→ Recently Established Centre of Excellence for Analytics

→ 21 academic staff, RFs and PhD Students

→ In the UK, partnerships with Co-op, Boots, NHS, Olio, BBC, ONS

→ Core Expertise is in use of "Big Data" + Machine Learning" for Social Good

❖ Diagnosing disease with shopping data
➢ https://www.nlab.org.uk/project/shopping-data-disease/

❖ Donating personal transactional data for research
➢ https://www.turing.ac.uk/research/research-projects/donating-personal-transactional-data-research

# What is respiratory disease? Respiratory Disease ICD 10 coding: **J00–J99**

Infections and diseases of the lungs and respiratory system https://icd.who.int/browse10/2016/en#/J00-J06
These can be due to **COVID-19**.

| Disease | ICD-10 Codes | Respiratory Disease ICD-10: J00-J99 | Notes |
|---|---|---|---|
| Asthma | J45-46 | ✔ | |
| Chronic Obstructive Pulmonary Disease (COPD) | J40-47 | ✔ | Includes Bronchitis, Emphysema |
| Lung Cancer | C33-34 | ✗ | |
| Pneumonia and Influenza | J10-18 | ✔ | |
| Respiratory Tuberculosis | A16-19 | ✗ | J65 Pneumoconiosis is associated with tuberculosis |
| Cystic Fibrosis | E84.8-84.9 | ✗ | |

- COVID-19 led to unparalleled pressure on healthcare services, with improved healthcare planning respiratory diseases becoming a key concern.

- We present the results of **two related studies** investigating the potential of using digital footprint data, in the form of non-prescription medication sales, to improve forecasting of weekly registered deaths of such diseases at local levels.

# Study: A new experimental design to compare the use of sales data against other datasets used in the prediction of respiratory deaths

Applying Hofman et al.'s [1] recommendations on computational social science:

- Baseline model
- Out-of-sample testing
- Combine prediction and explanation  - Here we use **a new AI Explainability Tool Model Class Reliance (MCR) for Random Forest Regressor**
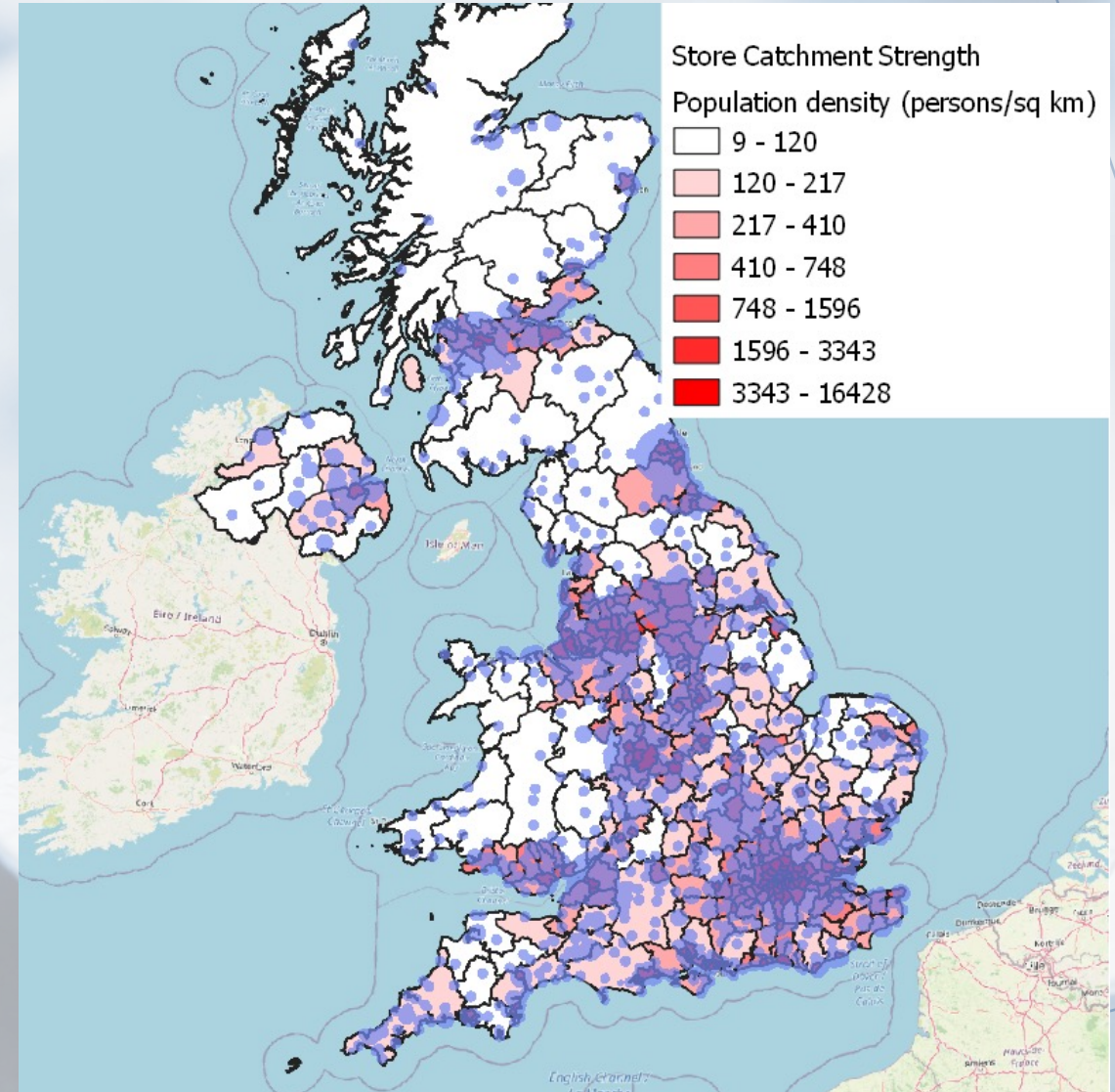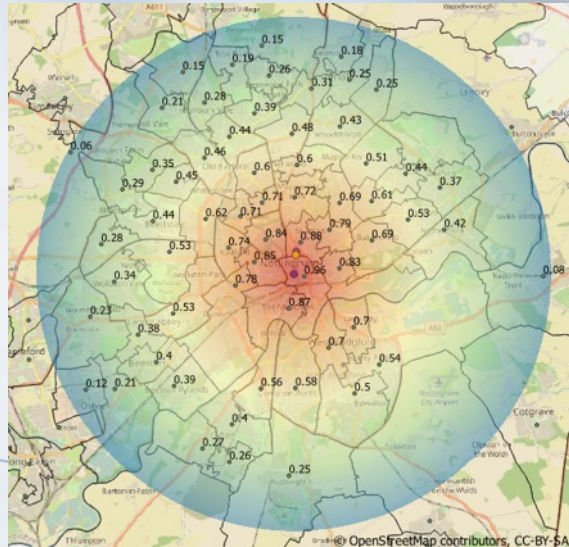


