



Social media sensors as early signals of influenza outbreaks at scale

David Martín-Corral^{1,6} , Manuel García-Herranz³, Manuel Cebrian⁴ and Esteban Moro^{1,2,5*}

*Correspondence:

esteban.moroegido@gmail.com

¹Department of Mathematics and GIS, Universidad Carlos III de Madrid, Leganes, 28911, Spain

²Institute for Data, Systems, and Society, Massachusetts Institute of Technology, Cambridge, MA 02139, USA

Full list of author information is available at the end of the article

Abstract

Detecting early signals of an outbreak in a viral process is challenging due to its exponential nature, yet crucial given the benefits to public health it can provide. If available, the network structure where infection happens can provide rich information about the very early stages of viral outbreaks. For example, more central nodes have been used as social network sensors in biological or informational diffusion processes to detect early contagious outbreaks. We aim to combine both approaches to detect early signals of a biological viral process (influenza-like illness, ILI), using its informational epidemic coverage in public social media. We use a large social media dataset covering three years in a country. We demonstrate that it is possible to use highly central users on social media, more precisely high out-degree users from Twitter, as sensors to detect the early signals of ILI outbreaks in the physical world without monitoring the whole population. We also investigate other behavioral and content features that distinguish those early sensors in social media beyond centrality. While high centrality on Twitter is the most distinctive feature of sensors, they are more likely to talk about local news, language, politics, or government than the rest of the users. Our new approach could detect a better and smaller set of social sensors for epidemic outbreaks and is more operationally efficient and privacy respectful than previous ones, not requiring the collection of vast amounts of data.

Keywords: Computational epidemiology; Social networks; Informational epidemics; Biological epidemics

1 Introduction

For many viral diseases, the early detection of when and where an outbreak will appear is critical. Public administrations responsible for public health management face public health risks such as the Avian flu [1], Zika [2], SARS [3, 4], Ebola [5, 6] or the latest SARS-COV-2 [7, 8] that can cause millions of deaths in a short period of time at global scale [9]. Traditional health surveillance systems require monitoring and detecting symptoms or case incidence in populations. However, their precision sometimes needs to be improved by the size and delayed testing methods on those populations. Combining those data sources with others about people's mobility, the spatial spreading structure of the disease, and even other data sources seem like a promising venue to establish appropriate warning models in the early epidemic stage [10]. Novel data streams like related web

© The Author(s) 2024. **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

search queries and web visits [11–15], weather data [16] or monitoring multiple digital traces at the same time [10] have proven to be complementary and even advantageous to traditional health monitoring systems. In the same way, social media traces have been demonstrated to be a good proxy for digital epidemiological forecasting models of ILI [17–19]. Online user activity exhibits some benefits like broader spatial and demographic reach or monitoring populations that have no easy access to health services [15].

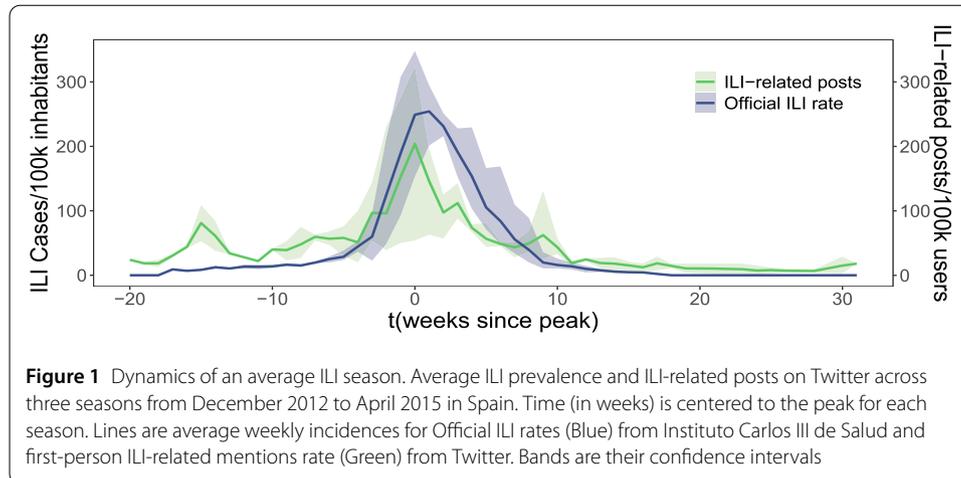
Since some viruses are transmitted by contact on face-to-face social networks, epidemiological methods that exploit the network structure are more effective in detecting, monitoring, and forecasting contagious outbreaks [20, 21], since they allow to anticipate more accurately the transmission dynamics. Furthermore, these methods can help public health decision-makers to enhance the adoption of public health interventions [22] like social distancing, vaccination, or behavior change campaigns, identifying those individuals most likely to get infected and spread an infectious disease or behavior (e.g., super-spreaders), or which places are more likely to be visited by those individuals [23]. This allows more efficient vaccination campaigns [24] when the vaccination of an entire population is not possible or recommended.

The key idea behind using high-connected individuals to monitor epidemic spreading is that they are more likely to be reached by the infection. In general, human social sensing, when carefully selected, can help predict and explain social dynamics better [25–27]. In the absence of complete detailed data about contact networks, simple approaches like the friendship paradox [28] can be used to identify more connected and central individuals (sensors) in the network that can give early signals and anticipate the spreading of information, behavior or disease before it reaches a significant fraction of the population. In particular, the friendship paradox has already been found advantageous to identify sensors for detecting influenza [29–31] or COVID-19 [32]. In social media, a previous study demonstrated the detection of global-scale viral outbreaks of information diffusion [33] by monitoring high-degree users on Twitter.

In this work, we address the question of how we can use sensors for information propagation in online social media to get better early warning signals of a biological epidemic. We hypothesize that social media connectivity and activity are a proxy of social interactions in the real world. Thus, highly-connected users in social media (online sensors) also mirror highly-connected individuals (offline sensors) in the physical contact network. This hypothesis is based on the wealth of literature showing that online networks mimic offline contacts' connections, similarity, and spatial organization [25, 34, 35]. Furthermore, we study if it is possible to identify better social media sensors automatically based on their centrality (degree) and mobility, and content behavior. We found that social media sensors can serve as early signals of the exponential growth of an epidemic several weeks before the peak. The current global pandemic threads make it vital to improve the efficiency of Early Warning Epidemiological Systems (EWES) by using operationally efficient methods to anticipate the exponential growth of a virus in a community, region, or country without compromising the citizens' privacy. Our method provides such a system in a fully privacy-preserving framework because it does not necessitate the collection of users' contact links; instead, it solely relies on the degree metric.

