# Our journey in digital epidemiology: modelling infectious diseases using online search activity

Department of Computer Science University College London (UCL)



## Vasileios Lampos



#### A. Estimating flu prevalence using web search activity

- (12760). doi:10.1038/srep12760
- doi:10.1145/3038912.3052622

#### **Transferring** a disease model from one country to another using web search activity B.

(WWW). doi:10.1145/3308558.3313477

#### Modelling COVID-19 prevalence using web search activity C.

#### D. Advanced models (neural network architectures) for disease prevalence forecasting

- trends. PLOS Computational Biology **19** (8). doi.org/10.1371/journal.pcbi.1011392
- doi.org/10.48550/arXiv.2406.07438

Modelling infectious diseases using online search

Lampos, Miller, Crossan, Stefansen (2015). Advances in nowcasting influenza-like illness rates using search query logs. Scientific Reports 5

Lampos, Zou, Cox (2017). Enhancing feature selection using word embeddings: The case of flu surveillance. Web Conference (WWW).

Zou, Lampos, Cox (2019). Transfer learning for unsupervised influenza-like illness models from online search data. Web Conference

Lampos et al. (2021). Tracking COVID-19 using online search. npj Digital Medicine 4 (17). doi:10.1038/s41746-021-00384-w

Morris, Hayes, Cox, Lampos (2023). Neural network models for influenza forecasting with associated uncertainty using Web search activity

Shu, Lampos (2024). DEFORMTIME: Capturing variable dependencies with deformable attention for time series forecasting. arXiv Preprint.





# Estimating flu prevalence using web search activity

Lampos *et al.* (2015), *Sci. Rep.* Lampos, Zou, Cox (2017), WWW '17

Modelling infectious diseases using online search

# Part A



### From web searches to influenza (flu) rates

**50000** 

#### flu treatment

flu treatment kids flu treatment otc flu treatment natural flu treatment medication flu treatment toddler



Eysenbach (2006), AMIA; Polgreen et al. (2008), Clin. Infect. Dis.; Ginsberg et al. (2009), Nature



#### **Complements** conventional syndromic surveillance systems

- Iarger cohort
- broader demographic coverage
- more granular geographic coverage
- not affected by closure days (weekends, holidays)
- ► timeliness
- Iower cost
- Applicable to locations that lack an established health surveillance infrastructure
- Track **novel** infectious diseases

confirmed infections, associated hospitalisations or deaths.

#### Why estimate disease rates from web search?

- Conventional (traditional) syndromic surveillance methods: disease prevalence, i.e. the % of infected people in a population, is determined via doctor (GP) visits and other related indicators, such as laboratory-
  - Wagner et al. (2018), Sci. Rep.; Budd et al. (2020), Nat. Med.
  - Modelling infectious diseases using online search







### Google Flu Trends (GFT) – discontinued

#### google.org Flu Trends

#### Google.org home

#### Dengue Trends

Flu Trends

Home

Select country/regior 🗘

#### How does this work?

<u>FAQ</u>

#### Flu activity Intense High Moderate Low Minimal

#### Explore flu trends around the world

We've found that certain search terms are good indicators of flu activity. Google Flu Trends uses aggregated Google search data to estimate flu activity. Learn more »



Language: English (United States)

-

#### Ginsberg et al. (2009), Nature



- P: percentage of doctor visits due to influnza-like illness (ILI)
- $\beta_0$ : regression intercept (bias)
- $\beta_1$ : regression weight (univariate regression)
  - $\epsilon$ : independent, zero-centered noise

#### Main issue

What if some of the selected queries are spurious or, in general, relate differently to flu rates compared to other selected search queries? This model makes a very naïve assumption.

Modelling infectious diseases using online search

#### Google Flu Trends (GFT) – regression function

 $logit(P) = \beta_0 + \beta_1 \times logit(Q) + \epsilon$ 

logit (x) = 
$$\ln\left(\frac{x}{1-x}\right)$$
  
where  $x \in (0,1)$ 

Q : aggregate frequency of a set of automatically selected search queries related to ILI

Ginsberg et al. (2009), Nature





### Google Flu Trends (GFT) – *shortcomings*





### Web search frequencies & flu rates: a *nonlinear* relationship



- Not all search queries have a linear relationship with flu rates
- Not all queries relate to flu rates in the same way either
- The same query may also have a bi-modal relationship which may fluctuate in time
- Very hard to model this using a single weight!

