

ARTICLE OPEN



Tracking COVID-19 using online search

Vasileios Lampos¹✉, Maimuna S. Majumder^{2,3}, Elad Yom-Tov⁴, Michael Edelstein^{5,6}, Simon Moura¹, Yohhei Hamada⁷, Molebogeng X. Rangaka^{7,8}, Rachel A. McKendry^{9,10} and Ingemar J. Cox^{1,11}

Previous research has demonstrated that various properties of infectious diseases can be inferred from online search behaviour. In this work we use time series of online search query frequencies to gain insights about the prevalence of COVID-19 in multiple countries. We first develop unsupervised modelling techniques based on associated symptom categories identified by the United Kingdom's National Health Service and Public Health England. We then attempt to minimise an expected bias in these signals caused by public interest—as opposed to infections—using the proportion of news media coverage devoted to COVID-19 as a proxy indicator. Our analysis indicates that models based on online searches precede the reported confirmed cases and deaths by 16.7 (10.2–23.2) and 22.1 (17.4–26.9) days, respectively. We also investigate transfer learning techniques for mapping supervised models from countries where the spread of the disease has progressed extensively to countries that are in earlier phases of their respective epidemic curves. Furthermore, we compare time series of online search activity against confirmed COVID-19 cases or deaths jointly across multiple countries, uncovering interesting querying patterns, including the finding that rarer symptoms are better predictors than common ones. Finally, we show that web searches improve the short-term forecasting accuracy of autoregressive models for COVID-19 deaths. Our work provides evidence that online search data can be used to develop complementary public health surveillance methods to help inform the COVID-19 response in conjunction with more established approaches.

npj Digital Medicine (2021)4:17; <https://doi.org/10.1038/s41746-021-00384-w>

INTRODUCTION

Over the past several years, numerous scientific studies have shown that user interactions with web applications generate latent health-related signals, reflective of individual as well as community level trends^{1–8}. Seasonal influenza and the H1N1 pandemic were used as case studies for the development and evaluation of machine learning models that produce disease rate estimates in a non-traditional way, using online search or social media as their input information^{1,2,5,9,10}. After an extensive scientific debate about their usefulness and accuracy^{11,12}, and with the manifestation of new findings that improved upon past methodological shortcomings^{13–15}, these techniques are now becoming part of public health systems, serving as complementary endpoints for monitoring the prevalence of infectious diseases, such as seasonal influenza^{8,16}. Compared to conventional health surveillance systems, online user trails exhibit certain advantages, including low latency, continuous operational capacity, denser spatial coverage, and broader demographic inclusion^{8,17}. Furthermore, during emerging epidemics they may offer community level insights that current monitoring systems are not equipped to obtain given a limited testing capacity¹⁸, and the physical distancing measures that discourage or prohibit people from interacting with health services¹⁹. Notably, the ongoing coronavirus disease (COVID-19) pandemic, caused by a novel coronavirus (SARS-CoV-2), has generated an unprecedented relative volume of online searches (Supplementary Fig. 1). This search behaviour could signal the presence of actual infections, but may also be due to general concern that is intensified by news media coverage, reported figures of disease incidence and

mortality across the world, and imposed physical distancing measures^{20,21}.

Previous research, which has used online search to predominantly model influenza-like illness (ILI) rates, has focused on supervised learning solutions, where “ground truth” information, in the form of historical syndromic surveillance reports, can be used to train machine learning models^{2,5,13,14,22}. These models learn a function that maps time series of online user-generated data, e.g. the frequency of web searches over time, to a noisy representation of a disease rate time series, indirectly minimising the impact of online activity that is not caused by infection. Typically, this training data spans multiple years and flu seasons^{13–15}. However, for most locations, if not all, no sufficient data—in terms of validity, representativeness, and time span—currently exists to apply supervised learning approaches to COVID-19. Therefore, unsupervised solutions that attempt to minimise the effect of concern should be sought, and fully supervised solutions should be used and interpreted with caution.

In this paper, we present models for COVID-19 using online search data based in both unsupervised and supervised settings. We first develop an unsupervised model for COVID-19 by carefully choosing search queries that refer to related symptoms as identified by a survey from the National Health Service (NHS) and Public Health England (PHE) in the United Kingdom (UK)²³. Symptom categories are weighted based on their reported ratio of occurrence in cases of COVID-19. The output from these new models provides useful insights, including early warnings for potential disease spread, and showcases the effect of physical distancing measures. Since online searches can also be driven by

¹Department of Computer Science, University College London, London, UK. ²Computational Health Informatics Program, Boston Children's Hospital, Boston, MA, USA. ³Department of Pediatrics, Harvard Medical School, Boston, MA, USA. ⁴Microsoft Research, Herzliya, Israel. ⁵National Infection Service, Public Health England, London, UK. ⁶Department of Population Health, Faculty of Medicine, Bar-Ilan University, Safed, Israel. ⁷Institute for Global Health, University College London, London, UK. ⁸Division of Epidemiology and Biostatistics, University of Cape Town, Cape Town, South Africa. ⁹London Centre for Nanotechnology, University College London, London, UK. ¹⁰Division of Medicine, University College London, London, UK. ¹¹Department of Computer Science, University of Copenhagen, Copenhagen, Denmark. ✉email: v.lampos@ucl.ac.uk

concern rather than by infection, we attempt to minimise this effect by incorporating a news media coverage time series. The influence of news media on the online search signal is quantified through an autoregressive task that resembles a Granger causality test between these two time series²⁴.

Recent research has demonstrated that a model estimating ILL rates from web searches, trained originally for a location where syndromic surveillance data is available, can be adapted and deployed to another location that cannot access historical ground truth data²⁵. The accuracy of the target location model depends on identifying the correct search queries and their corresponding weights via a transfer learning methodology. By adapting this method, we map supervised COVID-19 models from a source to a target country, in an effort to transfer noisy knowledge from areas that are ahead in their epidemic curves to areas that are at earlier stages. This supervised approach transfers the clinical reporting biases of the source location, but given the statistical supervision, it is not affected by concern to the same degree as the unsupervised models. Our analysis reaffirms the insights of the unsupervised approach, showcasing how early warnings could have been obtained from locations that had already experienced the impact of COVID-19.

Taking noisy supervision a step further, we conduct a correlation and regression analysis to uncover potentially useful online search queries that refer to underlying behavioural or symptomatic patterns in relation to confirmed COVID-19 cases. We show that generic COVID-19-related terms and less common symptoms are the most capable predictors. Finally, we use the task of forecasting confirmed COVID-19 deaths to illustrate that web search data, when added to an autoregressive forecasting model, significantly reduces prediction error.

We present results for a multilingual and multicultural selection of countries—namely, the United States of America (US), UK (including a comparative analysis for England), Australia, Canada, France, Italy, Greece, and South Africa.

RESULTS

Unsupervised models

The first part of the analysis presents unsupervised models of COVID-19 based on weighted search query frequencies. Queries and their weightings are determined using the first few hundred (FF100) survey on COVID-19 conducted by the NHS/PHE in the UK²³. FF100 identified 19 symptoms associated with confirmed COVID-19 cases and their probability of occurrence (Supplementary Table 1). There exist other studies, such as the findings by Carfi et al.²⁶, that corroborate the outcomes of the FF100 survey. Our choice to use it was based on the substantial size of the cohort (381 patients) and the comprehensive presentation of outcomes. In addition to the symptoms identified by FF100, we also include queries that mention COVID-19-related keywords (e.g. “covid-19” or “coronavirus”) as a separate category. Unsupervised models for COVID-19 in 8 countries are depicted in Fig. 1. We observe exponentially increasing rates that exceed the estimated seasonal average (previous 8 years) during their peak period in all investigated countries, as well as a steep drop of the score after the application of physical distancing or lockdown measures in most countries which is in concordance with clinical surveillance reports²⁷. In the time series where we have attempted to minimise the effect of news by maintaining a proportion of the signal (Eqs. (3), (4) and (5)), we observe more conservative estimates in all locations, including altered trends during the peak period. Compared to the original scores (no minimisation of news media effects), there is an average reduction of the signals by 16.4% (14.2%–18.7%) in a period of 14 days prior and after their peak moments; their corresponding average linear correlation during this period is equal to 0.822 (0.739–0.905). Outside of the peak

periods, the reduction is a moderate 3.3% (2.7%–4.0%), indicating that there is an association between the extent of media impact and the estimated disease prevalence. For the US, Australia, Canada, and Italy, the signal is visibly flattened during the time period around their respective peaks. Countries that took measures earlier in the epidemic curve (e.g., Greece) demonstrate a more pronounced pattern of decrease, with scores that go below the expected seasonal average. We also note that for Australia and the UK, search scores were already in decline after the application of physical distancing measures but before lockdowns. Outcomes based solely on the FF100 symptoms (Supplementary Fig. 2) or without using any weighting scheme (Supplementary Fig. 3) are available in the Supplementary Information (SI).

A comparison of the search scores with minimised media effects to the time series of confirmed cases is depicted in Fig. 2. If we exclude South Africa, as it displays an outlying behaviour, perhaps due to a limited testing capacity^{28,29} or demographically-skewed Internet access patterns³⁰, the correlation between these time series is maximised, reaching an average value of 0.826 (0.735–0.917) when clinical data is brought forward by 16.7 (10.2–23.2) days. This provides an indication of how much sooner the proposed unsupervised models could have signalled an early warning about these epidemics at a national level. It also shows numerically as well as visually that both the increase and the decline of confirmed cases are captured by the search signal. Replacing confirmed cases with deaths caused by COVID-19 (Supplementary Fig. 4) increases this period to 22.1 (17.4–26.9) days with a slightly greater maximised correlation ($r = 0.846$; 0.702–0.990).