2 Results

We used social media traces obtained from the micro-blogging site Twitter, where we collected more than 250 million tweets from December 2012 to April 2015 on Spain's

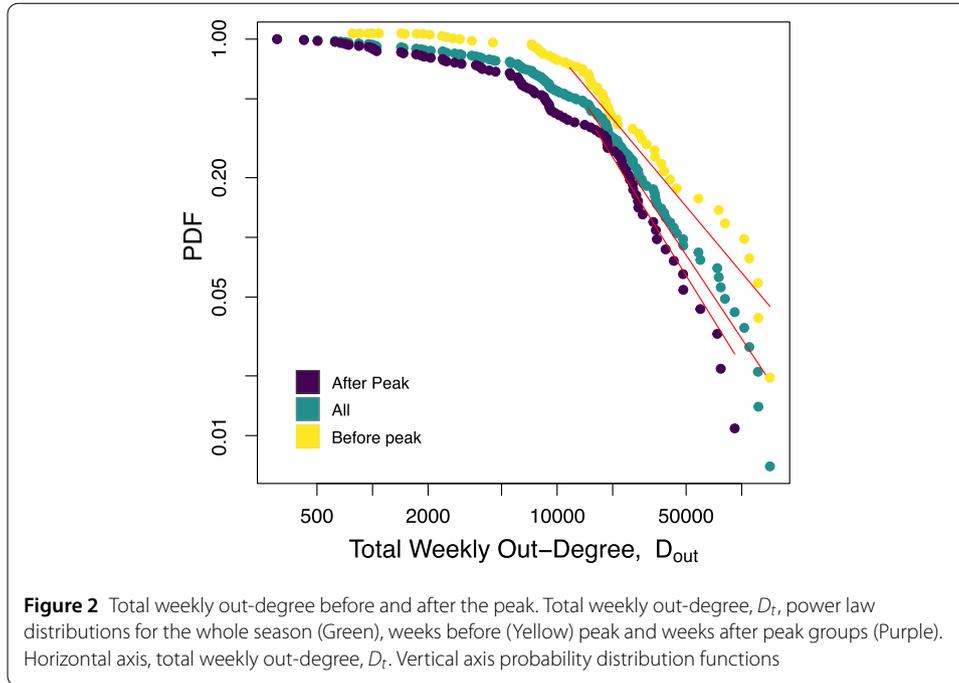


mainland. Using Natural Language processing techniques, we only included first-person ILI-related posts, summing up a population of 19,696 users with at least one first-person ILI-related mention, which comprised a total of 23,975 tweets (Sect. 4 & Additional file 1 Sect. 1 discusses our methodology). We also made use of official ILI cases from the surveillance system for influenza in Spain (ScVGE) [36] managed by the Instituto Carlos III de Salud [36]. This system reported weekly ILI cases in Spain for each province with two weeks of delay in the state of the seasonal flu epidemic based on the current European Union proposal that regulates ILI surveillance [37]. Our dataset of official ILI cases ranges from December 2012 to April 2015 and includes three different seasons of influenza outbreaks in Spain.

Figure 1 shows a generalized ILI season from the average of ILI cases and ILI-related mentions for the three seasons. ILI cases and ILI-related mentions time series have a Pearson correlation of 0.87 (CI [0.79, 0.93] and $p_{\text{value}} < 0.001$). Since different outbreaks happen at different times of the year, we have shifted each influenza outbreak to the time of its peak. We can see that ILI-related mentions precede the official ILI cases at the beginning of the growth stages before the peak. Previous studies have proved this [17–19]. Mentions of the outbreak in social media seem to precede the exponential growth in the total population. ILI-related posts peak at -15 weeks could be related to the start of the cold season and users mixing ILI symptoms with cold symptoms, stating that they are suffering from ILI. We found a similar pattern in Google trends data.

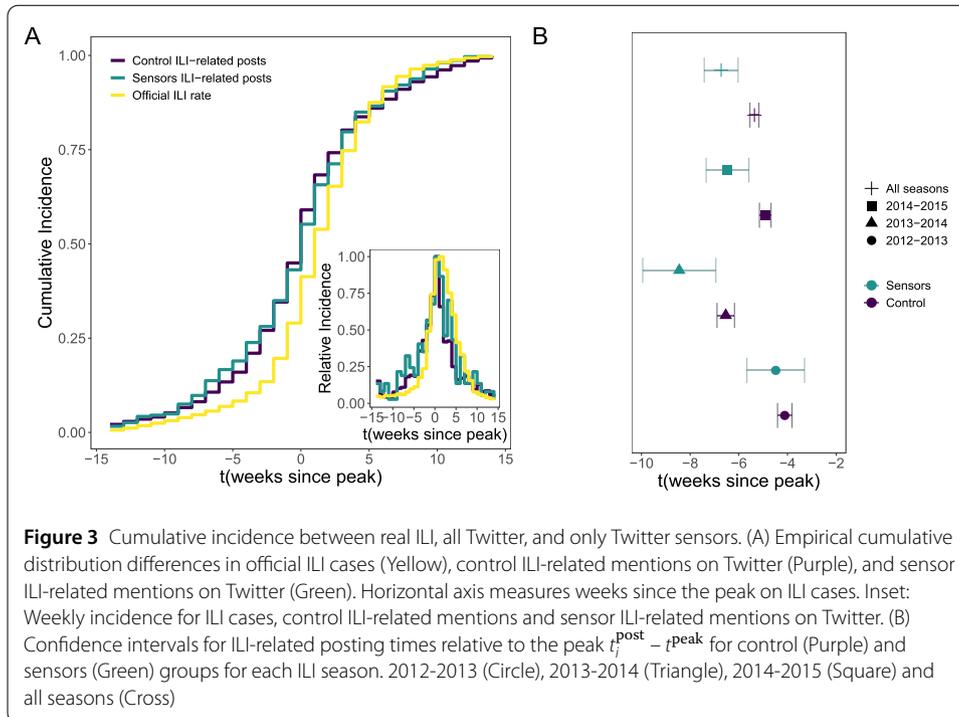
2.1 Validating high-degree individuals as sensors

However, here, we want to go a step further. Can we subset the users posting ILI-related posts to get better earlier signals about the outbreak than monitoring the whole social network platform? Similarly to [29], and [33], high-degree users could be better than the average individual on the platform. To test whether high centrality or degree correlates with early signals, we measure the total weekly out-degree, D_t , of users having social ILI-related mentions before and after the peak. To delineate the periods before and after the peak, we centered the epidemic curves around their respective maximum values for each season. Defining the pre-peak period as $-15 \leq t \leq -1$ and the post-peak period as $1 \leq t \leq 15$. Figure 2 shows distributions for D_t before the peak, after the peak and for the whole season. There is a statistically significant difference in the mean ($p_{\text{value}} < 0.01$). The average



total weekly out-degree is 31,108 (Confidence Interval, CI [21,539.03,40,677.32]) before the peak, while it is only 14,373 (CI [11,202.94,18,455.78]) after the peak. The difference is also present in extreme values. We modelled large values of D_t as power laws with an exponent of 2.56 (CI [2.51,2.62]) for the whole period. For the weeks before the peak, it follows an exponent of 2.10 (CI [1.91,2.29]). Finally, for the weeks after the peak, it follows an exponent of 2.86 (CI [2.48,3.25]). Thus, on the aggregated level, we indeed see that the users in social media that have ILI mentions before the peak have more social connections than after the peak. This result signals the possibility of using high-connected users as potential early sensors. This result is robust against other aggregated degree centrality variables (see Additional file 1, Sect. 2). For selecting sensors, we selected each individual with an out-degree greater than 1000 (see Additional file 1, Sect. 3).