Lampos et al. (2015), Sci. Rep.





#### Search query frequency time series

 $x_{ij} \in \mathbf{X} = \frac{\text{number of times query } j \text{ is issued during time step } i}{\text{total number of searches during time step } i}$ 

#### **Disease rates**

 $\mathbf{y} \in \mathbb{R}^n_{>0}$ where *n* is the number of time steps (days or weeks)

#### **Regression task**

 $f: \mathbf{X} \to \mathbf{y}$  linear case:  $f = \{ \mathbf{w} \in \mathbb{R}^m, \beta \in \mathbb{R} \}$  $\hat{\mathbf{y}}_t = \mathbf{X}_t \mathbf{w} + \beta$  assuming  $\mathbf{X}_t \in \mathbb{R}^{k \times m}$  holds k unseen test samples

Modelling infectious diseases using online search

- $X \in \mathbb{R}_{>0}^{n \times m}$  where *n* is the number of time steps / samples (days or weeks) and m is the number of search queries / variables (m > 1000)



10

## Multivariate Gaussian Process (GP) kernels on search query clusters

#### Composite Gaussian Process (GP) kernel

$$k(\mathbf{x}, \mathbf{x}') = \left(\sum_{i=1}^{C} k_{\text{SE}}\left(\mathbf{c}_{i}, \mathbf{c}_{i}'\right)\right) + \sigma_{n}^{2} \cdot \delta(\mathbf{x}, \mathbf{x}')$$

# $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^m_{>0}$ , where *m* is the number of search queries we consider

#### Squared Exponential (SE) kernel

$$k_{\text{SE}}(\mathbf{c}_i, \mathbf{c}'_i) = \sigma^2 \exp\left(-\frac{\|\mathbf{c}_i - \mathbf{c}'_i\|_2^2}{2\ell^2}\right)$$

Lampos et al. (2015), Sci. Rep.; Rasmussen, Williams (2006), MIT Press

Modelling infectious diseases using online search

**NB:** Queries are selected based on their **correlation** with ILI rates in the training data and an elastic net regression function

 $\mathbf{c}_i, \mathbf{c}'_i \in \mathbb{R}^z_{>0}, z < m, C$  query clusters based on frequency time series



11

### Modelling ILI rates with Gaussian Process (GP) kernels





## Modelling ILI rates with Gaussian Process (GP) kernels



.95 bivariate correlation (previously .89) with CDC rates

Modelling infectious diseases using online search

# 42% mean absolute error reduction compared to Google Flu Trends



# Modelling ILI rates with Gaussian Process (GP) kernels



Modelling infectious diseases using online search

#### 42% mean absolute error reduction compared to Google Flu Trends .95 bivariate correlation (previously .89) with CDC rates



### Autoregression (AR) with SARIMAX

$$y_{t} = \underbrace{\sum_{i=1}^{p} \phi_{i} y_{t-d}}_{\text{AR and seasonal AR}} + \underbrace{\sum_{i=1}^{J} \omega_{i} y_{t-52-i}}_{\text{MA and seasonal MA}} + \underbrace{\sum_{i=1}^{K} \nu_{i} \epsilon_{t-52-i}}_{\text{GP estimates}} + \underbrace{\sum_{i=1}^{D} w_{i} h_{t,i}}_{\text{GP estimates}} + \epsilon_{t}$$

- SARIMAX: Seasonal AutoRegressive Integrated Moving Average with eXogenous variables
- d weeks delay in including past ILI rates as reported by CDC
- Choose model parameters based on the Akaike Information Criterion (AIC)
  - sometimes past seasonal trends are helpful, but not always
  - the most important piece of information is the GP estimate for the ILI rate based on web search query frequencies

Lampos *et al.* (2015), *Sci. Rep.* 







### Modelling ILI rates with Gaussian Process (GP) kernels & SARIMAX



- .99 bivariate correlation with CDC

Modelling infectious diseases using online search

#### Incorporating historical CDC estimates into an autoregression (AR) using SARIMAX 27% MAE reduction compared to GFT with AR, 52% over the GP model without AR