Transfer learning

Models of confirmed COVID-19 cases are transferred from Italy (source) to all other countries (targets) in our analysis using a transfer learning methodology. In contrast to the unsupervised models, here we attempt to leverage information from a country that is ahead in terms of epidemic progression³¹. As a result, the obtained estimates are reflective of the clinical reporting systems in the source country, but not as influenced by user concern at the target countries given that they are derived from a supervised learning function. Our approach maps search queries about specific symptom categories, as identified by the FF100 survey, from the language of the source country (Italian) to the languages of the other countries using a temporal correlation metric (see Methods). We train 100,000 elastic net models for the source location, exploring the entire ℓ_1 -norm regularisation path, and transfer the subset of models (84,557 on average) that assign a nonzero weight to a minimum of 3 and up to a maximum of 49 queries (out of the 54 queries identified for Italy) on a daily basis from February 17 to May 24, 2020 (both inclusive). The transferred estimates for the last day in the considered time span and their 95% confidence intervals are depicted in Fig. 3. Intermediate outcomes showing estimates from models that are trained and transferred on a daily basis are available in Supplementary Fig. 5. With the exceptions of Australia and South Africa, the peak of the transferred signal appears approximately 2–3 weeks earlier than the confirmed cases. The decrease after the peak is also more steep in the transferred models and always occurs after the commencement of physical distancing or lockdown measures. The average correlation between the transferred and the unsupervised (with minimised news effects) estimates is equal to 0.654 (0.572–0.735), which implies that there exist some differences between the two approaches (Supplementary Fig. 6). In fact, their correlation is maximised to 0.789 (0.736–0.842) by bringing the transferred time series 5.14 (3.21–7.08) days forward, indicating that the unsupervised signal provides an earlier alert, and that both signals are much more similar when their temporal

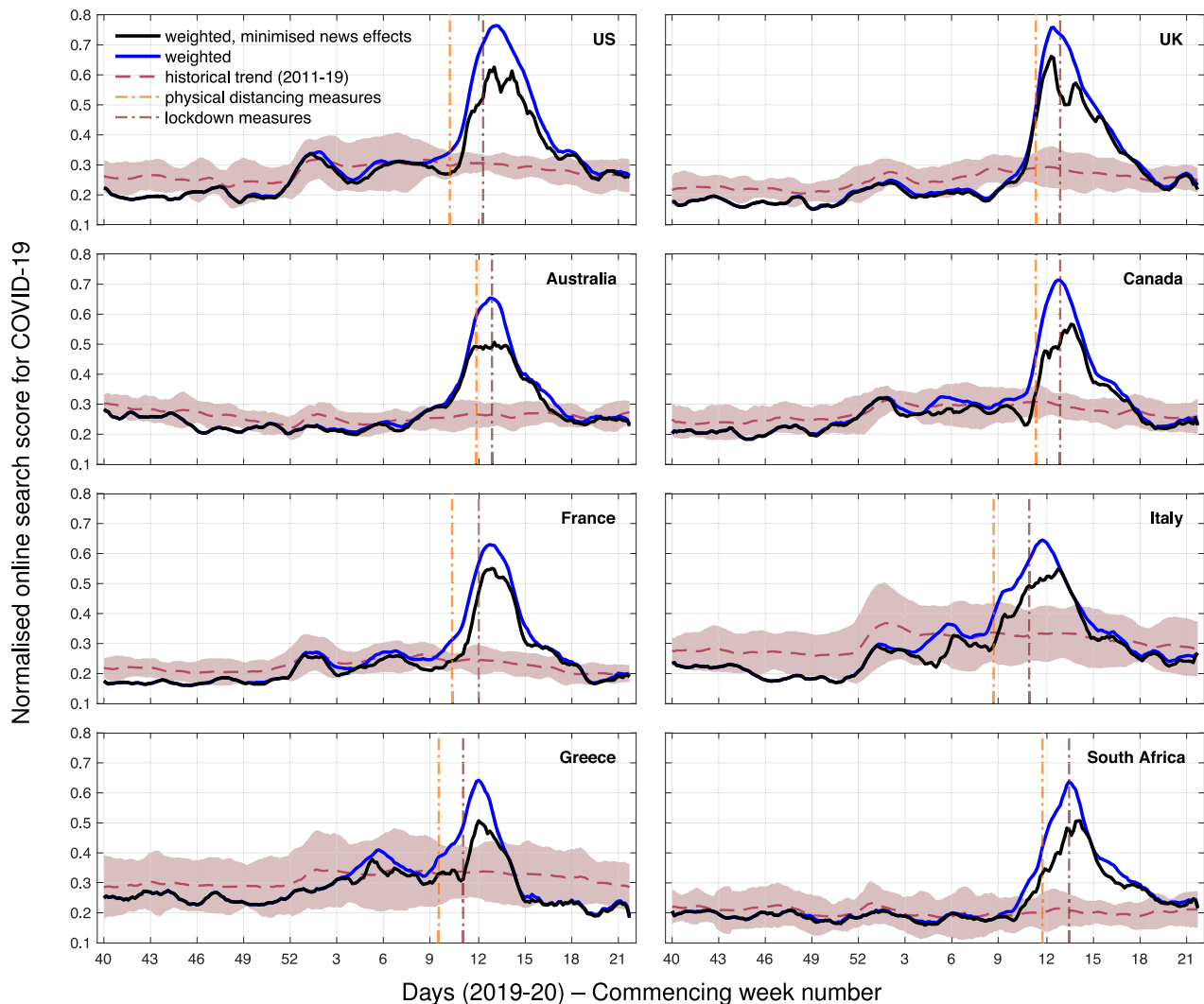


Fig. 1 Online search scores for COVID-19-related symptoms as identified by the FF100 survey, in addition to queries with coronavirus-related terms, for 8 countries from September 30, 2019 to May 24, 2020 (all inclusive). Query frequencies are weighted by symptom occurrence probability (blue line) and have news media effects minimised (black line). These scores are compared to an average 8-year trend of the weighted model (dashed line) and its corresponding 95% confidence intervals (shaded area). Application dates for physical distancing or lockdown measures are indicated with dash-dotted vertical lines; for countries that deployed different regional approaches, the first application of such measures is depicted. All time series are smoothed using a 7-point moving average, centred around each day.

discrepancy is corrected. Interestingly, during the temporal alignment of source and target search query frequencies, which is an intermediate step of the transfer learning technique, correlation increases when the target data is shifted forward. This partially confirms that Italy was indeed ahead by a few days in terms of either epidemic progression, user search behaviour, or both, and justifies our choice to use it as the source country. In particular, when we focus on the period from March 16 to May 24, 2020 (both inclusive)—dates that signify the beginning and end of high levels of transmission for Italy—and analyse all transferred models, the average shift in days that maximises these correlations per target country is: 13.76 (12.97–14.55) for the US, 12.67 (11.99–13.35) for the UK, 5.24 (4.30–6.18) for Australia, 8.06 (7.14–8.98) for Canada, 9.84 (8.62–11.06) for France, 1.44 (0.45–2.43) for Greece, and 10.99 (10.10–11.88) for South Africa.

Correlation and regression analysis

Aiming to uncover symptom-related patterns or associated behaviours, we examine the statistical relationship between web search frequencies and confirmed COVID-19 cases or deaths by

performing a correlation and regression analysis. Outcomes could also help inform the choice of search terms in follow-up models for COVID-19. To reduce the representation bias of clinical endpoints, we combine data from multiple countries to the extent possible. For a more comprehensive experiment, we aggregate data from 4 countries where English is the main spoken language, namely the US, UK, Australia, and Canada. We first estimate the linear correlation between query frequencies and clinical data at all locations (in an aggregate fashion), from December 31, 2019 up to and including May 24, 2020. Results illustrating the top-correlated and anti-correlated search queries are depicted in Fig. 4a; correlations with deaths, which are lower given the additional temporal difference, are depicted in Supplementary Fig. 7a. Queries about the disease (“covid”; $r = 0.72$) or the virus (“sars cov 2”; $r = 0.67$) are the top-correlated. Various related symptoms demonstrate strong correlations as well, although the amount of correlation is not necessarily reflective of the symptom’s occurrence probability; examples include rash ($r = 0.63$), pink eye ($r = 0.58$), blue face ($r = 0.56$), sneezing ($r = 0.55$), and loss of the sense of smell ($r = 0.49$). Associated behaviours,