[1] Hofman, J. M. et al. Integrating explanation and prediction in computational social science. Nature 595, 181-188 (2021).

# Explaining the Commercial Sales features

## BASIC STATS:

Data using in models (2016 to 2020) is over 2 billion in store sale units in England (not including online)

For example over 14 million cough medicine sales

**Sales are allocated from stores to 314 LTLAs in England**





Store Catchment Strength

Population density (persons/sq km)

- 9 - 120
- 120 - 217
- 217 - 410
- 410 - 748
- 748 - 1596
- 1596 - 3343
- 3343 - 16428

# Creating the PADRUS* model

**\*Predicting amount of deaths from respiratory disease using sales data**

**Model type:** Random Forest Regressor. Optimized using a time series cross-validation grid-search on training data to prevent over-fitting.

**Target:** Predict weekly respiratory deaths 17 days in advance for each of the 314 LTLA (Lower Tier Local Authority) areas in England from 18th March 2016 to 27th March 2020.

**Input features:** 56 features (static and dynamic including traditional variables associated with respiratory disease)

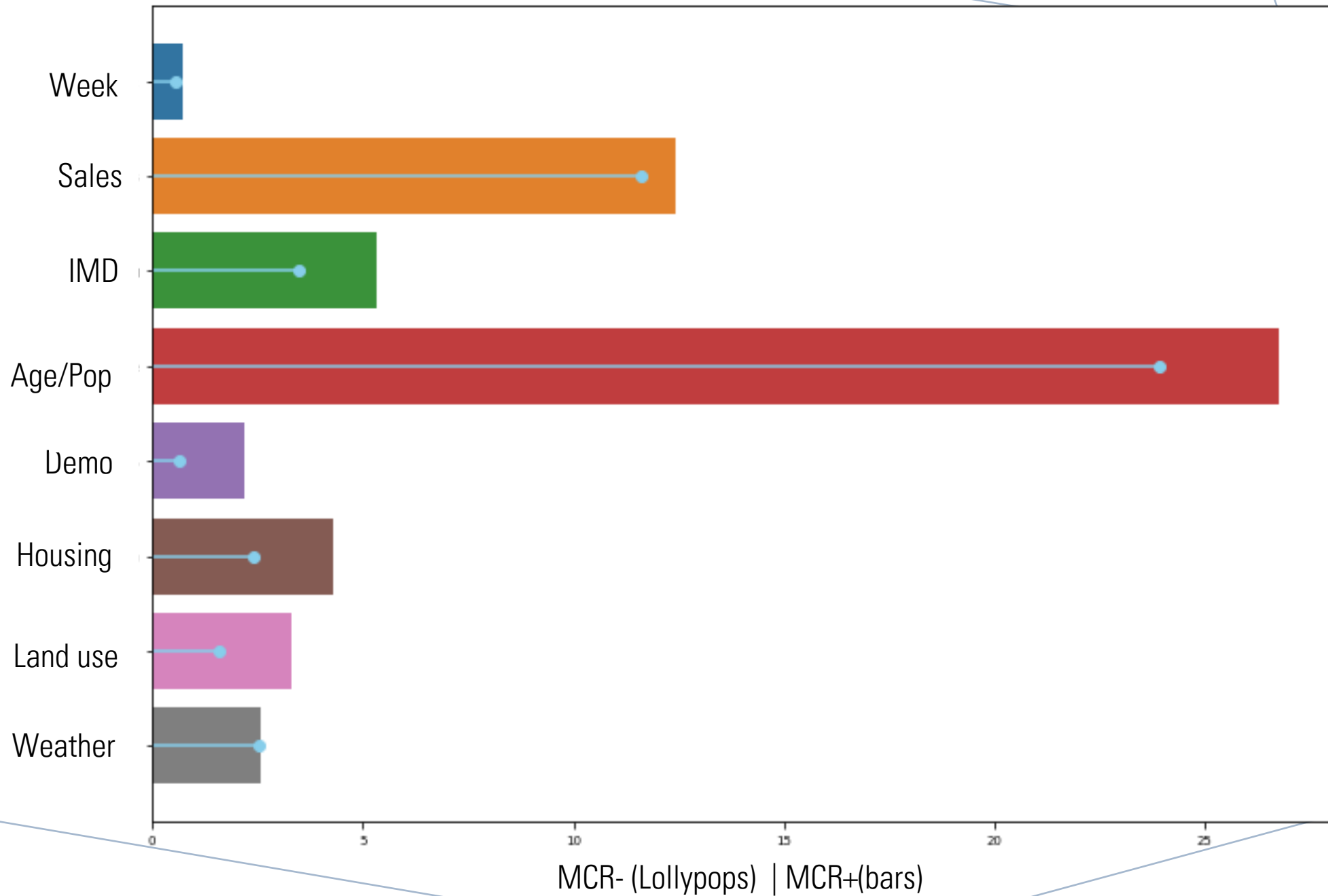**Training Datapoints:** 45844

*Out-of-sample Testing*
**Testing Datapoints:** 20410

| Results from predicting on test data: | | |
|---|---|---|
| Mean Absolute Error | Root Squared Mean Error | $R^2$ |
| 2.39 (Baseline model 2.78) | 3.42 (Baseline model 4) | 0.78 (Baseline model 0.71) |