- Feature selection was based on a temporal relationship
  - Is this sufficient? No / not always
- Spurious search queries such as "NBA injury report" or "muscle building supplements" were still included in the selection
  - query clustering: some guarantees for different treatment, but needs a more complex regression model
- Introduce a query filter based on distributional semantics using word embeddings
- Hybrid combination of this with previous feature selection regimes

Lampos et al. (2015), Sci. Rep.; Lampos, Zou, Cox (2017), WWW '17

Modelling infectious diseases using online search

17

### Query selection based on distributional semantics

# $\sin\left(q,\mathbb{C}\right) = -\frac{1}{\Sigma}$

 $\mathbf{e}_{(\cdot)}$ : embedding vector trained on Twitter data  $\mathbb{C} = \{\mathbb{C}_P, \mathbb{C}_N\} - a \text{ concept about influenza}$  $\mathbb{C}_P$ : phrases of a positive context for concept  $\mathbb{C}$  $\mathbb{C}_N$  : phrases of a negative context for concept  $\mathbb{C}$  $\theta = \cos(\cdot) \rightarrow \in [0,1]$  via  $(\theta + 1)/2$  to avoid negative components  $\gamma \in \mathbb{R}_{>0}$  to avoid, in theory, division by 0

$$\sum_{i=1}^{P} \cos\left(\mathbf{e}_{q}, \mathbf{e}_{p_{i}}\right)$$

$$\sum_{j=1}^{N} \cos\left(\mathbf{e}_{q}, \mathbf{e}_{n_{j}}\right) + \gamma$$

Lampos, Zou, Cox (2017), WWW '17; Levy, Goldberg (2014), CoNLL '14



### Query selection based on distributional semantics

Positive context	Negative context	Most similar queries
#flu fever flu flu medicine GP hospital	Bieber ebola Wikipedia	"cold flu medicine" "flu aches" "cold and flu" "cold flu symptoms" "colds and flu"
flu flu GP flu hospital flu medicine	ebola Wikipedia	"flu aches" "flu" "colds and flu" "cold and flu" "cold flu medicine"

Modelling infectious diseases using online search

Lampos, Zou, Cox (2017), WWW '17





Modelling infectious diseases using online search

Feature selection based on correlation and regularised regression



### Feature selection based on correlation and regularised regression



#### Examples of problematic query selections

prof. *surname*: 70% name surname: 27% heating oil: 21%

Modelling infectious diseases using online search

name surname recipes: 21% blood game: 12.3% swine flu vaccine side effects: 7.2%

21

### Hybrid feature selection: distributional semantics and correlation



- 12.3% accuracy improvement in terms of mean absolute error .913 bivariate correlation with the ground truth (RCGP ILI rates)



# Flu detector, part of UK's influenza surveillance



gov.uk/government/statistics/ national-flu-and-covid-19surveillance-reports-2023to-2024-season





# Why estimate disease rates from web search?

- **Complements** conventional syndromic surveillance systems
  - larger cohort
  - broader demographic coverage
  - broader, more granular geographic coverage
  - not affected by closure days and other temporal biases
  - ► timeliness
  - ► lower cost
- Track novel infectious diseases

Conventional (traditional) syndromic surveillance methods: disease prevalence, i.e. the % of infected people in a population, is determined via doctor (GP) visits and other related indicators, such as laboratoryconfirmed infections, associated hospitalisations or deaths.

Modelling infectious diseases using online search

#### Applicable to locations that lack an established health surveillance infrastructure

Wagner et al. (2018), Sci. Rep.; Budd et al. (2020), Nat. Med.





# Transferring a disease model from one country to another using web search activity

Modelling infectious diseases using online search

# Part B

Zou, Lampos, Cox (2019), WWW '19



- Transfer learning in general
  - Gain knowledge from a domain/task, and then apply it to another one
- Transfer learning for estimating flu rates across different countries
  - Locations: source (no missing data), target (no disease rates)
  - regularised regression model for a source location based on web search activity and historical disease rates
  - map search queries from the source to the target location - semantic similarity (bilingual if necessary)

    - temporal similarity
  - hybrid similarity (their linear combination controlled by  $\gamma$ ) transfer regression model (equivalent to zero-shot learning)



### Transferring a flu model based on web searches: from US to France



How similar are the flu rates between the US and France (FR)?

