reports of blood clots. *BMJ* 2021 Mar 11;372:n699. [doi: [10.1136/bmj.n699](https://doi.org/10.1136/bmj.n699)] [Medline: [33707182](https://pubmed.ncbi.nlm.nih.gov/33707182/)]
39. Klimek L, Novak N, Hamelmann E, Werfel T, Wagenmann M, Taube C, et al. Severe allergic reactions after COVID-19 vaccination with the Pfizer/BioNTech vaccine in Great Britain and USA. *Allergo J Int* 2021;30(2):51-55 [FREE Full text] [doi: [10.1007/s40629-020-00160-4](https://doi.org/10.1007/s40629-020-00160-4)] [Medline: [33643776](https://pubmed.ncbi.nlm.nih.gov/33643776/)]

40. CDC COVID-19 Response Team, FoodDrug Administration. Allergic Reactions Including Anaphylaxis After Receipt of the First Dose of Moderna COVID-19 Vaccine - United States, December 21, 2020-January 10, 2021. *MMWR Morb Mortal Wkly Rep* 2021 Jan 29;70(4):125-129 [FREE Full text] [doi: [10.15585/mmwr.mm7004e1](https://doi.org/10.15585/mmwr.mm7004e1)] [Medline: [33507892](https://pubmed.ncbi.nlm.nih.gov/33507892/)]
41. Cirillo N. Reported orofacial adverse effects of COVID-19 vaccines: The knowns and the unknowns. *J Oral Pathol Med* 2021 May;50(4):424-427 [FREE Full text] [doi: [10.1111/jop.13165](https://doi.org/10.1111/jop.13165)] [Medline: [33527524](https://pubmed.ncbi.nlm.nih.gov/33527524/)]
42. Interim clinical considerations for use of mRNA COVID-19 vaccines currently authorized in the United States. Centers for Disease Control and Prevention. URL: <https://www.cdc.gov/vaccines/covid-19/info-by-product/clinical-considerations.html> [accessed 2021-01-06]
43. McNeil MM, DeStefano F. Vaccine-associated hypersensitivity. *J Allergy Clin Immunol* 2018 Feb;141(2):463-472 [FREE Full text] [doi: [10.1016/j.jaci.2017.12.971](https://doi.org/10.1016/j.jaci.2017.12.971)] [Medline: [29413255](https://pubmed.ncbi.nlm.nih.gov/29413255/)]
44. CDC COVID-19 Response Team, FoodDrug Administration. Allergic reactions including anaphylaxis after receipt of the first dose of Pfizer-BioNTech COVID-19 vaccine - United States, December 14-23, 2020. *MMWR Morb Mortal Wkly Rep* 2021 Jan 15;70(2):46-51 [FREE Full text] [doi: [10.15585/mmwr.mm7002e1](https://doi.org/10.15585/mmwr.mm7002e1)] [Medline: [33444297](https://pubmed.ncbi.nlm.nih.gov/33444297/)]
45. Torjesen I. Covid-19: Norway investigates 23 deaths in frail elderly patients after vaccination. *BMJ* 2021 Jan 15;372:n149. [doi: [10.1136/bmj.n149](https://doi.org/10.1136/bmj.n149)] [Medline: [33451975](https://pubmed.ncbi.nlm.nih.gov/33451975/)]
46. Szabo L. Anti-vaccine activists peddle theories that COVID shots are deadly, undermining vaccination. Kaiser Family Foundation. 2021 Jan 25. URL: <https://khn.org/news/article/anti-vaccine-activists-peddle-theories-that-covid-shots-are-deadly-undermining-vaccination/> [accessed 2021-01-26]
47. Teixeira da Silva J. Corona exhaustion (CORONEX): COVID-19-induced exhaustion grinding down humanity. *Current Research in Behavioral Sciences* 2021 Nov;2:100014 [FREE Full text] [doi: [10.1016/j.crbeha.2021.100014](https://doi.org/10.1016/j.crbeha.2021.100014)]
48. Reicher S, Drury J. Pandemic fatigue? How adherence to covid-19 regulations has been misrepresented and why it matters. *BMJ* 2021 Jan 18;372:n137. [doi: [10.1136/bmj.n137](https://doi.org/10.1136/bmj.n137)] [Medline: [33461963](https://pubmed.ncbi.nlm.nih.gov/33461963/)]
49. Su Z, McDonnell D, Wen J, Kozak M, Abbas J, Šegalo S, et al. Mental health consequences of COVID-19 media coverage: the need for effective crisis communication practices. *Global Health* 2021 Jan 05;17(1):4 [FREE Full text] [doi: [10.1186/s12992-020-00654-4](https://doi.org/10.1186/s12992-020-00654-4)] [Medline: [33402169](https://pubmed.ncbi.nlm.nih.gov/33402169/)]
50. Paltiel AD, Schwartz JL, Zheng A, Walensky RP. Clinical outcomes of a COVID-19 vaccine: Implementation over efficacy. *Health Aff (Millwood)* 2021 Jan;40(1):42-52. [doi: [10.1377/hlthaff.2020.02054](https://doi.org/10.1377/hlthaff.2020.02054)] [Medline: [33211536](https://pubmed.ncbi.nlm.nih.gov/33211536/)]
51. de Figueiredo A, Simas C, Karafillakis E, Paterson P, Larson HJ. Mapping global trends in vaccine confidence and investigating barriers to vaccine uptake: a large-scale retrospective temporal modelling study. *The Lancet* 2020 Sep;396(10255):898-908. [doi: [10.1016/s0140-6736\(20\)31558-0](https://doi.org/10.1016/s0140-6736(20)31558-0)]
52. COCONEL Group. A future vaccination campaign against COVID-19 at risk of vaccine hesitancy and politicisation. *Lancet Infect Dis* 2020 Jul;20(7):769-770 [FREE Full text] [doi: [10.1016/S1473-3099\(20\)30426-6](https://doi.org/10.1016/S1473-3099(20)30426-6)] [Medline: [32445713](https://pubmed.ncbi.nlm.nih.gov/32445713/)]
53. Detoc M, Bruel S, Frappe P, Tardy B, Botelho-Nevers E, Gagneux-Brunon A. Intention to participate in a COVID-19 vaccine clinical trial and to get vaccinated against COVID-19 in France during the pandemic. *Vaccine* 2020 Oct 21;38(45):7002-7006 [FREE Full text] [doi: [10.1016/j.vaccine.2020.09.041](https://doi.org/10.1016/j.vaccine.2020.09.041)] [Medline: [32988688](https://pubmed.ncbi.nlm.nih.gov/32988688/)]
54. Loftus P. Wall Street Journal. 2021 Jan 29. URL: <https://www.wsj.com/articles/j-j-covid-19-vaccine-was-66-effective-in-late-stage-study-11611925201?mod=djemalertNEWS> [accessed 2021-01-29]
55. McKie R. Why is Britain delaying second doses of Covid vaccines? The Guardian. 2021 Jan 3. URL: <https://www.theguardian.com/world/2021/jan/03/why-britain-delaying-second-doses-covid-19-vaccines> [accessed 2021-05-31]
56. Wu K, Robbins R. In Europe, more countries delay second vaccine doses or mull plans to do so. The New York Times.: New York Times; 2021 Jan 4. URL: <https://www.nytimes.com/2021/01/04/world/second-covid-vaccine-delay.html> [accessed 2021-01-06]
57. Thomas K. Hospitals Discover a Surprise in Their Vaccine Deliveries: Extra Doses. The New York Times. 2020 Dec 16. URL: <https://www.nytimes.com/2020/12/16/health/Covid-Pfizer-vaccine-extra-doses.html> [accessed 2021-01-06]
58. Pierson B. U.S. may cut some Moderna vaccine doses in half to speed rollout, official says. Reuters. 2021 Jan 3. URL: <https://www.reuters.com/article/us-health-coronavirus-usa-moderna/u-s-may-cut-some-moderna-vaccine-doses-in-half-to-speed-rollout-official-says-idUSKBN2980NW> [accessed 2021-01-06]
59. Rowland C. Biden wants to squeeze an extra shot of vaccine out of every Pfizer vial. It won't be easy. The Washington Post. 2021 Jan 22. URL: <https://www.washingtonpost.com/business/2021/01/22/pfizer-vaccine-doses-syringes/> [accessed 2021-01-24]
60. Weiland N, Thomas K, LaFraniere S. Pfizer will ship fewer vaccine vials to account for extra doses. The New York Times. 2021 Jan 22. URL: <https://www.nytimes.com/2021/01/22/health/pfizer-vaccine.html> [accessed 2021-01-23]
61. Defense Production Act. FEMA - U.S. Department of Homeland Security. URL: <https://www.fema.gov/disasters/defense-production-act> [accessed 2021-01-24]
62. Mill JS. Utilitarianism. London: Longmans, Green and Company; 1879. URL: <https://www.worldcat.org/title/utilitarianism/oclc/5321953> [accessed 2021-05-31]
63. Mahase E. Covid-19: Pfizer vaccine efficacy was 52% after first dose and 95% after second dose, paper shows. *BMJ* 2020 Dec 11;371:m4826. [doi: [10.1136/bmj.m4826](https://doi.org/10.1136/bmj.m4826)] [Medline: [33310706](https://pubmed.ncbi.nlm.nih.gov/33310706/)]