With highest accuracy gains over models without digital footprint data (increases in $R^2$ between 0.09 to 0.11) occurring in periods of maximum risk to the general public.
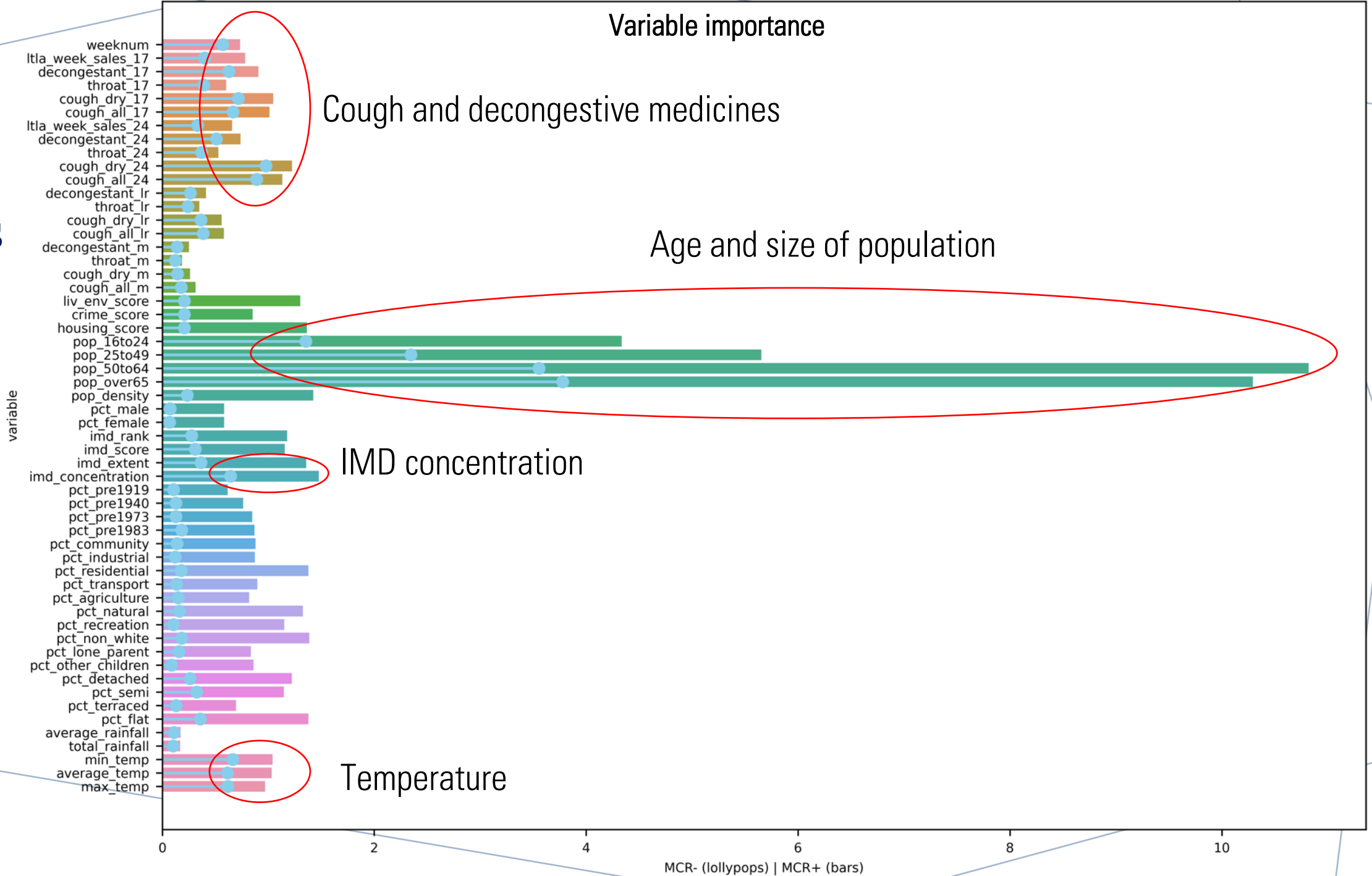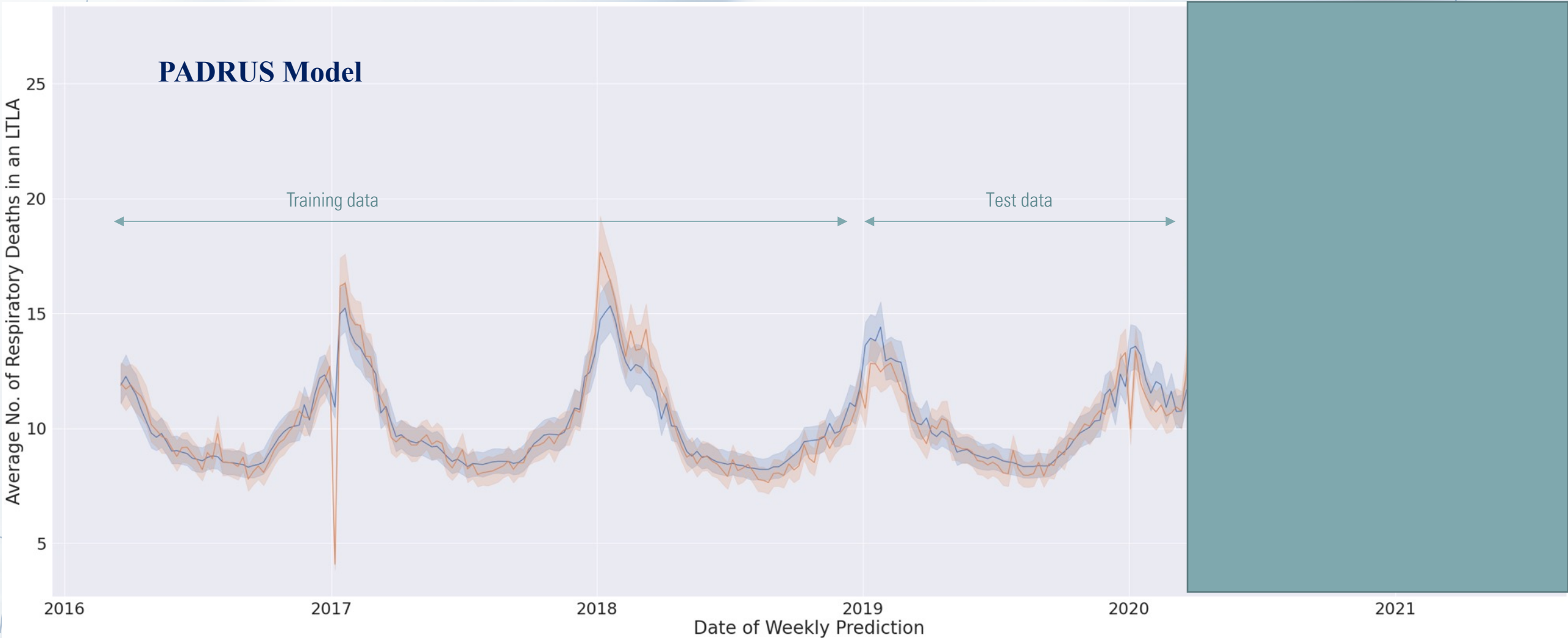
Grouped variable importance