Figure 3.A compares Twitter's cumulative ILI-related mentions of our control and sensor groups against the official ILI-related cases. As we said before, the activity in social media for both the control and sensor groups anticipates the cumulative incidence of ILI cases by one or two weeks. For each user i we define t_i^{post} as the time in which she has an ILI-related post on social media. Figure 3.B shows confidence intervals for ILI-related posting times for each group and ILI season, relative to the peak $t_i^{\text{post}} - t^{\text{peak}}$. For all ILI seasons, the control group has an average ILI-related posting time of $\Delta t_C = \langle t_i^{\text{post}} - t^{\text{peak}} \rangle_{i \in C} = -5.35$ (CI [-5.54, -5.17]) weeks before the peak. The sensor group has an average ILI-related posting time of $\Delta t_S = \langle t_i^{\text{post}} - t^{\text{peak}} \rangle_{i \in S} = -6.72$ (CI [-7.42, -6.02]) weeks before the peak. This yields that sensors are posting on average $\Delta t_S - \Delta t_C = -1.37$ (CI [-2.08, -0.64] and $p_{\text{value}} < 0.01$) weeks before the control group, during the exponential growth phase, between 8 to 4 weeks for all seasons. In more detail, the 2012-2013 season has a $\Delta t_S - \Delta t_C = -0.62$ (CI [-1.58, -0.84] and $p_{\text{value}} > 0.1$), the 2013-2014 season has a $\Delta t_S - \Delta t_C = -2.46$ (CI [-3.45, -0.36] and $p_{\text{value}} < 0.01$) and the 2014-2015 season has a $\Delta t_S - \Delta t_C = -1.54$ (CI [-2.45, -0.63] and $p_{\text{value}} < 0.01$). As we can see, the ILI-related mentions of sensors could



anticipate the epidemic's growth by 1 or 2 weeks with respect to other users in the platform.

2.2 Autoregressive models with sensors and its theoretical validation

To quantify statistically how valid our sensors in social media could be in a potential EWES model, we built an autoregressive model that considered different epidemiological and social media features (see Sect. 4). The models considered different combinations of the total number of weekly ILI cases at time t , I_t , the total weekly out-degree of all users from the social media platform ($D_{T,t}$) that posted ILI-related mentions, and the total weekly out-degree of the subset of those users in the sensor group ($D_{S,t}$). We have also considered different temporal week lags, $t - \delta$, for each variable to test their potential role as early warning signals. As a baseline, we have considered a model that only incorporates the ILI cases and their autoregressive power at $t - 1$. As we see in Table 1, that simple model is already quite accurate in explaining the evolution of the weekly ILI rate. On top of that baseline model, we built four others, including the degree centrality of all users and the sensor group at different lags. For each model, we predict the I_t number of ILI-related cases using the information of the I_{t-1} cases and the total out-degree of all users and sensors with ILI-related mentions at time t and $t - \delta$. We ran all models using a step-wise approach to keep only statistically significant regressors for $\delta = 1, 2, 3, 4, 5, 6$. Due to multicollinearity problems between variables, we also monitor the variance inflation factor (VIF) for each to choose the best δ . Improved: For δ values greater than 1, VIF values remain below 10. However, when considering $\delta = [5, 6]$, the VIF values drop below 5, albeit with a slight reduction in their predictive power.

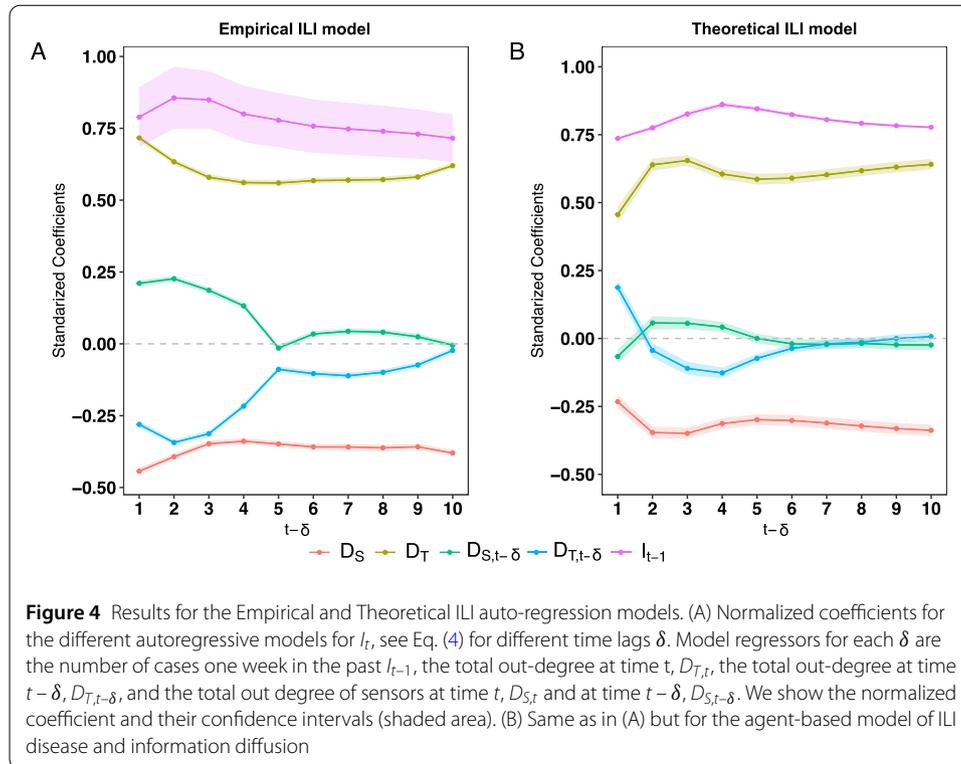
Results in Table 1 and Fig. 4A quantitatively show the importance of social media ILI-related mentions, especially those from the sensor group. As we can see, the predicting power (adjusted R^2) on next week's official ILI rate after incorporating social media men-

Table 1 Empirical ILI regression models. Regression table with normalized beta coefficients for each group of variables, Official (I)LI, (T)witter and (S)ensors, where X_t are weekly ILI related variables for each group. $D_{i,t}$ are weekly total out-degree variables from Twitter and Sensors

	Official weekly ILI rate						
	ILI (1)	T + S $t - 1$ (2)	T + S $t - 2$ (3)	T + S $t - 3$ (4)	T + S $t - 4$ (5)	T + S $t - 5$ (6)	T + S $t - 6$ (7)
I_{t-1}	0.925*** (0.041)	0.789*** (0.047)	0.856*** (0.049)	0.849*** (0.045)	0.800*** (0.044)	0.778*** (0.042)	0.757*** (0.042)
$D_{T,t}$		0.717*** (0.0003)	0.634*** (0.0002)	0.580*** (0.0002)	0.561*** (0.0002)	0.560*** (0.0002)	0.568*** (0.0002)
$D_{T,t-1}$		-0.281*** (0.0003)					
$D_{T,t-2}$			-0.344*** (0.0003)				
$D_{T,t-3}$				-0.313*** (0.0002)			
$D_{T,t-4}$					-0.217*** (0.0002)		
$D_{T,t-5}$						-0.089*** (0.0002)	
$D_{T,t-6}$							-0.104*** (0.0002)
$D_{S,t}$		-0.443*** (0.0004)	-0.393*** (0.0003)	-0.348*** (0.0003)	-0.339*** (0.0003)	-0.349*** (0.0004)	-0.359*** (0.0004)
$D_{S,t-1}$		0.211*** (0.0005)					
$D_{S,t-2}$			0.227*** (0.0004)				
$D_{S,t-3}$				0.186*** (0.0003)			
$D_{S,t-4}$					0.132*** (0.0003)		
$D_{S,t-5}$						-0.015*** (0.0003)	
$D_{S,t-6}$							0.034*** (0.0003)
Constant	4.951 (4.627)	0.000 (4.093)	0.000 (4.071)	0.000 (4.123)	0.000 (4.516)	0.000 (4.771)	0.000 (5.081)
Observations	87	87	86	85	84	83	82
R ²	0.854	0.925	0.932	0.935	0.929	0.927	0.923
Adjusted R ²	0.852	0.920	0.928	0.931	0.924	0.922	0.918
Maximum VIF	NA	15.16	9.36	6.77	5.28	4.59	4.37
Residual Std. Error	34.092 (df = 85)	25.042 (df = 81)	23.951 (df = 80)	23.529 (df = 79)	24.741 (df = 78)	25.166 (df = 77)	25.910 (df = 76)