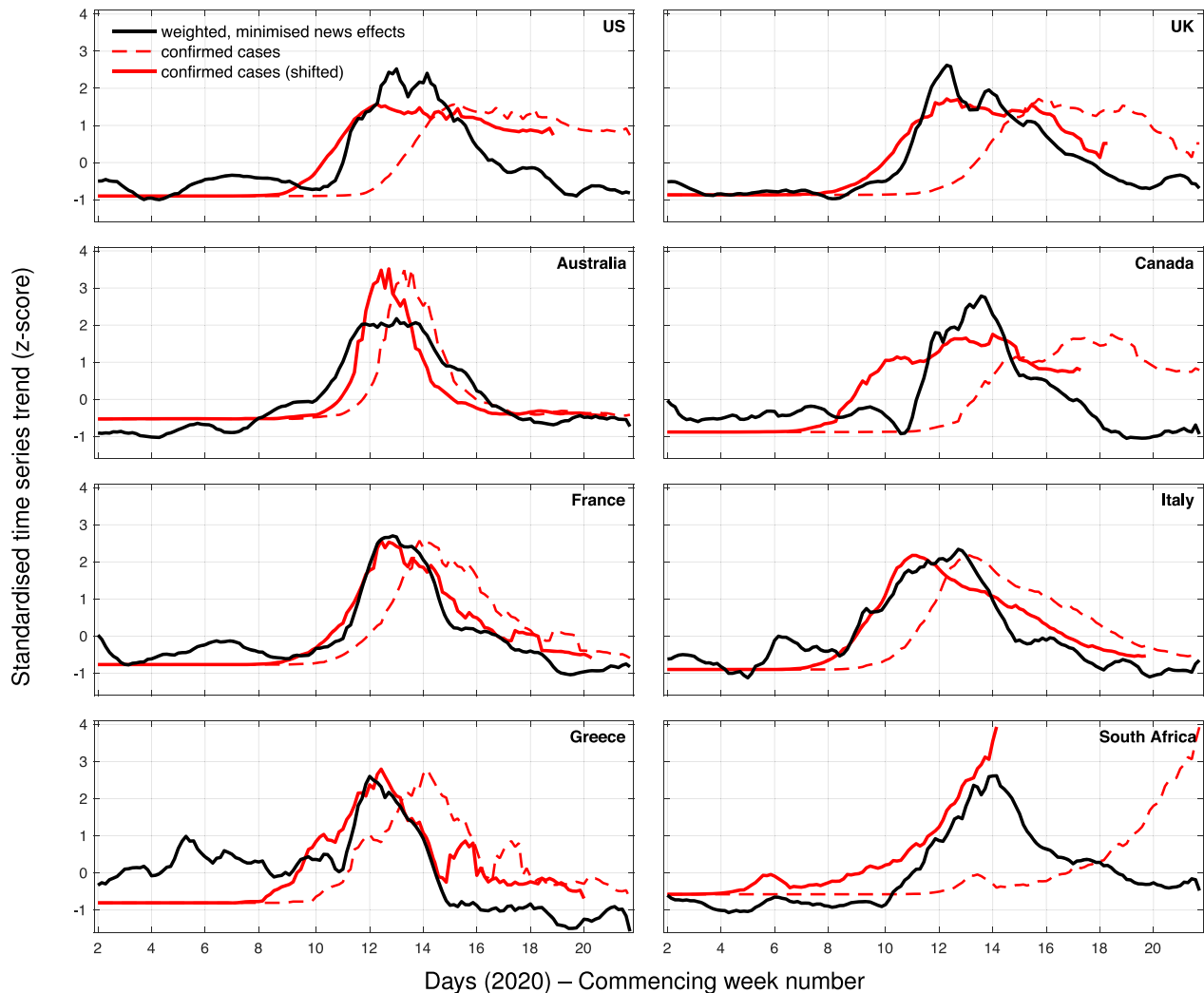


Fig. 2 Comparison between online search scores with minimised news media effects (black line) and confirmed cases (dashed red line), as well as confirmed cases shifted back (red line) such that their correlation with the online search scores is maximised. The confirmed cases time series are shifted back by a different number of days for each country: 20 days (US), 24 days (UK), 6 days (Australia), 31 days (Canada), 10 days (France), 14 days (Italy), 12 days (Greece), and 53 days (South Africa). All time series are smoothed using a 7-point moving average, centred around each day.

e.g. staying at home ($r = 0.61$), or measures, e.g. quarantine ($r = 0.60$), are also strongly correlated. On the other hand, symptoms such as vomiting ($r = -0.60$) and migraine ($r = -0.44$) show a significant anti-correlation. When we shift confirmed cases back by 19 days and deaths by 25 days, the average correlation across all considered web searches is maximised; results are depicted in Supplementary Fig. 7b, c. The membership and ranking is similar to the previously shown results, although there are marginal increases in the numbers of search queries that seek information about the disease (characteristics, testing, potential treatments). For the regression analysis, we focus on the last 84 days of the considered time span (March 2 to May 24, 2020), training models and assessing estimates for each day of that period. We consider elastic net models with up to a 50% feature density (nonzero weights for half of the considered search queries at most) to minimise the chance of overfitting. The outcomes of this analysis are depicted in Fig. 4b, c; for completeness, results where the entire regularisation path is explored (1%–100% feature density) are shown in Supplementary Fig. 7d, e. We observe that the most impactful feature, in terms of its normalised contribution to estimates of confirmed cases or deaths, are search queries that

include the term “covid” (29.32% and 36.29%, respectively), which is intuitive given the magnitude of this pandemic. Blue face (27.05%), loss of the sense of smell (12.57%), appetite loss (7.39%), pink eye (5.64%), and shortness of breath (5.42%) are the top-5 most impactful symptoms with regards to estimating confirmed cases. For mortality estimates, two new symptoms are introduced in the top-5: rash (8.80%) and loss of the sense of taste (6.5%). In addition, recommended behaviours, such as isolation and staying at home, and queries about masks, related diagnostics, or holidays are also impactful. Similarly to the correlation analysis, the search queries related to vomiting have a strong negative impact (-12.78%) in estimating confirmed cases; this is replaced by diarrhoea when estimating deaths (-15.26%). Queries about the common symptoms of cough, fatigue, and fever are not among the most correlated or impactful, suggesting that web searches about rarer symptoms may be more informative.

Forecasting deaths caused by COVID-19

As a final step in our analysis, we use fully supervised forecasting models for COVID-19 deaths to assess whether the inclusion of web searches can improve the accuracy of autoregressive (AR)

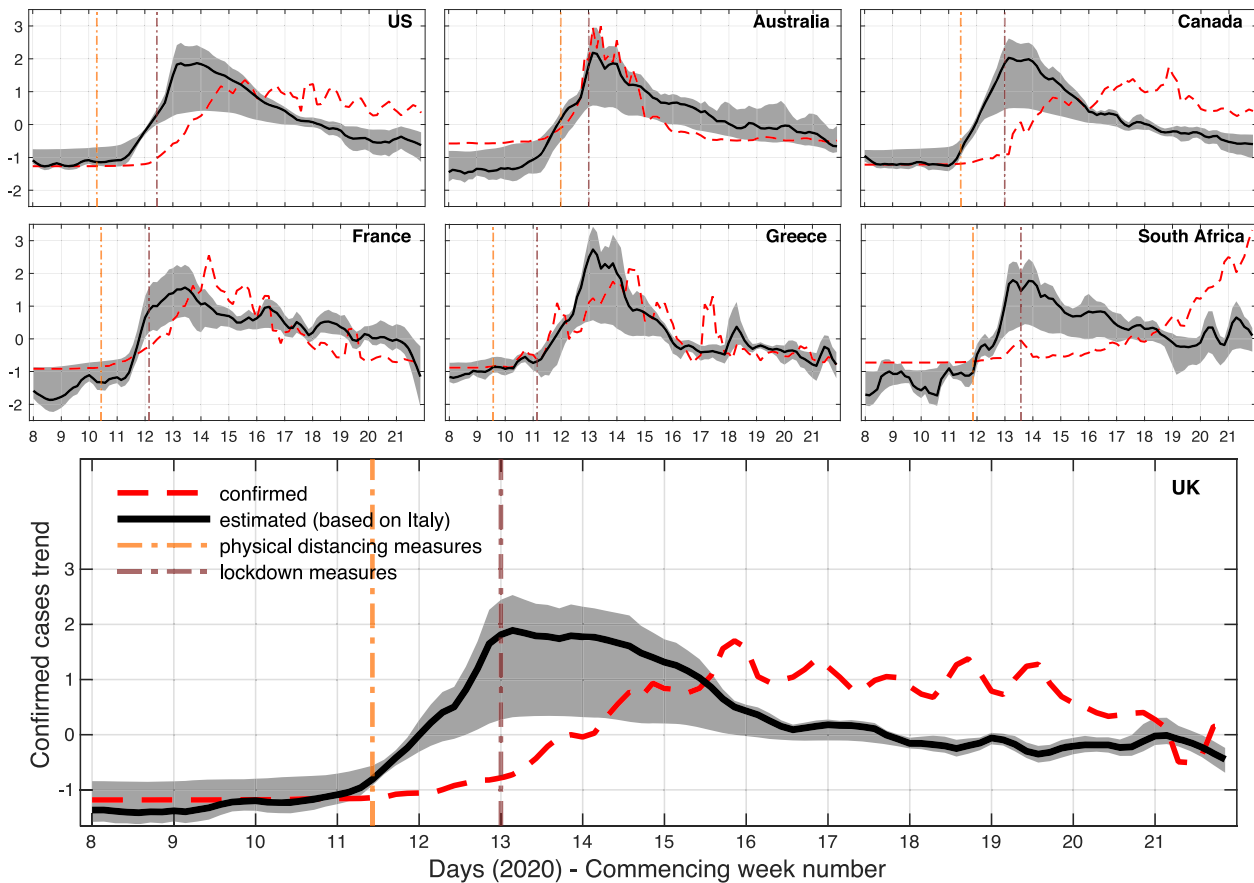


Fig. 3 Transfer learning models based on online search data for 7 countries using Italy as the source country. The figures show an estimated trend for confirmed COVID-19 cases compared to the reported one. The trend is derived by standardising the transferred estimates (raw values are reflective of the demographics and clinical reporting approach of the source country). The solid line represents the mean estimate from an ensemble of models. The shaded area shows 95% confidence intervals based on all model estimates. Application dates for physical distancing or lockdown measures are indicated with dash-dotted vertical lines; for countries that deployed different regional approaches, the first application of such measures is depicted. Time series are smoothed using a 3-point moving average, centred around each day. We use this minimum amount of smoothing to remove some of the noise for visualisation purposes and maintain our ability to compare the transferred models to the corresponding clinical data.

models. We provide the same analysis for confirmed cases in the SI (Supplementary Table 2, and Supplementary Fig. 8), but given the irregularities in the way laboratory diagnostics were conducted, the time series of deaths is more consistent, and hence more appropriate to use for this challenging supervised learning task. We first assess an AR model that uses only deaths data from the past $L = 6$ days (AR-F) to conduct forecasts 7 and 14 days ahead. We then expand on this by incorporating online search data as well (SAR-F). Both models are based on Gaussian Processes (GPs) as detailed in Methods. We also use a basic persistence model (PER-F) as a modest baseline. Our testing period starts from April 20, 2020, but we commence testing only when a cumulative number of 10 deaths is recorded in a country (this reduces the amount of test points for South Africa only). Models are retrained at every time step, and their accuracy is assessed using the mean absolute error (MAE) between forecasts and actual figures. Table 1 enumerates these results, including a normalised average MAE across locations, models, and forecasting tasks to allow for a fairer joint interpretation. We note that SAR-F performs considerably better than AR-F, decreasing MAE by 32.65% and 33.77% in the 7 and 14 days ahead forecasting tasks, respectively. It also improves upon the PER-F baseline in the more challenging task, which is a positive outcome given the small amount of available training data. The corresponding 14 days ahead forecasts are depicted in Supplementary Fig. 9. As expected from the empirical evaluation,

SAR-F estimates are visibly a better fit to the ground truth than the ones produced by AR-F, capturing the quantity and the overall trend more convincingly. These outcomes provide further evidence for the utility of incorporating online search information in disease models for COVID-19.