64. Mahase E. Covid-19: Reports from Israel suggest one dose of Pfizer vaccine could be less effective than expected. *BMJ* 2021 Jan 22;372:n217. [doi: [10.1136/bmj.n217](https://doi.org/10.1136/bmj.n217)] [Medline: [33483332](https://pubmed.ncbi.nlm.nih.gov/33483332/)]
65. Zion IB. Israel trades Pfizer doses for medical data in vaccine blitz. Associated Press. 2021 Jan 18. URL: <https://apnews.com/article/international-news-israel-coronavirus-vaccine-coronavirus-pandemic-benjamin-netanyahu-b30f9af2139e64794ce66c6c9b367b7b> [accessed 2021-01-24]
66. FDA statement on following the authorized dosing schedules for covid-19 vaccines. U.S. Food and Drug Administration. 2021 Jan 4. URL: <https://tinyurl.com/3cz3sx7x> [accessed 2021-01-06]
67. Media briefing on COVID-19. YouTube.: World Health Organization; 2021 Jan 5. URL: https://www.youtube.com/watch?v=57S_hkmUpHw&feature=emb_title [accessed 2021-01-06]
68. Folegatti PM, Ewer KJ, Aley PK, Angus B, Becker S, Belij-Rammerstorfer S, Oxford COVID Vaccine Trial Group. Safety and immunogenicity of the ChAdOx1 nCoV-19 vaccine against SARS-CoV-2: a preliminary report of a phase 1/2, single-blind, randomised controlled trial. *Lancet* 2020 Aug 15;396(10249):467-478 [FREE Full text] [doi: [10.1016/S0140-6736\(20\)31604-4](https://doi.org/10.1016/S0140-6736(20)31604-4)] [Medline: [32702298](https://pubmed.ncbi.nlm.nih.gov/32702298/)]
69. Strategic Advisory Group of Experts on Immunization (SAGE). World Health Organization. URL: <https://www.who.int/groups/strategic-advisory-group-of-experts-on-immunization> [accessed 2021-01-21]
70. SAGE Working Group on Covid-19 vaccines (established June 2020). World Health Organization. URL: https://www.who.int/immunization/sage/sage_wg_covid-19/en/ [accessed 2021-01-21]
71. Farge E, Revill J. WHO recommends two doses of Pfizer COVID-19 vaccine within 21-28 days. Reuters. 2021. URL: <https://www.reuters.com/article/us-health-coronavirus-who/who-recommends-two-doses-of-pfizer-covid-19-vaccine-within-21-28-days-idUSKBN29A25G> [accessed 2021-01-21]
72. Stolberg S, LaFraniere S. Warning of shortages, researchers look to stretch vaccine supply. *The New York Times*. 2021 Jan 5. URL: <https://www.nytimes.com/2021/01/05/us/politics/coronavirus-vaccine-supply.html> [accessed 2021-01-06]
73. Su Z, McDonnell D, Ahmad J. The need for a disaster readiness mindset: A key lesson from the coronavirus disease 2019 (COVID-19) pandemic. *Infect Control Hosp Epidemiol* 2021 Jan 25:1-2 [FREE Full text] [doi: [10.1017/ice.2021.26](https://doi.org/10.1017/ice.2021.26)] [Medline: [33487209](https://pubmed.ncbi.nlm.nih.gov/33487209/)]
74. Mellanby K. Medical experiments on human beings in concentration camps in Nazi Germany. *Br Med J* 1947 Jan 25;1(4490):148-150 [FREE Full text] [doi: [10.1136/bmj.1.4490.148](https://doi.org/10.1136/bmj.1.4490.148)] [Medline: [20244692](https://pubmed.ncbi.nlm.nih.gov/20244692/)]
75. Su Z, McDonnell D, Cheshmehzangi A, Abbas J, Li X, Cai Y. The promise and perils of Unit 731 data to advance COVID-19 research. *BMJ Glob Health* 2021 May 20;6(5):e004772 [FREE Full text] [doi: [10.1136/bmjgh-2020-004772](https://doi.org/10.1136/bmjgh-2020-004772)] [Medline: [34016575](https://pubmed.ncbi.nlm.nih.gov/34016575/)]
76. The Tuskegee timeline. Centers for Disease Control and Prevention. URL: <https://www.cdc.gov/tuskegee/timeline.htm> [accessed 2021-01-21]