Note: *p < 0.1; **p < 0.05; ***p < 0.01

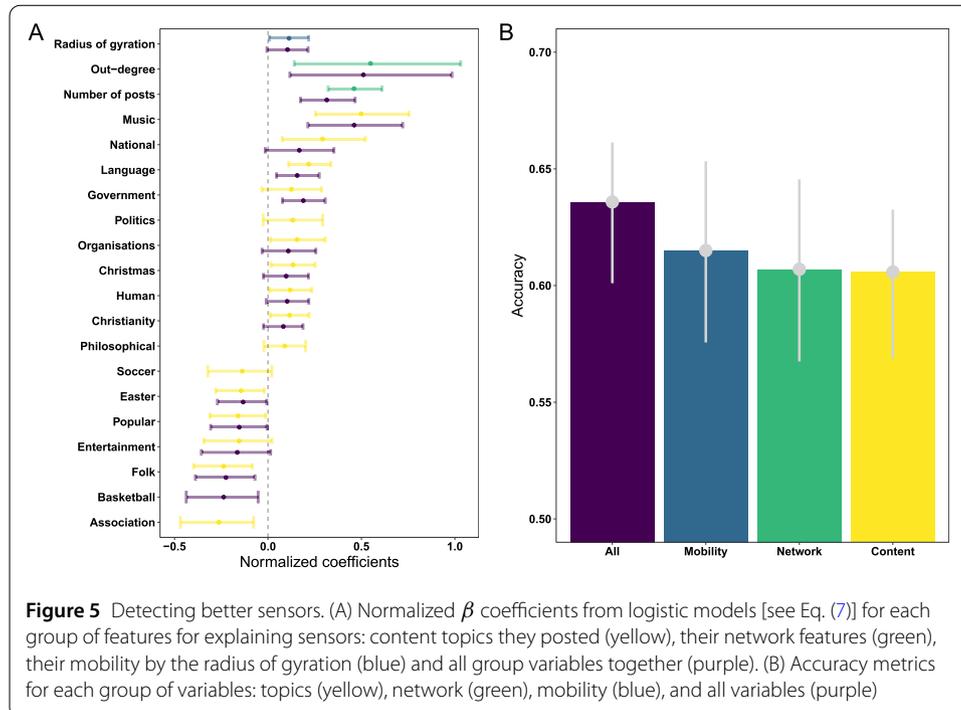
tions increases significantly (and we also reduced collinearity), especially at five- or six-week lags. In all those cases, the total degree of sensors at time T and time $t - \delta$ has a significant regression coefficient and role (in R^2) in the prediction. That is, social sensors can help anticipate official ILI cases five to six weeks before, a result consistent with previous similar analyses of ILI contagious outbreaks in small settings [29] or of information spreading in social media [33]. We also note that the signs of the variables of all users and sensors have different effects. For example, a higher total degree of sensors at times $t - \delta$ predicts more ILI-related cases (positive coefficient) at time t for $\delta > 0$, but a smaller number of cases (negative coefficient) for $\delta = 0$. As we will see below, this apparent con-



tradition comes from the auto-correlation of the time series of ILI-related cases and the total degree of users.

We investigated the predicting power of high-degree sensors in a synthetic model to validate that sensors anticipate ILI cases because social media connectivity mirrors social connections in the real world. Specifically, we built a base agent-based susceptible-infected-recovery (SIR) epidemic spreading on a random network mimicking real (face-to-face) social contacts between people (see Sect. 4 for details about the network and simulations details). Apart from their physical contacts, we also assumed that each person has acted on a social media platform and that the degree in both the real and online networks are correlated moderately. Assuming that agents post on social media when they are infected, we also constructed the time series $\hat{D}_{T,t}$ and $\hat{D}_{S,t}$ for the model and their autoregressive fits as in Table 1. Our results once again show that high-degree agents (sensors) carry some predicting power on the epidemic spreading.

Furthermore, the coefficients for the different models show the same regression structure as the empirical models in 1, see Fig. 4A. We can see that both coefficient structures are nearly the same, including their magnitude and signs. Although this is not direct proof of our hypothesis that the online and offline centrality of real users is similar, it shows that under that assumption, we not only get that the effect of sensors is the same as we found in our empirical analysis, but even the structure of coefficients (magnitude and sign) is similar. These results support the idea that sensors in an informational epidemic that mirrors a biological epidemic are also sensors of a biological epidemic, like ILI, that we can trace on Twitter.



2.3 Identification of sensors beyond out-degree

So far, we have seen that high out-degree users in social media can be early sensors of ILI cases. However, can we identify a better group of sensors beyond high degrees by looking at other traits? Are individuals that signal the epidemic's early stages defined just by their centrality degree, or do they have other behavioral or content traits? To do that, we define a sensor functionally as every user who posts an ILI-related tweet from fifteen weeks to two weeks before the epidemic's peak ($-15 \leq t \leq -2$). On the other hand, a control user was a random user who did not talk about ILI during the same period. (see Sect. 4 contextual features for more details).

To characterize users' content, behavior, and network traits in both groups, we analyzed every tweet they posted 30 days before their first ILI-related tweet (sensors) or a randomly chosen tweet (control). Specifically, we identify three groups of traits for each user. Firstly, we extract the content of each user's tweets and classify them into topics like sports, politics, entertainment and many other categories using the TextRazor classifier (see Sect. 4). Secondly, since our tweets are geolocalized, we extract the mobility features of each user, in particular, the radius of gyration, which measures the size of the area covered while moving around [38]. The radius of gyration could proxy the number of different and diverse people the user is in daily contact. Thus it might serve to estimate potential exposure to infected people [39]. Lastly, we also use their activity (number of posts) and, as before, their out-degree in the social network.