Modelling infectious diseases using online search

# - temporal differences (e.g. different onset/peak moments), intensity differences



### Transferring a flu model based on web searches: from US to France





### Transferring a flu model based on web searches: from US to Australia



How similar are the flu rates between the US and Australia (AU)? — different (≈opposite) seasons, significant intensity differences in more recent years



### Transferring a flu model based on web searches: from US to Australia





# Modelling COVID-19 prevalence using web search activity

Lampos et al. (2021), npj Digit. Med.

Modelling infectious diseases using online search

# Part C



# $x_{L,d}$

**Unprecedented** search frequency trends during the first COVID-19 pandemic waves



Modelling infectious diseases using online search

Google Health Trends: frequency  $x_{L,d}$  of web search query q for a location L during a day d

number of times q was issued by users in location L during day d

total number of searches by users in location L during day d





## Challenges in modelling COVID-19 using web search activity

- No reliable and not enough ground truth data
  - Supervised learning no longer possible can we use transfer learning?
  - Evaluation of any model will be problematic
- Unsupervised learning
  - Which search queries to use?

  - and media coverage rather than by infection?

#### How do we know our model is related to COVID-19 and not other infectious diseases?

How do we know our signal is not affected by other factors such as concern, curiosity,



# First few hundred (FF100) patient survey (NHS & UKHSA)



cough fatigue fever headache muscle ache appetite loss shortness of breath sore throat joint ache runny nose loss of the sense of smell diarrhoea sneezing nausea vomiting altered consciousness nose bleed rash seizure

Probability of occurrence in COVID-19 patients

Modelling infectious diseases using online search

	1			
		1		
_				-
L (	ר א ה	ר אר ר	רן 🔿	
			<i>.</i> , 0	

Boddington *et al.* (2021), *Bull. WHO* 



- cough: cough, coughing
- fatigue: fatigue
- fever: chills, fever, high temp fever, high temperature
- headache: head ache, headache, headaches, migraine
- muscle ache: muscle ache, muscular pain
- appetite loss: appetite loss, loss of appetite, lost appetite
- shortness of breath: breathing difficulties, breathing difficulty, cant breathe, shortness of breath, short breath
- loss of the sense of smell: anosmia, loss of smell, loss smell
- COVID-19 terms: coronavirus, covid, covid-19, covid19



- cough: tosse, tossire
- fatigue: affaticamento, fatica, spossatezza, stanchezza
- fever: alta temperatura, brividi, febbre
- headache: emicrania, mal di testa
- muscle ache: dolore muscolare, dolori muscolari, male ai muscoli, mialgia
- appetite loss: appetito perso, inappetenza, perdita appetito, perdita di appetito
- shortness of breath: difficoltà respiratoria, difficoltà respiratorie, fiato corto, mancanza di respiro, respiro corto
- ▶ ...
- Ioss of the sense of smell: perdita olfatto
- COVID-19 terms: coronavirus, covid, covid-19, covid19


Our analysis considered the following countries and corresponding languages:

- United States of America (US), United Kingdom (UK), Australia, Canada English
- France French
- ► Italy Italian
- South Africa Zulu, Afrikaans, English, and many more
- ► Greece Greek

Modelling infectious diseases using online search

### Symptom-related search terms — Locations (countries) & languages



- 1. Query frequencies are **noisy**
- linear detrending

Google Trends Explore	
• headache Search term	
United Kingdom 🔻 9	/1/11 - 8/31/19 ▼ All categories
Interest over time (?)	
100	
75	
50	
25	
Sep 1, 2011	Feb 1, 2014

### harmonic smoothing using the frequencies of the past 2 weeks 2. Query frequencies are not stationary (increasing or decreasing mean)





- 3. For each symptom category, obtain the frequency sum across all its search terms (cumulative symptom-related search frequency) on a daily basis
- 4. Apply min-max normalisation on the cumulative frequency of each symptom category; values become from 0 to 1 and all categories now share units
- 5. Compute a daily weighted score using the FF100 symptom probabilities as weights
- 6. Use the previous 8 years (2011-2019) to obtain a historical baseline of this scoring function



For a given *day* and *location* 

- proportion of COVID-19-related news articles:  $m \in [0,1]$
- COVID-19 score based on web searches:  $g \in [0,1]$