DISCUSSION

We have presented unsupervised (with minimised news media effects) and transfer learning models for COVID-19 based on online search data. The latter reaffirm the estimates of the former, albeit with an additional temporal delay of approximately 5 days. However, transfer learning provides a more statistically principled approach as it is based on a supervised model of confirmed cases from a source country (Italy). This is a component that makes it less prone to media influence, similarly to supervised^{14,15} or transferred²⁵ models for ILI. We have conducted a series of experiments, across different countries, to demonstrate the practical utility of online search in modelling the incidence of COVID-19. By comparing our outcomes to clinical endpoints, we argue that signals from web search data could have served as preliminary early indicators for COVID-19 prevalence at the national level. Our results also highlight the immediate impact that physical distancing or lockdown measures had in reducing disease rates. Furthermore, a qualitative analysis shows that rarer

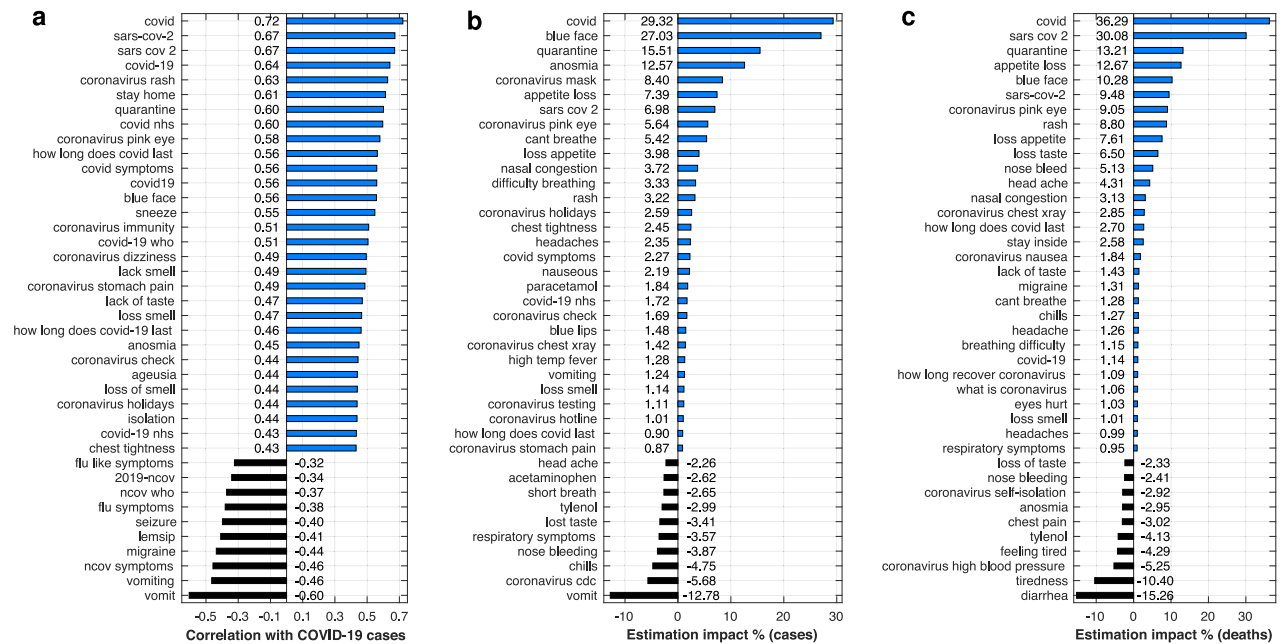


Fig. 4 Correlation and regression analysis of search query frequencies against confirmed COVID-19 cases or deaths in four English speaking countries (US, UK, Australia, and Canada). **a** Top-30 positively and top-10 negatively correlated search queries with COVID-19 confirmed cases; **b** Top-30 positively and top-10 negatively impactful queries in estimating COVID-19 confirmed cases; **c** Top-30 positively and top-10 negatively impactful queries in estimating deaths caused by COVID-19.

Table 1. Average mean absolute error and standard deviation (in parentheses) of forecasting models (7 and 14 days ahead) for daily deaths caused by COVID-19 in 8 countries.

Country	7 days ahead			14 days ahead		
	AR-F	SAR-F	PER-F	AR-F	SAR-F	PER-F
US	774.49 (544.70)	545.97 (447.31)	545.06 (667.84)	1049.07 (697.66)	591.55 (454.22)	613.66 (453.73)
UK	206.08 (183.25)	128.37 (115.20)	137.89 (82.19)	370.25 (242.63)	174.77 (181.53)	227.46 (129.86)
Australia	1.48 (0.87)	0.92 (0.69)	0.97 (1.01)	1.14 (0.68)	0.94 (0.63)	1.66 (1.35)
Canada	74.61 (46.70)	55.00 (34.39)	33.89 (24.47)	108.48 (50.75)	82.36 (43.96)	52.86 (35.79)
France	354.49 (119.69)	193.81 (120.62)	151.54 (156.93)	384.38 (177.43)	245.85 (147.42)	282.03 (243.91)
Italy	207.21 (97.77)	85.77 (54.34)	85.69 (66.30)	308.99 (134.86)	162.03 (95.31)	157.86 (86.09)
Greece	1.76 (1.09)	1.43 (0.96)	1.97 (1.64)	2.14 (1.23)	1.40 (1.07)	1.86 (1.78)
South Africa	6.92 (5.85)	7.33 (6.53)	6.14 (5.58)	9.14 (7.45)	8.96 (7.25)	7.68 (6.25)
Norm. mean	0.294 (0.192)	0.198 (0.161)	0.187 (0.187)	0.382 (0.241)	0.253 (0.198)	0.266 (0.215)

The last row contains min-max normalised averages across countries, methods, and forecasting tasks to account for the different ranges in different countries. AR-F autoregressive forecasting using past deaths, SAR-F combined online search and autoregressive forecasting, PER-F persistence model.

symptoms or generic queries about COVID-19 correlate better with and are more predictive of clinically reported metrics.

Assessing the correctness of our analysis is difficult because there exists no definitive ground truth representative of community level disease rates to compare against. However, intrinsic properties in our findings as well as quantitative comparisons with different clinical endpoints provide at least partial evidence of validity. It is important to note that the COVID-19 outbreak in Italy, which was the first major outbreak in Europe³², did not cause an increase in the frequency of the vast majority (>70%) of search terms at the other countries we considered based on a Granger causality analysis²⁴ (see SI). Nonetheless, this motivates our effort to minimise news media effects on the derived COVID-19 scores, but at the same time, it can also justify the moderate decrease that this process registers during peak periods (16.4% on average). For Italy itself, the COVID-19 score as estimated through online

searches (with minimised media effects) is preceding confirmed cases and deaths by 14 and 18 days, respectively. In addition, search terms that in most instances do not represent infection (e.g., “COVID unemployment”) clearly follow the corresponding trends of search terms about COVID-19 symptoms (Supplementary Fig. 10). These observations combined corroborate the hypothesis that the greatest portion of the online search signal based on the terms we have identified and the derivatives we develop is more representative of infection rather than concern.