To test how relevant those groups of traits are to define a sensor, we developed a straightforward logistic regression model (see Sect. 4) to classify users into the sensor or control groups using different variables. As we can see in Fig. 5, the accuracy of our models is above the primary level (0.5). While Network and Content groups independently achieve similar accuracies (~ 0.61) than the Mobility group (~ 0.62), we get better accuracy, including all types of traits (~ 0.64). This result signals that even different traits carry complemen-

tary information about who could be sensors in the social media platform. To understand this further, we looked into each trait's (normalized) coefficients in our model. As shown in Fig. 4A, the most crucial variable to predict a user in the sensor group is still the out-degree in the social network, even after controlling for the number of posts. This is important because it shows that our simple method of using high-connected Twitter users as sensors works much better than other traits. We also see a small but significant effect on the radius of gyration and high number of posts, meaning, all things equal, users that move further are more likely to be exposed to the virus, have a higher probability to get infected and sensors. A more pronounced effect is observable for the number of posts. According to our hypothesis asserting that social networks mirror physical networks, individuals with both a high radius of gyration and an elevated number of posts are likely to be highly socially active in the real world. Consequently, such individuals would possess a higher degree and serve as better sensors for detecting early signals of an epidemic. Regarding the content, we see a structure of topics that users in the sensor group are more likely to discuss, like National, Language, Politics, and Government. On the contrary, users that talk about Sports, Popular topics, or Entertainment are less likely to be in the sensor group. This finding could signal and be related to other unobserved user traits like income or educational attainment level, which also are known to be related to the activity in social media [40] and amount of real offline contacts [41].

3 Discussion

Early warning epidemiological systems (EWES) detect outbreaks weeks in advance to help public health decision-makers make more efficient allocations of public resources to avoid or minimize an overflow of contagious in the healthcare system. EWES are undergoing significant investments and changes due to the COVID-19 disruption. However, most of them harvest vast amounts of data and do not exploit the explanatory and predictive power of the network heterogeneity where a disease-informational epidemic is spreading.

In this study, we demonstrated that social media traces, like Twitter, could be used as a source of social-behavioral data to monitor disease-informational epidemics that mirror offline biological contagious disease epidemics, like ILI, by exploiting the network heterogeneity whenever social centrality measures of the network are available. By having a simple centrality metric, such as the out-degree, we can define suitable sensors for the disease-informational epidemic in the network. When aggregated correctly, we can use sensors to feed autoregressive models that could yield signals of an outbreak up to four weeks in advance. Although previous studies showed the advantage of using social network metrics to detect, monitor, and forecast contagious outbreaks [20, 21]. The usage of sensors in a network to detect early signals of an outbreak in a biological disease contagious epidemic [29, 30], or informational epidemics [33]. Furthermore previous studies have used digital traces for predicting epidemics like ILI [11, 12, 42]. However, these studies do not fully leverage the power of network heterogeneity. Consequently, our study stands as the first to integrate the use of social media sensors to forecast real-life epidemics, unlocking new potential in the field by leveraging an indirect metric of a user's network position, such as their degree, for detecting early signals of an epidemic. Our results are based on the hypothesis that social media networks are related to offline contact networks, which has been validated directly in other works [25, 34, 35]. Our empirical and theoretical results show that instead of harvesting large amounts of data and metrics from social networks

[19], we can track and anticipate early outbreaks of a disease-informational epidemic by inexpensively looking at a small set of specific users (sensors).

We also demonstrated that sensors could be profiled and detected automatically from social media raw data by using their topological network properties and based on the content posted by individuals and their mobility patterns. Explicitly, we found that sensors talk more about some topics like National, Politics, and Government and less about Sports and Entertainment. The fact that those topics could also be related to their income, educational attainment [40], but also to other traits like more extroversion personality traits [43] opens the possibility to investigate the potential overlapping reasons why sensors not only are more prone to get infected earlier but also that they would like to post about it on social media. For instance, the Music topic requires further investigation; previous literature suggests individual differences in personality in the way we use and experience music [44], possibly having a social component.

Finally, our method uses the out-degree in the social media platform as a proxy for centrality. Better knowledge of the network structure could yield more optimized methods to detect highly-central users. Our approach exhibits additional limitations. Specifically, our dataset is confined to a particular epidemic within a specific country, covering flu-related mentions from 2012 to 2015 in a given social media platform. The method has not undergone testing across diverse regions, with more recent data, or against global, contemporary epidemics or different social media platforms, such as the COVID-19 pandemic. However, given that our findings rely on the collective behavior of people in social media and the observed relationship between offline and online networks [45, 46], we think that our findings could be extrapolated to other epidemics, regions and social media platforms. We hope our research can help study the role of sensors in other pandemics, specially COVID-19, where more information about real-world offline contact networks exists due to better mobility data [47] or contact tracing applications.

In summary, this study proposes a feasible approach to exploit the network heterogeneity underneath social media sites, like Twitter, to detect more efficiently and earlier outbreaks from a disease-informational epidemic that mirror a biological disease contagious epidemic, like ILI. Furthermore, the sensors approach we used to detect early outbreaks within informational epidemics and biological contagious disease epidemics, but this is the first time in a disease-informational epidemic as we have done in this study. Finally, novel epidemiological systems have been developed for other pathogens such as Zika, SAR, or COVID-19, among others, in addition to influenza, using conventional and non-conventional data sources such as the official public cases, online searches, or health forums. For instance, for the COVID-19 pandemic, some studies used social media traces to try to predict the dynamics of the pandemic [48, 49]. Such approaches, along with our findings about the power of the network structure, could improve the results of their predictions.

Also, health systems and health organizations initiatives, like the Global Outbreak Alert and Response Network (GOARN) [50] from WHO that is composed of 250 technical institutions and networks globally and projects like the Integrated Outbreak Analytics (IOA) [51], Epidemic Intelligence from Open Sources (EIOS) [52], and Epi-Brain [53] that respond to acute public health events. This network is already moving in a double direction of incorporating early signals from Big Data, social sciences techniques and behavioral data into epidemic response systems [54] to control outbreaks and public health emergen-

cies across the globe. Also, syndromic surveillance platforms like InfluenzaNet could ask for Twitter profiles or the number of people an individual interacted with in the last week to reweight the impact of different users in the prediction. Our innovative approach might help detect early outbreaks without having to monitor and harvest data from a whole population, making EWES more accurate in time prediction of an outbreak, more efficient in resources, and more respectful of citizens' data privacy.

4 Methods

4.1 Data collection

We extracted Twitter data through their streaming API [55] that allowed us to collect data programmatically on the Spanish mainland by using a geolocated query in Spanish for minimizing data inconsistencies. The official ILI rate data was extracted through a web crawler built ad-hoc for the web of the Institute Carlos III of Health since there was no access to the raw data from an open data portal or a programmatic interface.

4.2 ILI-related keywords based search and tweets classification

To get ILI-related mentions from users in the social media platform, we first filtered tweets by keeping those that mentioned simple terms like “flu” or other ILI-related words (see Additional file 1). After that, we only kept first-person ILI-related mentions to exclude general or not directly-related posts like ‘The Spanish flu was an unusually deadly influenza pandemic’. This was done using Natural Language Processing methods. We employed a text classifier utilizing the Stochastic Gradient Descent algorithm, implemented through the `scikit-learn` library [56]. We handpicked and labeled a set of 7836 tweets to train our classifier, containing 3918 true positive (first-person) tweets and 3918 true negative tweets. We performed a vectorization on the labeled data using the Term Frequency-Inverse Document Frequency (TF-IDF) procedure in the `scikit-learn` library. Subsequently, we reduced the TF-IDF matrix using the TruncatedSVD procedure, also provided by `scikit-learn`. Finally, we hyper-parameterized the Stochastic Gradient Descent classifier with $\alpha = 0.0001$ and a regularization L2 norm, and applied it to the processed matrix, achieving an accuracy of (~ 0.94) and kappa value of (~ 0.83). We then applied our classifier to identify first-person mentions in the remaining tweets (see Additional file 1 for more details about our pipeline). After this process, we ended up with $N = 19,696$ users and 23,975 tweets classified as first-person ILI-related posts.