**Decompose** g such that  $g = g_p + g_c$ 

- $-g_p$  represents 'infection'
- $-g_c$  represents 'concern'
- Then  $\gamma \in [0,1]$  exists such that

$$-g_p = \gamma g$$

$$-g_c = (1 - \gamma)g$$





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

Linear autoregressive model to forecast COVID-19 score g at a time point t based on its past values and the current and past values of *m* 

$$\arg\min_{\mathbf{w},\mathbf{v},b_2} \frac{1}{N} \sum_{t=1}^{N} \left( g_t - w_1 g_{t-1} - w_2 g_{t-2} - v_1 m_t - v_2 m_{t-1} - v_3 m_{t-2} - b_2 \right)^2 \to \text{prediction error } \epsilon_2$$

- expected to not have a causal effect on the estimated COVID-19 scores
- $\epsilon_1 \geq \epsilon_2 : \gamma = \epsilon_2/\epsilon_1$  (crude estimation of % of impact of news media)

Modelling infectious diseases using online search



Linear autoregressive model to forecast COVID-19 score g at a time point t based on its past values

 $\rightarrow$  prediction error  $\epsilon_1$ 

•  $\epsilon_1 < \epsilon_2$ : the media signal does not help COVID-19 score predictions  $\rightarrow \gamma \approx 1$ , i.e. the media is

- Data obtained from the Media Cloud database mediacloud.org
- Number of news media sources per country

US	225
UK	93
Australia	61
Canada	79
France	360
Italy	178
Greece	75
South Africa	135

title or main text e.g. "covid" or "coronavirus"

### News media coverage corpus

### Obtain the daily ratio of articles that include basic COVID-19-related keywords in their



- > 0 frequency from ~January, 2020 onwards
- ~2.5 million COVID-19-related articles from a total of ~10 million





Modelling infectious diseases using online search

#### Data obtained from September 30, 2019 to May 24, 2020

#### Average proportion of COVID-19-related news articles in the 8 countries of our analysis





Normalised online search score for COVID-19





Normalised online search score for COVID-19









Normalised online search score for COVID-19

Reducing news media effects:

- Altered trend during peak periods
- Average reduction by 16.4% (14.2%–18.7%) in a period of 14 days prior and after their peak moments, r = .822 (.739–.905)
- ► Reduction of 3.3% (2.7%-4%) outside peak periods





## Comparison with confirmed COVID-19 cases



Modelling infectious diseases using online search

Web search activity based models provide an early warning

 $r_{\rm max} = .83 (.74 - .92)$ when cases are brought forward by 16.7 (10.2–23.2) days

(South Africa is excluded)





### Comparison with *deaths of people with COVID-19*



Modelling infectious diseases using online search

Web search activity based models provide an early warning

 $r_{\rm max} = .85 (.70 - .99)$ 

when deaths of people with COVID-19 are brought forward by 22.1 (17.4–26.9) days

(South Africa is excluded)





- stages of the epidemic
- "Supervised" learning approach
  - corroborate our previous unsupervised findings
  - reporting system
- Source country: Italy
  - first major outbreak in Europe and among the countries in our study

• Transfer an incidence model — trained on web search activity — for a source country that has already experienced a COVID-19 epidemic to other *target* countries that are on earlier

will also transfer characteristics/biases of the source country, and especially of its clinical

Modelling infectious diseases using online search



# Transfer learning for COVID-19 incidence models

- **Source model**: regularised regression (*elastic net*)
  - use daily search query frequencies to estimate confirmed cases
  - Italy is our source country

$$\arg\min_{\mathbf{w},\beta} \left( \|\mathbf{y} - \mathbf{Sw} - \beta\|_2^2 + \lambda_1 \|\mathbf{w}\|_1 + \lambda_2 \|\mathbf{w}\|_2^2 \right)$$

$$\mathbf{S} \in \mathbb{R}^{M \times N}: M \text{ daily freq}$$
$$\mathbf{w} \in \mathbb{R}^{N}, \beta \in \mathbb{R}: \text{ regress}$$
$$\lambda_{1}, \lambda_{2} \in \mathbb{R}_{\geq 0}: \text{ regularisa}$$

Many regression models (~80K) — different regularisation amount

- sparsity levels from 5.5% to 91%
- 3 to 49 selected queries from the 54 we considered for Italy use this as crude quantification of model's uncertainty

- $\mathbf{u}$  uencies of N search terms sion weights and intercept ation parameters