As described in the results, the correlation between the unsupervised search-based signal with minimised news effects and confirmed COVID-19 cases is maximised when the latter is brought forward by 16.7 days. This temporal difference between online searches and clinical information could be partly explained by the amount of time between the onset of common symptoms and of more severe ones that warrant a hospital admission³³, the

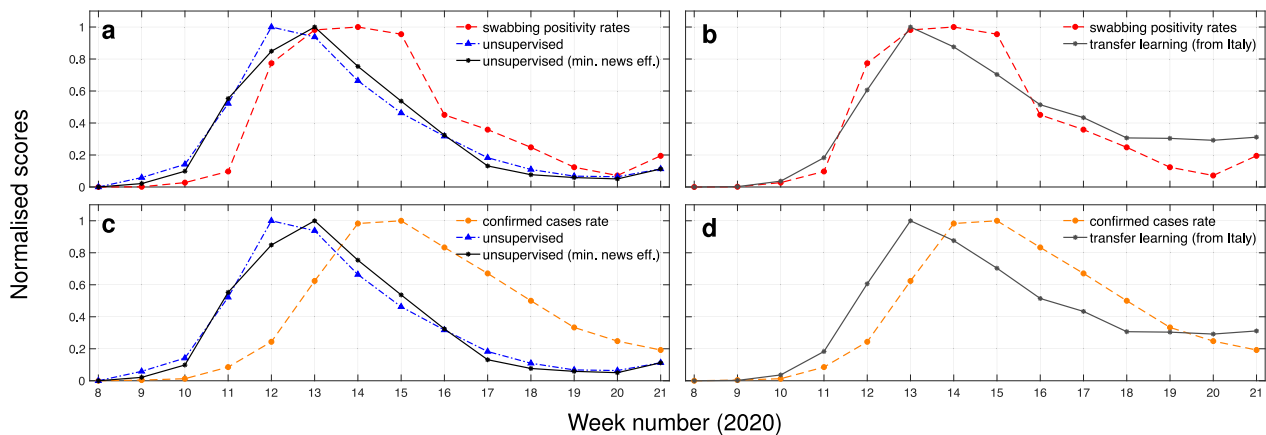


Fig. 5 Comparison of weekly online search-based signals for COVID-19 to different clinical endpoints for England. **a** Estimates from the unsupervised models with or without minimising news media effects are compared to COVID-19 positivity rates obtained through a swabbing test scheme operated by the RCGP. **b** Estimates for COVID-19 obtained via transfer learning (source model of confirmed cases is based on data from Italy) are compared to COVID-19 positivity rates obtained through a swabbing test scheme operated by the RCGP. **c** Estimates from the unsupervised models with or without minimising news media effects are compared to confirmed cases rates as reported by PHE. **d** Estimates for COVID-19 obtained via transfer learning (source model of confirmed cases is based on data from Italy) are compared to confirmed cases rates as reported by PHE.

overall delay of some health systems to respond to the pandemic³⁴, and national policies that recommended testing only after the persistence of symptoms in milder cases³⁵. Interestingly, if no minimisation of news effects is conducted, the temporal distance for a maximal correlation increases by 1.7 days. The fact that most countries in our analysis exhibited concern about COVID-19 prior to the rapid increase of infections could be a justification for this. By repeating this analysis for mortality figures, we obtain that correlation is maximised by bringing death time series 22.1 (17.4–26.9) days forward. The added amount of days corroborates with findings about the time interval between the commencement of a hospitalisation and a death outcome (5–11 days)³³.

Drawing our focus to England, we were able to compare outcomes from the unsupervised (with and without minimised effects) and transfer learning approaches to more conventional metrics. The first is a swabbing scheme led by the Royal College of General Practitioners (RCGP), which included many patients who did not have COVID-19-related symptoms and thus allowing us to obtain a more representative metric for community-level spread³⁶; results are depicted in Fig. 5a, b. We observe that the weekly COVID-19 swabbing positivity rates correlate strongly with the unsupervised ($r = 0.816$, $p < 0.001$), unsupervised with minimised news media effects ($r = 0.855$, $p < 0.001$), and transferred ($r = 0.954$, $p < 0.001$) scores. We find that the unsupervised models could have provided an early warning at least a week before the RCGP swabbing scheme. The second metric is a normalised version of the confirmed cases that takes into consideration relative population figures (coronavirus.data.gov.uk); comparisons are depicted in Fig. 5c, d. Here, the time discrepancy is more visible, and correlations without any form of shifting are within a similar range as in our previous analysis (Figs. 2 and 3), showing that search-based estimates precede the clinical estimates by 1 and 2 weeks for the transfer learning and unsupervised approaches respectively.

PHE in the UK has incorporated the unsupervised models described in this paper into their weekly surveillance reports about COVID-19 (gov.uk/government/publications/national-covid-19-surveillance-reports). According to PHE, estimates from these models provided important insights into community transmission at a time where conventional surveillance had limited reach due to the fact that testing was restricted to hospitalised cases. They were analysed alongside other surveillance systems, with the

broad expectation that any changes in trends would be first detected through the analysis of web searches and would then be seen at a later stage in more conventional surveillance sources. The decrease in the web search scores was interpreted as an early signal of the impact of physical distancing³⁷. The ability to quantify the lag between search engine trends and incidence trends in confirmed cases and mortality further improves the utility of this approach.

Related work has also indicated the potential value in using online search as an early indicator of COVID-19^{38,39}. However, the provided analysis is not entirely comprehensive, in that it focuses in a very narrow subset of search terms (e.g., just the term “coronavirus” for Effenberger et al.³⁹), only considers a few COVID-19 symptoms and does not weight them according to their frequency in patient cohorts, and does not account for potential biases in the obtained signal from concerned users. In fact, the reported lags compared to confirmed cases are not in par with our findings that add 5–6 more days of early warning. In addition, a study reporting that searches have been mostly influenced by the news media again focuses on a very narrow set of search terms, and does not conduct a causal analysis to support that claim⁴⁰. We show that although some of the search terms could have been driven by concern, this becomes a minority when they are curated carefully (see SI), and we provide a method to reduce some of the remaining biases in the signal.

Nonetheless, our study is not without limitations. Primarily, it is currently impossible to provide a solid evaluation for our findings. This will require a better representation of disease incidence from a clinical or epidemiological perspective, which is an ongoing challenge. In addition, wherever a form of statistical supervision has been introduced (transfer learning, correlation and regression analysis, forecasting), reliable outcomes depend on the existence of a consistent sampling technique of the target information, which in this case is confirmed cases or deaths. We made a concerted effort to reduce the impact of such inconsistencies by training models jointly across multiple countries, focusing on mortality figures rather than the testing-dependent confirmed cases for more demanding tasks, and assessing qualitatively our outcomes across different locations. More detailed comparisons based on public health data from England reaffirm our broader findings. From a practical perspective, this work might not directly be applicable to healthcare resource allocation. This requires more geographically granular estimates that are not supported by the

currently available online search data. However, our method can provide an alert at the national level that subsequently could motivate a more thorough examination using different disease surveillance systems and lead to the identification of outbreak hotspots. Online search behaviour could be influenced by many factors not related to infection. Although we have attempted to mitigate those effects by carefully selecting search terms and developing a method to minimise news effects, some biases may still be present in the obtained signals. Our Granger-causal analysis indicates that only a minority of the searches might have been affected (see SI). Finally, similarly to other studies that are based on online search data^{13,14}, our approach is likely to have limited applicability to locations with lower rates of Internet access.

METHODS

We first present the methodological approaches deployed in our analysis, and then describe the data sets used.

Community-level surveillance with an unsupervised symptom-based online search model

We generate k symptom-based search query sets using the $k=19$ identified symptoms from the FF100 survey for COVID-19 (Supplementary Table 1)²³. We also consider one additional set that includes specific COVID-19 terminology, i.e. the “COVID-19” keyword itself among others. Query sets may include different wordings for the same symptom or queries with minor grammatical differences (especially for queries in Greek and French). If a symptom is represented by more than one search query, then we obtain the total frequency (sum) across these queries. Each query set time series is smoothed using a harmonic mean over the past 14 days (see Supplementary Eq. 4 for a definition of the harmonic mean), and any trends across the entire period of the analysis are removed using linear detrending. We then apply a min-max normalisation (Supplementary Eq. 1) to the frequency time series of each query set to obtain a balanced representation between more and less frequent searches. We divide our data into two periods of interest, the current one (from September 30, 2019 until May 24, 2020) and a historical one (from September 30, 2011 to September 29, 2019). The corresponding data sets are denoted by $\mathbf{X} \in \mathbb{R}_{>0}^{N_1 \times k}$ and $\mathbf{H} \in \mathbb{R}_{>0}^{N_2 \times k}$, where N_1, N_2 represent the different numbers of days in the current and historical data, respectively. We use the symptom conditional probability distribution from the FF100 to assign weights ($\mathbf{w} \in \mathbb{R}_{>0}^k$) to each query category, and compute weighted time series ($\mathbf{x} = \mathbf{X}\mathbf{w}$, $\mathbf{h} = \mathbf{H}\mathbf{w}$), which are subsequently divided by the sum of \mathbf{w} (weighted average). The additional set that includes specific COVID-19 terms is assigned a weight of 1. Our rationale is that people who experience COVID-19-related symptoms in addition to searching about them, will also issue queries about the disease as it is a broadly discussed topic. For the historical data, we divide their time span into yearly periods, and compute an average time series trend, \mathbf{h}_μ , using two standard deviations as upper and lower confidence intervals. We deploy a considerable historical period (8 years) and not a shorter and more recent one based on previous work on ILI rate modelling from online search data^{14,15} which indicated that the inclusion of an increased amount of past instances generally improves model accuracy. By doing this, we also capture search trends during past years that were affected by the H1N1 (swine flu), Ebola, and Zika virus epidemics.