Group MCR PADRUS model

MCR- (Lollypops) | MCR+(bars)

MCR PADRUS model (run on training data)…

Variable importance

Cough and decongestive medicines

Age and size of population

IMD concentration

Temperature

MCR- (lollypops) | MCR+ (bars)
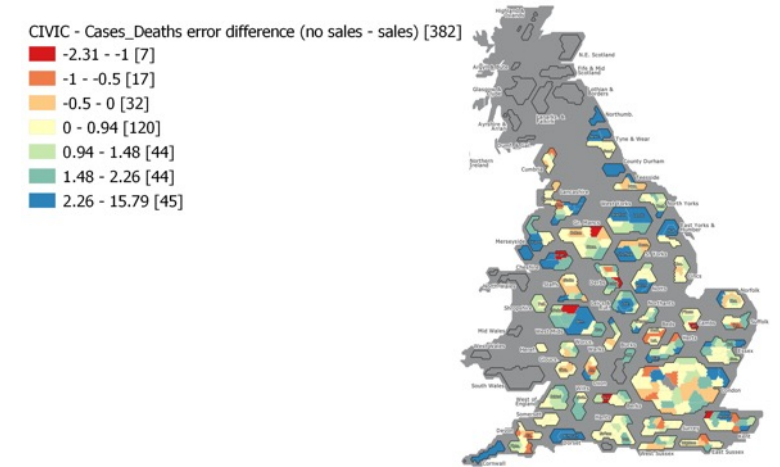
Running the model during COVID-19

**Study: Creating the CIVIC model:** Predicting Covid-19 Impact on Vulnerable Individuals and Communities via health, deprivation, and transactional sales data

- Addressed the feature drift of sales variables due to government lockdowns

- Used additional data resources available in the pandemic,
  i.e. easily accessible and regularly updated COVID-19 test, case and mortality data,
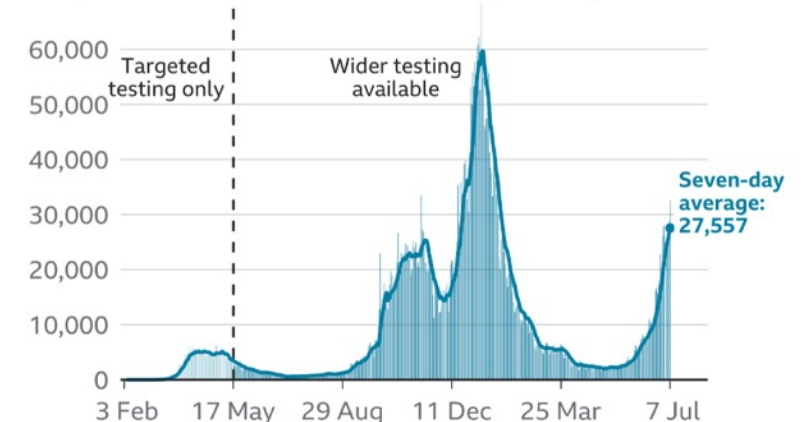  newly available mobility data