Modelling infectious diseases using online search

- Establish search query pairs between the source and the target countries
  - Iookup for query pairs within the same symptom category
  - pair a source query to the target query with the greatest bivariate correlation, after identifying an optimal shifting period
- Transfer the regression weights from the source to the target feature space for all ~80K elastic net models
  - Final estimate of COVID-19 incidence is the mean over all elastic net models • .025 and .975 quantiles are used to form 95% confidence intervals
- Perform this daily from Feb. 17 to May 24, 2020, training models on increasing data from the source country

Modelling infectious diseases using online search



# Transfer learning for COVID-19 incidence models





# Transfer learning for COVID-19 incidence models — In practice





## Transfer learning vs. unsupervised learning





# Transfer learning vs. unsupervised learning



Modelling infectious diseases using online search

**Correlation** between the transferred models and the unsupervised models with reduced media effects

• 
$$r_{\rm avg}$$
 = .66

•  $r_{\text{max-avg}}$  = .80, when the *transferred* time series are brought 5 days forward





- Examine the statistical relationship between web search frequencies and confirmed COVID-19 cases (or deaths)
- Jointly for 4 English-speaking countries (US, UK, Australia, Canada)
  - attempt to reduce the bias of clinical endpoints in these different countries
  - focus on English-speaking countries for more comprehensive outcomes (without the need to translate searches)
- Use a broader set of search terms, not just symptom-related — figshare.com/projects/Tracking\_COVID-19\_using\_online\_search/81548
- Compute the joint bivariate correlation between search frequency and clinical indicators (cases or deaths) without any shifting and after shifting data so as to maximise it



### Correlation between web searches and COVID-19 cases

covid SARS-CoV-2 SARS CoV 2 COVID-19 coronavirus rash stay home quarantine covid NHS coronavirus pink eye how long does covid last covid symptoms COVID19 blue face sneeze coronavirus immunity lemsip migraine nCoV symptoms vomiting vomit

-0.4

-0.8





## Maximised correlation between web searches and COVID-19 cases



coronavirus dizziness SARS-CoV-2 SARS CoV 2 quarantine COVID-19 coronavirus rash covid NHS COVID-19 WHO coronavirus stomach pain covid symptoms coronavirus test covid how long does COVID-19 last coronavirus drugs COVID19 muscular pain feeling tired seizure vomit migraine

-0.7

Maximised correlation with confirmed COVID-19 cases





- Same 4 English speaking countries (US, UK, Australia, Canada)
- Joint approach again
- Multivariate regression analysis
  - Learn many elastic net models for different levels of sparsity (50%-99% to reduce the chance of overfitting) to jointly estimate cases or deaths based on web search data in these 4 countries
  - Train on data up to day d, test performance on the next day, d+1
  - Repeat this daily from the 2nd of March to the 24th of May, 2020
  - Use ground truth to find the best solution at each sparsity level
  - Compute the impact (average across all days) of each search term in the best solution at each density level



	(

covid blue face quarantine anosmia coronavirus mask appetite loss SARS CoV 2 coronavirus pink eye cant breathe loss appetite nasal congestion difficulty breathing rash coronavirus holidays chest tightness respiratory symptoms nose bleeding chills coronavirus cdc vomit

-6

-20

Estimation impact % (confirmed COVID-19 cases)

Modelling infectious diseases using online search

### Regression analysis – confirmed COVID-19 cases





covid SARS CoV 2 quarantine appetite loss blue face SARS-CoV-2 coronavirus pink eye rash loss appetite loss taste nose bleed head ache nasal congestion coronavirus chest xray how long does covid last tylenol feeling tired coronavirus high blood pressure tiredness diarrhea



-22

Estimation impact % (confirmed COVID-19 deaths)

11

27

43

Modelling infectious diseases using online search

-6





# COVID-19-related symptoms $\rightarrow$ better capturing community-level spread



The Royal College of General Practitioners (**RCGP**) swabbing scheme included people with no





Modelling infectious diseases using online search

## Translation and impact — Part of UK's COVID-19 surveillance



gov.uk/government/statistics/ national-flu-and-covid-19surveillance-reports-2023to-2024-season



Modelling infectious diseases using online search

#### Estimated COVID-19 prevalence score using Google search data for the UK

lockdown measures lockdown measures 2020 Google search score with reduced media effects Google search score Historical Trend (2011-19)





# Advanced models for disease prevalence forecasting

Morris, Hayes, Cox, Lampos (2023), PLOS Comput. Biol. Shu, Lampos (2024), Under review