Minimising the effect of news media using autoregression

On any given day the proportion of news articles about the COVID-19 pandemic is $m \in [0, 1]$, and the weighted score of symptom-related online searches (see previous paragraph) is equal to g ; we can apply a min-max normalisation so that $g \in [0, 1]$ as well. We hypothesise that g incorporates two signals based on infected (g_p) and concerned (g_c) users, respectively, i.e.

$$g = g_p + g_c. \quad (1)$$

Then, there exists a constant $\gamma \in [0, 1]$ such that

$$g_p = \gamma g \quad \text{and} \quad g_c = (1 - \gamma)g. \quad (2)$$

Our goal is to approximate γ using the observed variables m and g . The rationale of the deployed approach is similar to the logic behind a

Granger causality test²⁴. First, we train a linear AR model for forecasting the online search score at a time point (day) t , g_t , using its previous values; this is denoted by $\text{AR}(g)$. We also train a linear AR model with the same forecasting target, but an expanded space of observations that includes current (m_t) and previous (e.g. m_{t-1}) values of the news articles ratio; this is denoted by $\text{AR}(g, m)$. We then use the error ratio between the two models in forecasting g_t as an estimate of γ for time point t . In particular, we first solve

$$\arg \min_{\mathbf{w}, b_1} \frac{1}{N} \sum_{t=1}^N (g_t - w_1 g_{t-1} - w_2 g_{t-2} - b_1)^2, \quad (3)$$

to learn a pair of weights (\mathbf{w}) and an intercept term (b_1) for $\text{AR}(g)$. We use 2 lags (past values) to keep the complexity of the task tractable given the small amount of samples at our disposal, N . In our experiments we use the previous $N=56$ days, an intermediate value to account for the trade-off between the amount of samples and maintaining recency. We then solve

$$\arg \min_{\mathbf{w}, \mathbf{v}, b_2} \frac{1}{N} \sum_{t=1}^N (g_t - w_1 g_{t-1} - w_2 g_{t-2} - v_1 m_t - v_2 m_{t-1} - v_3 m_{t-2} - b_2)^2, \quad (4)$$

to learn the weights ($(\mathbf{w}; \mathbf{v})$) and an intercept term (b_2) for $\text{AR}(g, m)$. Using both models, we forecast the next (unseen) value of g , which is denoted by \hat{g}_{t+1} , and compute the absolute error from its known true value, g_{t+1} . This yields errors, ϵ_1 and ϵ_2 for $\text{AR}(g)$ and $\text{AR}(g, m)$, respectively. If $\epsilon_1 < \epsilon_2$, then the news media signal does not help to improve the accuracy of $\text{AR}(g)$, and hence we assume that it does not affect the online searches. Otherwise, we estimate its effect to be represented by

$$\gamma = \frac{\epsilon_2}{\epsilon_1}. \quad (5)$$

After obtaining a time series of γ 's for the all days in our analysis, we smooth each one of them using a harmonic mean (Supplementary Eq. 4) over the values of the previous 6 days (7 days in total including the day of focus).

Transferring supervised models of confirmed COVID-19 cases to different countries

Previous work has shown that it is possible to transfer a model for seasonal flu, based on online search query frequency time series, from one country that has access to historical syndromic surveillance data to another that has not²⁵. Here, we adapt this method to transfer a model for COVID-19 incidence from a source country where the disease spread has progressed significantly to a target country that is still in earlier stages of the epidemic. The main motivation for this is our assumption that a supervised model based on data from the source country might be able to capture disease dynamics sooner and therefore better than unsupervised or supervised models based solely on data from the target countries. The steps and data transformations that are required to apply this technique are detailed below.

Search query frequency time series are denoted by $\mathbf{S} \in \mathbb{R}_{>0}^{M \times n_s}$ and $\mathbf{T} \in \mathbb{R}_{>0}^{M \times n_t}$, for the source and target countries respectively; M denotes the number of days considered, and n_s, n_t the number of queries for the two locations. As these time series are quite volatile for some locations in our study, something that does not help in cross-location mapping of the data, we have smoothed them using a harmonic query frequency mean based on a window of $D=14$ past days (Supplementary Eq. 4). We subsequently train an elastic net model on data from the source location⁴¹, similarly to previous work on ILI^{14,15,17} or other text regression tasks^{42,43}. In particular, we solve the following optimisation task

$$\arg \min_{\mathbf{w}, b} (\|\mathbf{y} - \mathbf{S}\mathbf{w} - b\|_2^2 + \lambda_1 \|\mathbf{w}\|_1 + \lambda_2 \|\mathbf{w}\|_2^2), \quad (6)$$

where $\mathbf{y} \in \mathbb{R}^M$ denotes the daily number of confirmed COVID-19 cases in the source location, $\lambda_1, \lambda_2 \in \mathbb{R}_{>0}$ are the ℓ_1 - and ℓ_2 -norm regularisation parameters, and $\mathbf{w} \in \mathbb{R}^{n_s}$, $b \in \mathbb{R}$ denote the query weights and regression intercept, respectively. Prior to deploying elastic net, we apply a min-max normalisation on both \mathbf{S} and \mathbf{y} . We fix the ratio of λ 's, and then train q models for different values of λ_1 starting from the largest value that does not yield a null model and exploring the entire regularisation path. From all the different regression models represented by the columns of $\mathbf{W} \in \mathbb{R}^{n_s \times q}$, and the elements of $\mathbf{b} \in \mathbb{R}^q$, we use the ones that satisfy a sparsity level from 5.56% to 90.74% (3 to 49 search queries out of a total of 54) as an ensemble to accomplish a more inclusive transfer.

To generate an equivalent feature space for the target location (same dimensionality, similar feature attributes), we first establish query set pairs between the source and the target location using the symptom categories in the FF100 survey. We map a source query to the target query from the same symptom category that maximises their linear correlation based on their frequency time series. As the time series of source and target queries may not be temporally aligned, prior to computing correlations, we look at a window of $z = 45$ days (backwards and forwards) and identify at which temporal shift the average correlation between search query frequencies in \mathbf{S}' and \mathbf{T} is maximised; \mathbf{S}' here denotes a subset of \mathbf{S} that includes only the search queries that have been assigned a non zero weight by the elastic net (Eq. (6)). In the rare event that no active target search query exists for a certain symptom category, to maintain model consistency we use the best correlated one from all target queries available (irrespective of the symptom category) as its mapping. After this process, we end up with a subset $\mathbf{Z} \in \mathbb{R}^{M \times n_s}$ of the target feature space \mathbf{T} . Notably, \mathbf{Z} does not necessarily hold data for n_s distinct queries as different source queries may have been mapped to the same target query. \mathbf{Z} is subsequently normalised using min-max. To reduce the distance in the query frequency range between the two feature spaces (\mathbf{S}, \mathbf{Z}), we scale the latter based on their mean, column-wise (per search query) ratio $\mathbf{r} \in \mathbb{R}_{>0}^{n_s}$, i.e. $\mathbf{Z}_S = \mathbf{Z} \odot \mathbf{r}$. We then deploy the ensemble source models to the target space, making multiple inferences (for different λ_1 values) held in $\mathbf{Y} \in \mathbb{R}^{n_s \times q}$:

$$\mathbf{Y} = \mathbf{Z}_S \mathbf{W} + \mathbf{b}. \quad (7)$$

We reverse the min-max normalisation for each one of the inferred time series (columns of \mathbf{Y}) using values from the source model's ground truth \mathbf{y} (prior to its normalisation). Finally, we compute the mean of the ensemble (across the rows of \mathbf{Y}) as our target estimate, and also use the 0.025 and 0.975 quantiles to form 95% confidence intervals.

Correlation and regression analysis

The relationship of search frequency time series and confirmed cases or deaths can uncover symptoms or behaviours related to COVID-19. However, as it is hard to find resources that are representative of community level disease rates, looking at this relationship separately for each country might produce misleading outcomes. To mitigate this to the extent possible, we combine data from C countries and produce an aggregate set of query frequencies, $\mathbf{Z}_a \in \mathbb{R}^{CM \times n}$, where M, n denote the considered days and search queries, respectively. We denote the aggregated daily confirmed COVID-19 cases or deaths for these countries with $\mathbf{y}_a \in \mathbb{R}^{CM}$. Prior to the aggregation, we apply min-max normalisation on the query frequency, confirmed cases, and deaths time series separately for each country to balance out local properties. Initially, we compute the linear correlation between the columns of \mathbf{Z}_a and \mathbf{y}_a . Correlation is an informative metric, but considers each search query in isolation. Therefore, we also perform a multivariate regression analysis to more rigorously estimate the impact of each search query in estimating \mathbf{y}_a . To do this, we apply elastic net regularised regression (Eq. (6)), training and testing K models, one for each day of an identified test period of interest. During each of the K training phases, assuming it contains η days, we use data up to and including the past $\eta - 1$ days to train, and test only on the last day (η_{th}) which is unseen; this results in daily test sets of size C (one value for each country). We explore elastic net's regularisation path to consider L models that maintain (by assigning a nonzero weight) up to a reasonable percentage of the features (e.g. 50%), so that a solution is not overfitting. We do this gradually, selecting first 1% of the features and moving towards the maximum considered percentage. In this experiment, we use the test set to identify the most accurate (in terms of mean squared error) model at each density level. For this model, we determine the impact of each one of the features (search queries) by considering both its frequency and allocated weight. The impact $\Theta(\cdot)$ of a query q is equal to

$$\Theta(q) = \left(\sum_{\ell=1}^L \sum_{t=1}^K \sum_{j=1}^C f_{t,j} w_{\ell,t} \right) / \left(\sum_{\ell=1}^L \sum_{t=1}^K \sum_{j=1}^C \hat{y}_{\ell,t,j} \right), \quad (8)$$

where $f_{t,j}$ denotes the query frequency at time point (day) t and for country j , $w_{\ell,t}$ the corresponding weight at sparsity level ℓ , and $\hat{y}_{\ell,t,j}$ the respective estimated confirmed cases or deaths. Impacts are summed across all the considered days, and model densities, and normalised at the end by the sum of all the corresponding estimates.