4.3 ILI-related post time series

We added up and normalized the number of weekly users mentioning the flu by the total number of users in the system. We followed equation

$$\hat{x}_{\text{users},t} = \frac{x_{\text{ILI Users},t}}{x_{\text{Total Users},t}}, \quad (1)$$

where t is the week. This time series is shown in Fig. 1, together with the prevalence of ILI cases.

4.4 Centrality features

Each tweet at time t has information about the out-degree (followees), $d_{\text{out},i,t}$, and in-degree (followers), $d_{\text{in},i,t}$, for each Twitter user i posting it. We used them as proxies of the

network centrality for each user. Only 5% of users have more than one ILI-related mention and their in and out degrees do not change dramatically so we take $d_{out,i,t} \simeq d_{out,i}$ (similarly for $d_{in,i,t}$) with t being their first (or most of the times only) ILI-related mention. We tested out several aggregated centrality features for the selection of sensors. We found that the weekly total out-degree was the best centrality metric to apply with a Pearson correlation of 0.91 (CI [0.87, 0.93] and $p_{value} < 0.001$), compare against the weekly total in-degree with a Pearson correlation of 0.77 (CI [0.68, 0.82] and $p_{value} < 0.001$). We also calculated the weekly total, mean, median, maximum and minimum out-degree of individuals before and after the peak making first-person ILI-related mentions to test if other out-degree statistics had more explanatory power. The centrality metrics are solely based on twitter users metrics and we do not build the real network between users. See Additional file 1, Sect. 2 for further details.

The weekly total out-degree is defined by

$$D_{T,t} = \sum_{i \in \Omega_t} d_{out,i}, \quad (2)$$

where Ω_t is the set of users that made an ILI-related mention at week t .

Sensors are selected as the group of users with $d_{out,i} > 1000$. For that group, we also define the time series of their centrality as

$$D_{S,t} = \sum_{i \in \Omega_t^*} d_{out,i}, \quad (3)$$

where Ω_t^* is the set of users in the sensor group that made an ILI-related mention at week t .

4.5 Linear autoregressive model

The following equation represents a linear autoregressive model for explaining and now-casting the dependent variable, I_t , being the Official ILI rate for each week. $D_{T,t}$ are total weekly out-degree for the whole twitter population, and $D_{S,t}$ are total weekly out-degree for the whole sensor population. We followed

$$I_t = \beta_0 + \beta_1 I_{t-1} + \sum_{\delta \geq 0} (\alpha_\delta D_{T,t-\delta} + \gamma_\delta D_{S,t-\delta}) + \epsilon_t. \quad (4)$$

4.6 Agent-based model of ILI disease and information diffusion

To understand our empirical findings, we compare them with the simulations of epidemic spreading on a physical and online network through an agent-based model (ABM). We model the offline (physical) contacts using a random heavy-tailed network. Specifically, we created a synthetic population of $N = 150k$ agents which are connected through a scale-free network with degree distribution $P(k) \sim k^{-3}$ obtained through the Barabasi-Albert model. [57]. The network was built using the R package `igraph` [58].

At the same time, we supposed that each agent participates in a social media platform. We hypothesize that the online degree of the agents is related to the offline degree in the complex network. To account for some variability, we assumed that the degree in the social media platform was modified by a random uniform distributed number (See Additional

file 1 Sect. 4 for more details). Thus, the degree in the social media platform is given by $d_{out,i}^{Twitter} = d_{out,i}^{Offline}(1 + \nu_i)$, where ν_i is a random number uniformly distributed between 0 and 1. This way we account for potential variability between offline and online degrees.

We simulate the ILI spreading using a simple Susceptible-Infected-Recovered (SIR) epidemic model. In particular, at each time-step t , the infectious (I) agents can transmit the disease to their susceptible (S) neighbors in the contact network with probability β . If the transmission is successful, the susceptible node will move to the (I) state. An individual will move independently to the recovery (R) state with a probability α . We initialized the model with two initial infected seeds. After getting infected, we assumed that the agent immediately posted an ILI-related tweet on the social media platform. In our model, we considered a user to be sensors if she has an out-degree in the platform higher than four times the average degree in the Barabasi-Albert model. We also calibrated the time unit in this model so that the epidemic curves have a similar time scale as the real ILI rate (See Additional file 1 Sect. 4 for further details on the simulation's parameters).

4.7 User traits

To characterize the different traits of Twitter users, we analyzed each user's tweets during a time window of 30 days before the initial event. For the sensor group, we selected individuals with an out-degree $d_{out,i} \geq 1000$ and that made at least an ILI-related mention during the weeks $-15 \leq t \leq -2$ before the peak of the epidemic. The initial event is their first post with the ILI-related mention. For the control group, we picked individuals that made an ILI-related mention after the $-15 \leq t \leq -2$, then we picked a random post of them as an initial event in weeks $-15 \leq t \leq -2$, before the peak of the epidemic. Using that 30 days period, we computed different Mobility, Content, and Network traits to characterize each user.

4.7.1 Mobility traits

We worked out the mobility pattern from a user by looking at geolocations from tweets. To characterize their mobility, we used the radius of gyration [38], which measures the size of the area covered while moving around:

$$R_g^i = \sqrt{\frac{1}{N} \sum_{i=1}^N (r_i - r_{\text{mean}})^2}, \quad (5)$$

where variable r_i represents the user's position at time instant i .

4.7.2 Content topics

We extracted topics from the texts in each user's tweets. To this end, we use the TextRazor classifier trained against the IPTC news-codes [59], which classify each tweet into approximately 1400 high-level categories organized into a three-level tree hierarchy. Each tweet is given a probability of containing such a topic. Thus each user is characterized by a content vector of n topics

$$C^i = \{C_1^i, C_2^i, \dots, C_n^i\}, \quad (6)$$

where the components C_m^i are the aggregated probability of topic m in all her tweets.

4.7.3 Network traits

Apart from the out-degree for each user i we also took into account the total user activity in the social network platform by computing the number of tweets generated during the observation period. This variable is called the number of posts.

4.8 Linear logistic regression model

The following equation represents a linear logistic regression model for explaining the probability of an individual being a sensor by different features, where $\{M^i\}$ are the mobility features (we only consider the radius of gyration variable, R_g), $\{N^i\}$ the group of network variables, out-degree, $d_{out,i}$, and the number of posts, and $\{C^i\}$ is the group of content variables for each individual i . Our model is

$$\Pr(i \in \Omega^*) = \text{logit}^{-1} \left[\beta_0 + \sum_l \alpha_l M_l^i + \sum_n \beta_n N_n^i + \sum_m \gamma_m C_m^i \right], \quad (7)$$

where Ω^* is the set of users defined as sensors, and $\text{logit}^{-1}(x) = e^x / (1 + e^x)$. In the model, each individual variable in the different groups is standardized to have zero mean and unit variance.