Modelling infectious diseases using online search

# Part D



# Neural networks for disease forecasting – Feedforward baseline

#### Inputs

$$[\tau (m + 1) + 1] \times L_1 \qquad (L$$

$$F_{t_0 - \tau : t_0}, \mathbf{Q}_{t_0 - \tau + \delta : t_0 + \delta} \longrightarrow \mathbf{FC}_1 \longrightarrow \mathbf{ReLu}$$

- window of  $\tau + 1$  days
- Output: mean forecasted ILI rate and its standard deviation  $\gamma \delta$  days ahead
- mean and a standard deviation for each forecast
- Model / epistemic uncertainty: by training the BNN layer using variational inference

Morris et al. (2023), PLOS Comput. Biol.

Modelling infectious diseases using online search



• Input: web search activity (Q), previous ILI rates (F) with a temporal delay  $\delta$  flattened over a

BNN denotes a fully connected Bayesian layer with a probability distribution over its weights Data / aleatoric uncertainty: by using negative log likelihood as our loss function to obtain a

Combine data and model uncertainties by sampling the posterior of the NN's parameters

Multiple output estimates (samples) are used to derive a forecast and its confidence intervals



# Neural networks for disease forecasting — Simple RNN (SRNN)



Morris et al. (2023), PLOS Comput. Biol.

- Replace FF layers with a GRU layer
- Input is not flattened as it becomes a time series sequence



# Neural networks for disease forecasting – Iterative RNN (IRNN)



- Feeds this data back to itself, unlimited forecasting horizon
- not the predicted ones)
- Limitation: IRNN does not use forecasting distance to calibrate uncertainty

Morris et al. (2023), PLOS Comput. Biol.

Modelling infectious diseases using online search

Fully autoregressive, i.e. the network predicts all the input data for the next time step

Initially for a certain some of the data (Google) is known to us (we feed the actual data)







### **CRPS**: Continuous Ranked Probability Score **MAE**: Mean Absolute Error

#### r: bivariate (linear) correlation

 $\gamma$ :  $\gamma$  days-ahead compared to the last ILI rate in the input (autoregressive),  $\gamma$ -14 days ahead compared the last search query frequency

Morris et al. (2023), PLOS Comput. Biol.

Forecasting horizon	Accuracy metrics	FF	SRNN	IRNN
γ = 21 7 days ahead	CRPS	0.39	0.41	0.30
	MAE	0.51	0.55	0.42
	r	0.85	0.83	0.87
γ = 28 14 days ahead	CRPS	0.50	0.50	0.38
	MAE	0.63	0.64	0.53
	r	0.76	0.78	0.84







#### **Feedforward NN**

#### **Simple RNN**

Morris et al. (2023), PLOS Comput. Biol.







#### **Simple RNN**



Morris et al. (2023), PLOS Comput. Biol.

Modelling infectious diseases using online search

### Uncertainty calibration



Proportion of confidence intervalAverage over all test seasons

Morris et al. (2023), PLOS Comput. Biol.


- State of the art performance based on a CDC competition
- "Dante" leverages information from US regions, IRNN model does not
- IRNN provides better accuracy and more meaningful uncertainty bounds (see next slide!)
- ► NB: In these experiments, we have removed the ability of our model to use more recent web search activity

Morris et al. (2023), PLOS Comput. Biol.

Modelling infectious diseases using online search

Forecasting horizon	Accuracy metrics	Dante	IRNN
γ = 21	MAE	0.53	0.47
21 days ahead	r	0.73	0.81
γ = 28	MAE	0.61	0.60
28 days ahead	r	0.68	0.78

Osthus & Moran (2021), Nat. Commun.









### **Iterative RNN**

Dante

Morris et al. (2023), PLOS Comput. Biol.