Short-term forecasting of COVID-19 deaths

Moving further into supervised models, we develop forecasting solutions for COVID-19 to assess the potential value of incorporating web search data in AR models. In this case, we focus on modelling related deaths as opposed to confirmed cases because the reporting of the latter has been dependent on testing approaches, and thus has a less reliable time series structure from a statistical perspective. Short-term forecasts of COVID-19 deaths can be conducted using the time series of past records. Augmenting this AR signal with online user trails could help to improve accuracy, if we draw a parallel with ILL rate modelling⁴⁴. Let $\mathbf{z} \in \mathbb{R}_{>0}^M$ denote the average frequency of N min-max normalised search queries for M days, and let $\mathbf{y} \in \mathbb{N}^M$ be the corresponding COVID-19 deaths. We first solve a strictly AR task using L past values, meaning that at a time point t we use $\mathbf{y}_{AR}(t, L) = [y_t, y_{t-1}, \dots, y_{t-L}] \in \mathbf{y}$ to forecast y_{t+D} , performing D days ahead forecasting. We denote this forecasting model as AR-F. We also augment our observations by incorporating search query frequency data for the same time window resulting to an input $\mathbf{x}_{AR}(t, L) = [\mathbf{z}_{AR}(t, L); \mathbf{y}_{AR}(t, L)]$ that is held in the concatenated matrix $\mathbf{X} = [\mathbf{Z}_{AR}(L); \mathbf{Y}_{AR}(L)] \in \mathbb{R}_{>0}^{M \times 2(L+1)}$. We denote this search AR forecasting model as SAR-F. For simplicity, we drop the subscript notation of the input variables ($\mathbf{x}, \mathbf{z}, \mathbf{y}$) in subsequent references to them.

We develop forecasting models using Gaussian Processes (GP), training a different model for each D days ahead forecasting task. The choice of GPs is justified by previous work on modelling infectious diseases^{14,15,45}. GPs are defined as random variables any finite number of which have a multivariate Gaussian distribution⁴⁶. GP methods aim to learn a function $f: \mathbb{R}^m \rightarrow \mathbb{R}$ drawn from a GP prior. They are specified through a mean and a covariance (or kernel) function, i.e. $f(\mathbf{x}) \sim \text{GP}(\mu(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$, where \mathbf{x} and \mathbf{x}' (both $\in \mathbb{R}_{>0}^{2(L+1)}$) denote rows of the input matrix \mathbf{X} for the SAR-F model. By setting $\mu(\mathbf{x}) = 0$, a common practice in GP modelling, we focus only on the kernel function. The specific kernel function used in SAR-F is given by

$$k(\mathbf{x}, \mathbf{x}') = k_{SE}(\mathbf{z}, \mathbf{z}'; \sigma_1, \ell_1) + k_{SE}(\mathbf{y}, \mathbf{y}'; \sigma_2, \ell_2) + k_{SE}(\mathbf{x}, \mathbf{x}'; \sigma_3, \ell_3) + \sigma_4^2 \delta(\mathbf{x}, \mathbf{x}'), \quad (9)$$

where $k_{SE}(\cdot, \cdot)$ denotes a squared exponential (SE) covariance function (Supplementary Eq. 5), $\delta(\cdot, \cdot)$ denotes a Kronecker delta function used for an independent noise component, ℓ 's and σ 's are lengthscale and scaling (variance) parameters, respectively. This composite covariance function uses separate kernels for online search and deaths data, as well as a kernel where they are modelled jointly. This provides the GP with more flexibility in adapting its decisions to one class of observations over the other. The covariance function of the AR-F model is a simplified version of Eq. (9), where search data is not used, i.e.

$$k(\mathbf{y}, \mathbf{y}') = k_{SE}(\mathbf{y}, \mathbf{y}'; \sigma_1, \ell_1) + k_{SE}(\mathbf{y}, \mathbf{y}'; \sigma_2, \ell_2) + \sigma_3^2 \delta(\mathbf{y}, \mathbf{y}'). \quad (10)$$

The hyperparameters of both covariance functions (σ 's, ℓ 's) are optimised using Gaussian likelihood and variational inference⁴⁶.

In our analysis, we compare the aforementioned forecasting models with a basic persistence model (PER-F). For a D days ahead forecasting task, PER-F uses the most recently observed value of the ground truth, y_t , as the forecasting estimate for the time instance $t + D$ ($\hat{y}_{t+D} = y_t$). Given the limited time span and the irregularities in reporting, this is a competitive baseline to improve upon.

Data sets

Our analysis focuses on a set of 8 countries, namely the US, UK, England, Australia, Canada, France, Italy, Greece, and South Africa. This choice of countries covers various geographical locations across the world (Europe, North America, Australia, Africa), and allowed the curation of COVID-19-related search and news media terms from native speakers with a good understanding of this research project. Given the various logistic constraints related to data collection and translation from experts, we did not expand the analysis to other countries. However, it is expected that outcomes should replicate to other countries with similar socioeconomic and cultural characteristics. China was not included in our study because data from the Google search engine is not considered as a representative sample due to its modest market share (wikipedia.org/wiki/Google_China).

Online search data is obtained from Google Health Trends, a non public application interface offered by Google for research on health-related topics. Data represent daily online search query frequencies for specific areas of interest. Query frequencies are defined as the sum of search sessions that include a target search term divided by the total number of search sessions (for a day and area of interest). Hence, in the paper when

we refer to a certain search query, we also refer to all other search queries that include all of its words. Google defines a search session as a grouping of consecutive searches by the same user within a short time interval. We have obtained data from September 30, 2011 to May 24, 2020 for the aforementioned 8 countries and separately for the nation of England. The list of search terms is determined by COVID-related symptoms and keywords. For each country, we mainly used queries in its native language(s). Each daily search query frequency is smoothed using a harmonic mean over the past $D=14$ days including itself (Supplementary Eq. 4).

We are also using a global news corpus to extract media coverage trends for COVID-19 in all the countries in our study. This is estimated by counting the proportion of articles mentioning a COVID-19-related term (see SI for the list of terms). In particular, daily counts of the total of news media articles published, and the subset that included at least one relevant keyword anywhere in the body of the text were collected from the Media Cloud database (mediacloud.org) for the UK (93), US (225), Australia (61), Canada (79), France (360), Italy (178), Greece (75), and South Africa (135), where in the parentheses we state the number of media sources considered per country. These counts were collected from September 30, 2019 through May 24, 2020. Within this time-span we identified 2,535,735 articles as COVID-19-related from a total of 10,093,349. Supplementary Fig. 11 depicts the average daily ratio across all countries, as soon as it started being above zero (Jan. 2020), with two standard deviations as confidence intervals. Initially, there exists a distinctive pattern of (exponential) increase, but from April onwards the average media coverage has a moderate decreasing trend. In addition, we observe a certain variance across locations or time periods, that adds to the potential value of this signal.

For all countries in our analysis, we obtained daily confirmed COVID-19 cases and deaths time series from the European Centre for Disease Prevention and Control (ECDC). Links are provided in the SI.

Ethics

Ethical approval was not required for the analysis presented in this paper. Data was obtained in an anonymised and aggregated (national level) format.

Reporting summary

Further information on research design is available in the Nature Research Reporting Summary linked to this article.

DATA AVAILABILITY

The online search and RCGP swabbing scheme data sets that support the findings of this study are available from Google and RCGP/PHE, respectively. Restrictions apply to the availability of these data sets, which were used under license for the current study, and so are not publicly available. These data sets are however available from the authors upon reasonable request and with the respective permission of Google and RCGP/PHE. The rest of the data sets that this study is based on are available at figshare.com/projects/Tracking_COVID-19_using_online_search/81548.

CODE AVAILABILITY

We have used standard and contributed MATLAB[®] (version R2019a) functions or libraries to conduct experiments. In particular, we have used function `regress` for linear regression without regularisation, function `lasso` for linear regression with elastic net regularisation, function `getest` for the Granger causality analysis, and the Gaussian Process Regression and Classification Toolbox (version 4.2, gpmll/code/matlab/doc/) for deriving GP models.

Received: 2 July 2020; Accepted: 24 December 2020;