Abbreviations

EU: European Union

SAGE: Strategic Advisory Group of Experts on Immunization

WHO: World Health Organization

Edited by T Sanchez; submitted 27.11.20; peer-reviewed by LT Sen, M Das, F Jijrees; comments to author 23.12.20; revised version received 24.12.20; accepted 01.02.21; published 04.08.21.

Please cite as:

Su Z, McDonnell D, Cheshmehzangi A, Li X, Maestro D, Šegalo S, Ahmad J, Hao X

With Great Hopes Come Great Expectations: Access and Adoption Issues Associated With COVID-19 Vaccines

JMIR Public Health Surveill 2021;7(8):e26111

URL: <https://publichealth.jmir.org/2021/8/e26111>

doi: [10.2196/26111](https://doi.org/10.2196/26111)

PMID: [33560997](https://pubmed.ncbi.nlm.nih.gov/33560997/)

©Zhaohui Su, Dean McDonnell, Ali Cheshmehzangi, Xiaoshan Li, Daniel Maestro, Sabina Šegalo, Junaid Ahmad, Xiaoning Hao. Originally published in *JMIR Public Health and Surveillance* (<https://publichealth.jmir.org>), 04.08.2021. This is an open-access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in *JMIR Public Health and Surveillance*, is properly cited. The complete bibliographic information, a link to the original publication on <https://publichealth.jmir.org>, as well as this copyright and license information must be included.

Publisher:
JMIR Publications
130 Queens Quay East.
Toronto, ON, M5A 3Y5
Phone: (+1) 416-583-2040
Email: support@jmir.org

<https://www.jmirpublications.com/>