- MODEL TYPE : A range of ML regression models (Linear, SVR, Random Forest & XGBoost) were developed and tested across the initial waves of the COVID-19 pandemic

- TARGET : Forecast of COVID-19 cases/deaths in 7/14/21 days at the 314 Lower Tier Local Authorities in England

- INPUT FEATURES: Grouped feature sets(e.g. cases, sales of over-the-counter medications, demographics, mobility, etc.) combined across multiple runs to test predictive performance and identify optimal set of predictors. (Retail sales attributed to LTLAs using simulated catchment model)

- Datapoints : Weekly time series covering period from April 2020 to June 2021. Split into two waves (April-July 2020, July 2020-March 2021)



CIVIC - Cases_Deaths error difference (no sales - sales) [382]
- -2.31 - -1 [7]
- -1 - -0.5 [17]
- -0.5 - 0 [32]
- 0 - 0.94 [120]
- 0.94 - 1.48 [44]
- 1.48 - 2.26 [44]
- 2.26 - 15.79 [45]



Number of new cases rising
Daily confirmed coronavirus cases by date reported

Targeted testing only

Wider testing available

Seven-day average: 27,557
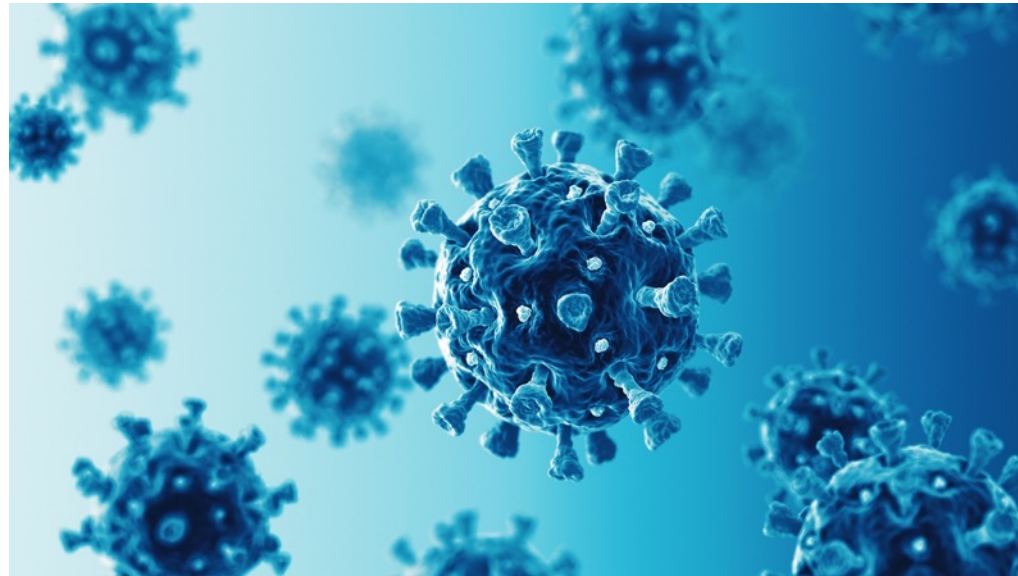
Source: Gov.uk dashboard, updated 7 Jul

BBC

**Results**

- XGBoost models containing sales data optimally predicted the number of COVID deaths 21 days in advance (R2=0.68**), significantly outperforming models based on official COVID case data alone at local-area levels (R2=0.44**)

- Demographics and mobility inputs less useful than retail sales.

(Without demographics – 0.61**)

**Conclusions:**

- Over-the-counter medication purchases related to management of respiratory illness are correlated with registered deaths at a 17-21 day window.

- Results demonstrate the potential for sales data to support early warning population health mechanisms at local area levels.

Any questions?

@3lizabeth_Dolan

Useful Links:

https://www.nlab.org.uk/project/shopping-data-disease/

https://www.researchsquare.com/article/rs-2226531/v1

https://github.com/nhsx/commercial-data-healthcare-predictions

**Forecasting local COVID-19/Respiratory Disease mortality via national longitudinal shopping data:** the case for integrating digital footprint data into early warning systems