Published online: 08 February 2021

REFERENCES

- Polgreen, P. M., Chen, Y., Pennock, D. M., Nelson, F. D. & Weinstein, R. A. Using internet searches for influenza surveillance. *Clin. Infect. Dis.* **47**, 1443–1448 (2008).
- Ginsberg, J. et al. Detecting influenza epidemics using search engine query data. *Nature* **457**, 1012–1014 (2009).
- Eysenbach, G. Infodemiology and infoveillance: framework for an emerging set of public health informatics methods to analyze search, communication and publication behavior on the internet. *J. Med. Internet Res.* **11**, <https://doi.org/10.2196/jmir.1157> (2009).
- Lampos, V. & Cristianini, N. Tracking the flu pandemic by monitoring the Social Web. In *Proc. of the 2nd International Workshop on Cognitive Information Processing*, 411–416, <https://doi.org/10.1109/CIP.2010.5604088> (2010).
- Culotta, A. Towards detecting influenza epidemics by analyzing Twitter messages. In *Proc. of the 1st Workshop on Social Media Analytics*, 115–122, <https://doi.org/10.1145/1964858.1964874> (2010).
- Choudhury, M. D., Gamon, M., Counts, S. & Horvitz, E. Predicting depression via social media. In *Proc. of the 7th International AAAI Conference on Weblogs and Social Media*, 128–137 (2013).
- Yom-Tov, E. *Crowdsourced health: How what you do on the Internet will improve medicine* (MIT Press, 2016).
- Wagner, M., Lampos, V., Cox, I. J. & Pebody, R. The added value of online user-generated content in traditional methods for influenza surveillance. *Sci. Rep.* **8**, <https://doi.org/10.1038/s41598-018-32029-6> (2018).
- Chew, C. & Eysenbach, G. Pandemics in the Age of Twitter: Content Analysis of Tweets during the 2009 H1N1 Outbreak. *PLoS ONE* **5**, <https://doi.org/10.1371/journal.pone.0014118> (2010).
- Lampos, V. & Cristianini, N. Nowcasting events from the social web with statistical learning. *ACM Trans. Intell. Syst. Technol.* **3**, 1–22 (2012).
- Olson, D. R., Konty, K. J., Paladini, M., Viboud, C. & Simonsen, L. Reassessing google flu trends data for detection of seasonal and pandemic influenza: a comparative epidemiological study at three geographic scales. *PLOS Comp. Biol.* **9**, <https://doi.org/10.1371/journal.pcbi.1003256> (2013).
- Lazer, D., Kennedy, R., King, G. & Vespignani, A. The parable of google flu: traps in big data analysis. *Science* **343**, 1203–1205 (2014).
- Yang, S., Santillana, M. & Kou, S. C. Accurate estimation of influenza epidemics using google search data via ARGO. *PNAS* **112**, 14473–14478 (2015).
- Lampos, V., Miller, A. C., Crossan, S. & Stefansen, C. Advances in nowcasting influenza-like illness rates using search query logs. *Sci. Rep.* **5**, <https://doi.org/10.1038/srep12760> (2015).
- Lampos, V., Zou, B. & Cox, I. J. Enhancing feature selection using word embeddings: the case of flu surveillance. In *Proc. of the 26th International Conference on World Wide Web*, 695–704, <https://doi.org/10.1145/3038912.3052622> (2017).
- Reich, N. G. et al. A collaborative multiyear, multimodel assessment of seasonal influenza forecasting in the United States. *PNAS* **116**, 3146–3154 (2019).
- Lampos, V., Yom-Tov, E., Pebody, R. & Cox, I. J. Assessing the impact of a health intervention via user-generated internet content. *Data Min. Knowl. Discov.* **29**, 1434–1457 (2015).
- Roser, M., Ritchie, H., Ortiz-Ospina, E. & Hasell, J. Coronavirus pandemic (COVID-19). *Our World in Data*. <https://ourworldindata.org/coronavirus> (2020).
- Lipsitch, M., Swerdlow, D. L. & Finelli, L. Defining the epidemiology of covid-19: studies needed. *NEJM* **382**, 1194–1196 (2020).
- Baumgartner, S. E. & Hartmann, T. The role of health anxiety in online health information search. *Cyberpsychol. Behav. Soc. Netw.* **14**, 613–618 (2011).
- Qiu, J. et al. A nationwide survey of psychological distress among Chinese people in the COVID-19 epidemic: implications and policy recommendations. *Gen. Psychiatry* **33**, <https://doi.org/10.1136/gpsych-2020-100213> (2020).
- Santillana, M. et al. Combining Search, Social Media, and Traditional Data Sources to Improve Influenza Surveillance. *PLoS Comp. Biol.* **11**, <https://doi.org/10.1371/journal.pcbi.1004513> (2015).
- Boddington, N. L. et al. COVID-19 in Great Britain: epidemiological and clinical characteristics of the first few hundred (FF100) cases: a descriptive case series and case control analysis. [Preprint]. *Bull. World Health Organ.* <https://doi.org/10.2471/BLT.20.265603> (2020).
- Granger, C. W. J. Investigating causal relations by econometric models and cross-spectral methods. *Econometrica* **37**, 424–438 (1969).
- Zou, B., Lampos, V. & Cox, I. J. Transfer Learning for Unsupervised Influenza-like Illness Models from Online Search Data. In *Proc. of the 28th International Conference on World Wide Web*, 2505–2516, <https://doi.org/10.1145/3308558.3313477> (2019).
- Carfi, A., Bernabei, R., Landi, F. & for the Gemelli Against COVID-19 Post-Acute Care Study Group. Persistent symptoms in patients after acute COVID-19. *JAMA* **324**, 603–605 (2020).
- Lau, H. et al. The positive impact of lockdown in Wuhan on containing the COVID-19 outbreak in China. *J. Travel Med.* **27**, <https://doi.org/10.1093/jtm/taaa037> (2020).
- Hopman, J., Allegranzi, B. & Mehtar, S. Managing COVID-19 in low- and middle-income countries. *JAMA* **323**, 1549–1550 (2020).
- Gilbert, M. et al. Preparedness and vulnerability of African countries against importations of COVID-19: a modelling study. *The Lancet* **395**, 871–877 (2020).
- Mahler, D. G., Montes, J. & Newhouse, D. L. Internet Access in Sub-Saharan Africa. Tech. Rep., The World Bank. <http://documents1.worldbank.org/curated/en/518261552658319590/pdf/Internet-Access-in-Sub-Saharan-Africa.pdf> (2019).

31. Remuzzi, A. & Remuzzi, G. COVID-19 and Italy: what next? *The Lancet* **395**, 1225–1228 (2020).
32. Grasselli, G., Pesenti, A. & Cecconi, M. Critical care utilization for the COVID-19 outbreak in Lombardy, Italy: early experience and forecast during an emergency response. *JAMA* **323**, 1545–1546 (2020).
33. Zhou, F. et al. Clinical course and risk factors for mortality of adult inpatients with COVID-19 in Wuhan, China: a retrospective cohort study. *The Lancet* **395**, 1054–1062 (2020).
34. Sohrabi, C. et al. World Health Organization declares global emergency: a review of the 2019 novel coronavirus (COVID-19). *Int. J. Surg.* **76**, 71–76 (2020).
35. Hale, T. et al. Variation in Government Responses to COVID-19. *Blavatnik School of Government Working Paper*. <https://www.bsg.ox.ac.uk/sites/default/files/2020-12/BSG-WP-2020-032-v10.pdf> (2020).
36. de Lusignan, S. et al. Emergence of a Novel Coronavirus (COVID-19): Protocol for Extending Surveillance Used by the Royal College of General Practitioners Research and Surveillance Centre and Public Health England. *JMIR Public Health Surveill.* **6**, <https://doi.org/10.2196/18606> (2020).
37. Jarvis, C. I. et al. Quantifying the impact of physical distance measures on the transmission of COVID-19 in the UK. *BMC Med.* **18**, <https://doi.org/10.1186/s12916-020-01597-8> (2020).
38. Li, C. et al. Retrospective analysis of the possibility of predicting the COVID-19 outbreak from Internet searches and social media data, China, 2020. *Euro-surveillance* **25** (2020).
39. Effenberger, M. et al. Association of the COVID-19 pandemic with Internet Search Volumes: A Google Trends™ Analysis. *Int. J. Infect. Dis.* **95**, 192–197 (2020).
40. Sousa-Pinto, B. et al. Assessment of the impact of media coverage on COVID-19-related google trends data: Infodemiology Study. *J. Med. Internet Res.* **22**, <https://doi.org/10.2196/19611> (2020).
41. Zou, H. & Hastie, T. Regularization and variable selection via the elastic net. *J. Royal Stat. Soc.: Series B* **67**, 301–320 (2005).
42. Lampos, V., Preotjuc-Pietro, D. & Cohn, T. A user-centric model of voting intention from social media. In *Proc. of the 51st Annual Meeting of the Association for Computational Linguistics*, 993–1003 (2013).
43. Lampos, V., Aletras, N., Preotjuc-Pietro, D. & Cohn, T. Predicting and Characterising User Impact on Twitter. In: *Proc. of the 14th Conference of the European Chapter of the Association for Computational Linguistics*, 405–413, <https://doi.org/10.3115/v1/E14-1043> (2014).
44. Paul, M. J., Dredze, M. & Broniatowski, D. Twitter Improves Influenza Forecasting. *PLoS Curr.* **6**, <https://doi.org/10.1371/currents.outbreaks.90b9ed0f59bae4ccaa683a39865d9117> (2014).
45. Zou, B., Lampos, V. & Cox, I. Multi-Task Learning Improves Disease Models from Web Search. In *Proc. of the 27th International Conference on World Wide Web*, 87–96, <https://doi.org/10.1145/3178876.3186050> (2018).
46. Rasmussen, C. E. & Williams, C. K. I. *Gaussian Processes for Machine Learning* (MIT Press, 2006).

ACKNOWLEDGEMENTS

V.L., S.M., R.A.M., and I.J.C. would like to acknowledge all levels of support from the EPSRC projects “EPSRC IRC in Early-Warning Sensing Systems for Infectious Diseases” (EP/K031953/1), “i-sense: EPSRC IRC in Agile Early Warning Sensing Systems for Infectious Diseases and Antimicrobial Resistance” and its COVID-19 plus award “EPSRC i-sense COVID-19: Harnessing digital and diagnostic technologies for COVID-19” (EP/R00529X/1). V.L. and I.J.C. would also like to acknowledge the support from the MRC/NIHR project “COVID-19 Virus Watch: Understanding community incidence,

symptom profiles, and transmission of COVID-19 in relation to population movement and behaviour” (MC_PC_19070) as well as from a Google donation funding the project “Modelling the prevalence and understanding the impact of COVID-19 using web search data”. The authors would like to thank the NHS/PHE FF100 team for sharing early results, and the RCGP for sharing swabbing data for COVID-19. We also appreciate the contribution of Ettore Severi, Anna Odone, and Daniela Paolotti in the translation of search queries from English to Italian. Finally, V.L. would like to thank Sam J. Gilbert for interesting discussions and pointers during the development of this work.

AUTHOR CONTRIBUTIONS

V.L. conceived this research project, formed the majority of the data sets, conceived and developed the methods, ran the experiments, wrote the first drafts of this manuscript, and led the write-up thereafter. M.S.M. provided news coverage data for all countries in our analysis. M.E. provided a translation of search query groups from English to Italian, S.M. from English to French, and M.X.R. and Y.H. from English to various languages spoken in South Africa. M.S.M., E.Y.-T., M.E., S.M., Y.H., and I.J.C. contributed in the write-up of the manuscript. All authors provided feedback during the development of this work.

COMPETING INTERESTS

The authors declare no competing interests.

ADDITIONAL INFORMATION

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41746-021-00384-w>.

Correspondence and requests for materials should be addressed to V.L.

Reprints and permission information is available at <http://www.nature.com/reprints>

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



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 license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license 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 license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2021