### Modelling infectious diseases using online search





Shu, Lampos (2024), Under review

Modelling infectious diseases using online search



# DEFORMTIME — Influenza-like illness (ILI) forecasting accuracy (England)



Shu, Lampos (2024), Under review

Modelling infectious diseases using online search

### **Forecasting horizon (days ahead)**





## DEFORMTIME — Influenza-like illness (ILI) forecasting accuracy (US region 9)



Shu, Lampos (2024), Under review

Modelling infectious diseases using online search

Shu, Lampos (2024)

DeformTime

### **Forecasting horizon (days ahead)**



77



- "SOTA" forecasting models make similar forecasts to a persistence model - **Exceptions:** DeformTime, ModernTCN, Crossformer, LightTS - Forecasting influenza-like illness 28 days ahead seems possible

Shu, Lampos (2024), Under review

Modelling infectious diseases using online search

## DEFORMTIME vs. other models – 28 days ahead, 2018/19, England



- Web search activity can be used for infectious disease monitoring
  - the original Google Flu Trends model was definitely not the right approach
  - ... but there are ways to get this right!
- Disease models can be transferred to locations where historical disease rates are not available
- Unsupervised models based on web search activity
  - require a careful design
- **Deep learning** for time series forecasting
  - uncertainty / further improvements / interpretability
  - combine with mechanistic models

could be helpful for novel infectious diseases (COVID-19), esp. when everything else fails

Modelling infectious diseases using online search







- Lampos, Miller, Crossan, Stefansen. Advances in nowcasting influenza-like illness rates using search query logs. Scientific Reports 5 (12760), 2015. 1. doi:10.1038/srep12760
- Zou, Lampos, Cox. Transfer learning for unsupervised influenza-like illness models from online search data. WWW '19, pp. 2505-2516, 2019. 2. doi:10.1145/3308558.3313477
- Lampos et al. Tracking COVID-19 using online search. npj Digital Medicine 4 (17), 2021. doi:10.1038/s41746-021-00384-w 3.
- Eysenbach. Infodemiology: tracking flu-related searches on the web for syndromic surveillance. AMIA, pp. 244-248, 2006. 4.
- Polgreen, Chen, Pennock, Nelson. Using internet searches for influenza surveillance. Clinical Infectious Diseases 47 (11), pp. 1443-1448, 2008. 5. doi:10.1086/593098
- Ginsberg, Mohebbi, Patel et al. Detecting influenza epidemics using search engine query data. Nature 457, pp. 1012–1014, 2009. 6. doi:10.1038/nature07634
- 7. doi:10.1038/s41598-018-32029-6
- Budd, Miller, Manning et al. Digital technologies in the public-health response to COVID-19. Nature Medicine 26, pp. 1183-1192, 2020. 8. doi:10.1038/s41591-020-1011-4
- Rasmussen, Williams. *Gaussian Processes for Machine Learning*. MIT Press, 2006. 9.
- 10. Lampos, Zou, Cox. Enhancing feature selection using word embeddings: The case of flu surveillance. WWW '17, pp. 695-704, 2017. doi:10.1145/3038912.3052622
- 11. Levy, Goldberg. Linguistic regularities in sparse and explicit word representations. CoNLL '14, pp. 171-180, 2014. doi:10.3115/v1/W14-1618
- *analysis*. Bulletin WHO **99**, pp. 178-189, 2021. doi:10.2471/BLT.20.265603
- Biology **19** (8), 2023. doi:10.1371/journal.pcbi.1011392
- 14. Osthus, Moran. Multiscale influenza forecasting. Nature Communications **12** (2991), 2021. doi:10.1038/s41467-021-23234-5
- 15. Shu, Lampos. DEFORMTIME: Capturing Variable Dependencies with Deformable Attention for Time Series Forecasting. arXiv preprint, 2024. doi:10.48550/arXiv.2406.07438
- 16. Zeng et al. Are Transformers Effective for Time Series Forecasting? AAAI '23, pp. 11121-11128, 2023. doi:10.1609/aaai.v37i9.26317
- 17. Zhang, Yan. Crossformer: Transformer Utilizing Cross-Dimension Dependency for Multivariate Time Series Forecasting. ICLR '23, 2023. Link: openreview.net/forum?id=vSVLM2j9eie

Modelling infectious diseases using online search

## References

Wagner, Lampos, Cox, Pebody. The added value of online user-generated content in traditional methods for influenza surveillance. Scientific Reports 8 (13963), 2018.

12. Boddington 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

13. Morris, Hayes, Cox, Lampos. Neural network models for influenza forecasting with associated uncertainty using Web search activity trends. PLOS Computational

18. Luo, Wang. ModernTCN: A Modern Pure Convolution Structure for General Time Series Analysis. ICLR '24, 2024. Link: openreview.net/forum?id=vpJMJerXHU