Abbreviations

ABM, Agent based model; EWES, Early warning epidemiological systems; ILI, Influenza-like illness; IPTC, International Press Telecommunications Council; NLP, Natural Language Processing; SIR, Susceptible-infected-recovery.

Acknowledgements

E.M. acknowledges support by Ministerio de Ciencia e Innovación/Agencia Española de Investigación (MCIN/AEI/10.13039/501100011033) through grant PID2019-106811GB-C32. M.C. wishes to acknowledge the following funding: the Ministry of Universities of the Government of Spain, under the program "Convocatoria de Ayudas para la recualificación del sistema universitario español para 2021-2023, de la Universidad Carlos III de Madrid, de 1 de Julio de 2021". Furthermore, he is thankful for the support of project "Ayuda PID2022-137243OB-I00 financiada por MCIN/AEI/10.13039/501100011033" and by "FEDER Una manera de hacer Europa".

Author contributions

M.G.-H. built the code to access Twitter data from their API. D.M.-C., E.M. and M.-C. designed the research. D.M.-C. performed research. D.M.-C. analysed the results. D.M.-C. and E.M. wrote the first draft of the manuscript. D.M.-C., E.M. M.-C. and M.G.-H. discussed results and edited the manuscript. All authors approved the final version.

Funding

Open access funding provided by Northeastern University Library.

Data availability

All data needed to evaluate the conclusions in the paper are present in the paper, the Supplementary Materials and the following repository: <https://github.com/dmartincc/sensors>. Additional data related to this paper may be requested from the authors.

Declarations

Competing interests

The authors declare no competing interests.

Author details

¹Department of Mathematics and GISC, Universidad Carlos III de Madrid, Leganes, 28911, Spain. ²Institute for Data, Systems, and Society, Massachusetts Institute of Technology, Cambridge, MA 02139, USA. ³UNICEF, New York City, USA. ⁴Center for Automation and Robotics, Spanish National Research Council, Madrid, Spain. ⁵Network Science Institute, Northeastern University, Boston, MA, 02115, USA. ⁶Department of Telematics and Computing, ICAI Engineering School, Universidad Pontificia Comillas, Madrid, 28015, Spain.

References

1. Parry J (2013) H7n9 avian flu infects humans for the first time. *BMJ* 346
2. Petersen LR, Jamieson DJ, Powers AM, Honein MA (2016) Zika virus. *N Engl J Med* 374(16):1552–1563
3. Stadler K, Masignani V, Eickmann M, Becker S, Abrignani S, Klenk H-D, Rappuoli R (2003) Sars—beginning to understand a new virus. *Nat Rev Microbiol* 1(3):209–218
4. Fouchier RA, Kuiken T, Schutten M, Van Amerongen G, Van Doornum GJ, Van Den Hoogen BG, Peiris M, Lim W, Stöhr K, Osterhaus AD (2003) Koch's postulates fulfilled for sars virus. *Nature* 423(6937):240–240
5. Feldmann H, Geisbert TW (2011) Ebola haemorrhagic fever. *Lancet* 377(9768):849–862
6. Briand S, Bertherat E, Cox P, Formenty P, Kieny M-P, Myhre JK, Roth C, Shindo N, Dye C (2014) The international Ebola emergency. *N Engl J Med* 371(13):1180–1183
7. Riou J, Althaus CL (2020) Pattern of early human-to-human transmission of Wuhan 2019 novel coronavirus (2019-ncov), December 2019 to January 2020. *Euro Surveill* 25(4):2000058
8. Li Q, Guan X, Wu P, Wang X, Zhou L, Tong Y, Ren R, Leung KS, Lau EH, Wong JY et al (2020) Early transmission dynamics in Wuhan, China, of novel coronavirus-infected pneumonia. *N Engl J Med*
9. Pastor-Satorras R, Vespignani A (2001) Epidemic spreading in scale-free networks. *Phys Rev Lett* 86(14):3200
10. Kogan NE, Clemente L, Liautaud P, Kaashoek J, Link NB, Nguyen AT, Lu FS, Huybers P, Resch B, Havas C et al (2021) An early warning approach to monitor covid-19 activity with multiple digital traces in near real time. *Sci Adv* 7(10):6989
11. Ginsberg J, Mohebbi MH, Patel RS, Brammer L, Smolinski MS, Brilliant L (2009) Detecting influenza epidemics using search engine query data. *Nature* 457(7232):1012–1014
12. Yuan Q, Nsoesie EO, Lv B, Peng G, Chunara R, Brownstein JS (2013) Monitoring influenza epidemics in China with search query from baidu. *PLoS ONE* 8(5):64323
13. McIver DJ, Brownstein JS (2014) Wikipedia usage estimates prevalence of influenza-like illness in the United States in near real-time. *PLoS Comput Biol* 10(4):1003581
14. Generous N, Fairchild G, Deshpande A, Del Valle SY, Priedhorsky R (2014) Global disease monitoring and forecasting with Wikipedia. *PLoS Comput Biol* 10(11):1003892
15. Lamos V, Majumder MS, Yom-Tov E, Edelstein M, Moura S, Hamada Y, Rangaka MX, McKendry RA, Cox IJ (2021) Tracking covid-19 using online search. *npj Digit Med* 4(1):1–11
16. Soebiyanto RP, Adimi F, Kiang RK (2010) Modeling and predicting seasonal influenza transmission in warm regions using climatological parameters. *PLoS ONE* 5(3):9450
17. Culotta A (2010) Towards detecting influenza epidemics by analyzing Twitter messages. In: Proceedings of the first workshop on social media analytics. ACM, New York, pp 115–122
18. Santillana M, Nguyen AT, Dredze M, Paul MJ, Nsoesie EO, Brownstein JS (2015) Combining search, social media, and traditional data sources to improve influenza surveillance. *PLoS Comput Biol* 11(10):1004513
19. Chen L, Hossain KT, Butler P, Ramakrishnan N, Prakash BA (2014) Flu gone viral: syndromic surveillance of flu on Twitter using temporal topic models. In: 2014 IEEE international conference on data mining. IEEE, pp 755–760
20. Ma L-l, Ma C, Zhang H-F, Wang B-H (2016) Identifying influential spreaders in complex networks based on gravity formula. *Phys A, Stat Mech Appl* 451:205–212
21. Christley RM, Pinchbeck G, Bowers R, Clancy D, French N, Bennett R, Turner J (2005) Infection in social networks: using network analysis to identify high-risk individuals. *Am J Epidemiol* 162(10):1024–1031
22. Alexander M, Forastiere L, Gupta S, Christakis NA (2022) Algorithms for seeding social networks can enhance the adoption of a public health intervention in urban India. *Proc Natl Acad Sci* 119(30):2120742119
23. Aleta A, Martín-Corral D, Bakker MA, Pastore y Piontti A, Ajelli M, Litvinova M, Chinazzi M, Dean NE, Halloran ME, Longini IM Jr et al (2022) Quantifying the importance and location of sars-cov-2 transmission events in large metropolitan areas. *Proc Natl Acad Sci* 119(26):2112182119
24. Wang Z, Bauch CT, Bhattacharyya S, d'Onofrio A, Manfredi P, Perc M, Perra N, Salathé M, Zhao D (2016) Statistical physics of vaccination. *Phys Rep* 664:1–113
25. Galesic M, Bruine de Bruin W, Dalege J, Feld SL, Kreuter F, Olsson H, Prelec D, Stein DL, van Der Does T (2021) Human social sensing is an untapped resource for computational social science. *Nature* 595(7866):214–222
26. Ghosh R, Mareček J, Griggs WM, Souza M, Shorten RN (2021) Predictability and fairness in social sensing. *IEEE Int Things J*
27. Rashid MT, Wang D (2021) Covidsens: a vision on reliable social sensing for covid-19. *Artif Intell Rev* 54(1):1–25
28. Hodas NO, Kooti F, Lerman K (2013) Friendship paradox redux: your friends are more interesting than you. *ICWSM* 13:8–10
29. Christakis NA, Fowler JH (2010) Social network sensors for early detection of contagious outbreaks. *PLoS ONE* 5(9):12948
30. Farrahi K, Emonet R, Cebrian M (2015) Predicting a community's flu dynamics with mobile phone data. In: Proceedings of the 18th ACM conference on computer supported cooperative work & social computing. ACM, New York, pp 1214–1221
31. Shao H, Hossain K, Wu H, Khan M, Vullikanti A, Prakash BA, Marathe M, Ramakrishnan N (2016) Forecasting the flu: designing social network sensors for epidemics. *ArXiv preprint. arXiv:1602.06866*
32. Kianersi S, Ahn Y-Y, Rosenberg M (2020) Association between sampling method and covid-19 test positivity among undergraduate students: Testing friendship paradox in covid-19 network of transmission. *medRxiv*
33. Garcia-Herranz M, Moro E, Cebrian M, Christakis NA, Fowler JH (2014) Using friends as sensors to detect global-scale contagious outbreaks. *PLoS ONE* 9(4):92413
34. Dunbar RI, Arnaboldi V, Conti M, Passarella A (2015) The structure of online social networks mirrors those in the offline world. *Soc Netw* 43:39–47
35. Zhang J, Centola D (2019) Social networks and health: new developments in diffusion, online and offline. *Annu Rev Sociol* 45:91–109
36. Grupo de Vigilancia de Gripe del Centro Nacional de Epidemiología. Instituto de Salud Carlos III: Sistema de Vigilancia de la Gripe en España. <http://vgripe.isciii.es/gripe/inicio.do> Accessed 22-06-2019
37. Commission E (2018) Commission implementing decision (eu) 2018/945 of 22 June 2018 on the communicable diseases and related special health issues to be covered by epidemiological surveillance as well as relevant case definitions. *Off J Eur Union* 61:1–74

38. González MC, Hidalgo CA, Barabási A-L (2008) Understanding individual human mobility patterns. *Nature* 453(7196):779–782
39. Kishore N, Taylor AR, Jacob PE, Vembar N, Cohen T, Buckee CO, Menzies NA (2022) Evaluating the reliability of mobility metrics from aggregated mobile phone data as proxies for SARS-CoV-2 transmission in the USA: a population-based study. *Lancet Digit Health* 4(1):27–36
40. Preoțiuc-Pietro D, Volkova S, Lampos V, Bachrach Y, Aletras N (2015) Studying user income through language, behaviour and affect in social media. *PLoS ONE* 10(9):0138717
41. Nelson KN, Siegler AJ, Sullivan PS, Bradley H, Hall E, Luisi N, Hipp-Ramsey P, Sanchez T, Shioda K, Lopman BA (2022) Nationally representative social contact patterns among U.S. adults, August 2020–April 2021. *Epidemics* 40:100605
42. Achrekar H, Gandhe A, Lazarus R, Yu S-H, Liu B (2011) Predicting flu trends using Twitter data. In: 2011 IEEE conference on computer communications workshops (INFOCOM WKSHPs). IEEE, pp 702–707
43. Golbeck J, Robles C, Edmondson M, Turner K (2011) Predicting personality from Twitter. In: 2011 IEEE third international conference on privacy, security, risk and trust and 2011 IEEE third international conference on social computing. IEEE, pp 149–156
44. Chamorro-Premuzic T, Furnham A (2007) Personality and music: can traits explain how people use music in everyday life? *Br J Psychol* 98(2):175–185
45. Reich SM, Subrahmanyam K, Espinoza G (2012) Friending, iming, and hanging out face-to-face: overlap in adolescents' online and offline social networks. *Dev Psychol* 48(2):356
46. Huang GC, Unger JB, Soto D, Fujimoto K, Pentz MA, Jordan-Marsh M, Valente TW (2014) Peer influences: the impact of online and offline friendship networks on adolescent smoking and alcohol use. *J Adolesc Health* 54(5):508–514
47. Aleta A, Martín-Corral D, Pastore y Piontti A, Ajelli M, Litvinova M, Chinazzi M, Dean NE, Halloran ME, Longini IM Jr, Merler S et al (2020) Modelling the impact of testing, contact tracing and household quarantine on second waves of covid-19. *Nat Hum Behav* 4(9):964–971
48. Li C, Chen LJ, Chen X, Zhang M, Pang CP, Chen H (2020) Retrospective analysis of the possibility of predicting the covid-19 outbreak from Internet searches and social media data, China, 2020. *Euro Surveill* 25(10):2000199
49. Qin L, Sun Q, Wang Y, Wu K-F, Chen M, Shia B-C, Wu S-Y (2020) Prediction of number of cases of 2019 novel coronavirus (covid-19) using social media search index. *Int J Environ Res Public Health* 17(7):2365
50. World H Organization: Global Outbreak Alert and Response Network (GOARN). <https://extranet.who.int/goarn/>
51. World H Organization: Integrated Outbreak Analytics (IOA). <https://extranet.who.int/goarn/content/integrated-outbreak-analytics-delivers-holistic-understanding-outbreak-dynamics>
52. World H Organization: Epidemic Intelligence from Open Sources (EIOS). <https://www.who.int/initiatives/eios>
53. World Health Organization: Epi-Brain. <https://www.epi-brain.com/>
54. Carter SE, Gobat N, Zambruni JP, Bedford J, Van Kleef E, Jombart T, Mossoko M, Nkakarande DB, Colorado CN, Ahuka-Mundeki S (2020) What questions we should be asking about covid-19 in humanitarian settings: perspectives from the social sciences analysis cell in the democratic republic of the Congo. *BMJ Glob Health* 5(9):003607
55. Twitter: Twitter Developer Documentation. <https://dev.twitter.com/streaming/overview>
56. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R, Dubourg V, Vanderplas J, Passos A, Cournapeau D, Brucher M, Perrot M, Duchesnay E (2011) Scikit-learn: machine learning in Python. *J Mach Learn Res* 12:2825–2830
57. Albert R, Barabási A-L (2002) Statistical mechanics of complex networks. *Rev Mod Phys* 74(1):47
58. Csardi G, Nepusz T (2006) The igraph software package for complex network research. *InterJournal Complex Systems* 1695
59. Troncy R (2008) Bringing the iptc news architecture into the semantic web. In: International semantic web conference. Springer, Berlin, pp 483–498

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Submit your manuscript to a SpringerOpen[®] journal and benefit from:

- Convenient online submission
- Rigorous peer review
- Open access: articles freely available online
- High visibility within the field
- Retaining the copyright to your article

Submit your next manuscript at ► [springeropen.com](https://www.springeropen.com)
