



# Predicting Pedestrian Counts using Machine Learning

Molly Asher<sup>1</sup>, Yannick Oswald<sup>2</sup>, Nick Malleson<sup>2</sup>

<sup>1</sup> School of Earth and Environment, University of Leeds

<sup>2</sup> School of Geography, University of Leeds

These slides: <https://urban-analytics.github.io/dust/presentations.html>



**UNIVERSITY OF LEEDS**







# Context

Accurately predicting number of pedestrians is both *important* and *challenging*

## **Aims and objectives:**

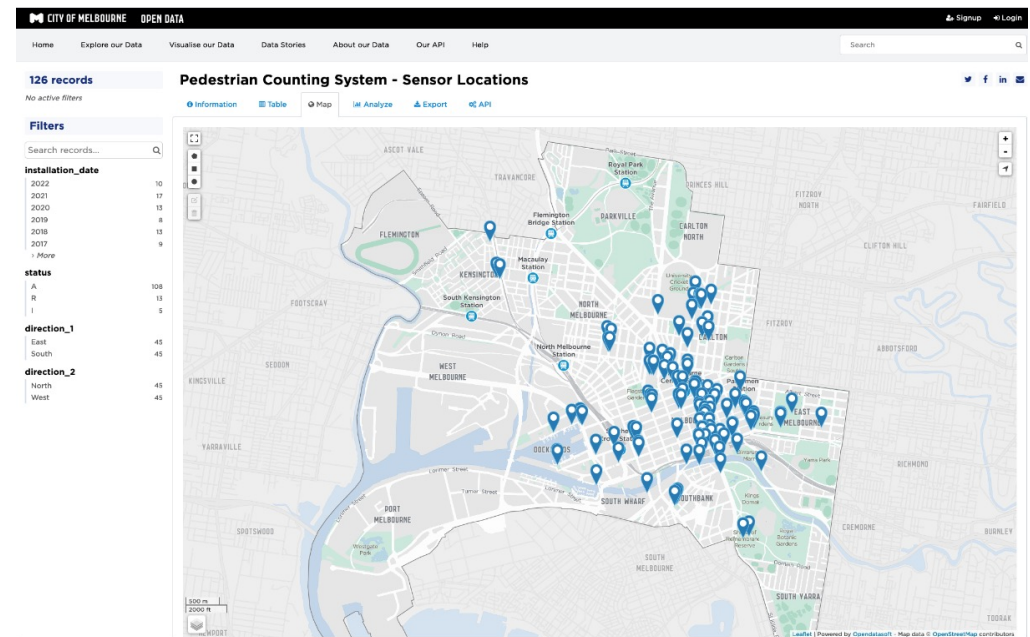
This work is using machine learning to:

- Better understand the impact of the built environment and other contextual factors on pedestrian counts
- Predict the number of pedestrians at un-sampled locations under different conditions
- Evaluate the success of past events



# Melbourne Open Data

- Melbourne Open Data Portal for open data:<sup>1</sup>
  - Land-use, litter, built environment, roads, bike sharing, air quality, etc.
- Network of pedestrian sensors:
  - 18 sensors in 2009
  - 82 sensors in 2022
- Record bi-directional pedestrian movements 24h/day every hour



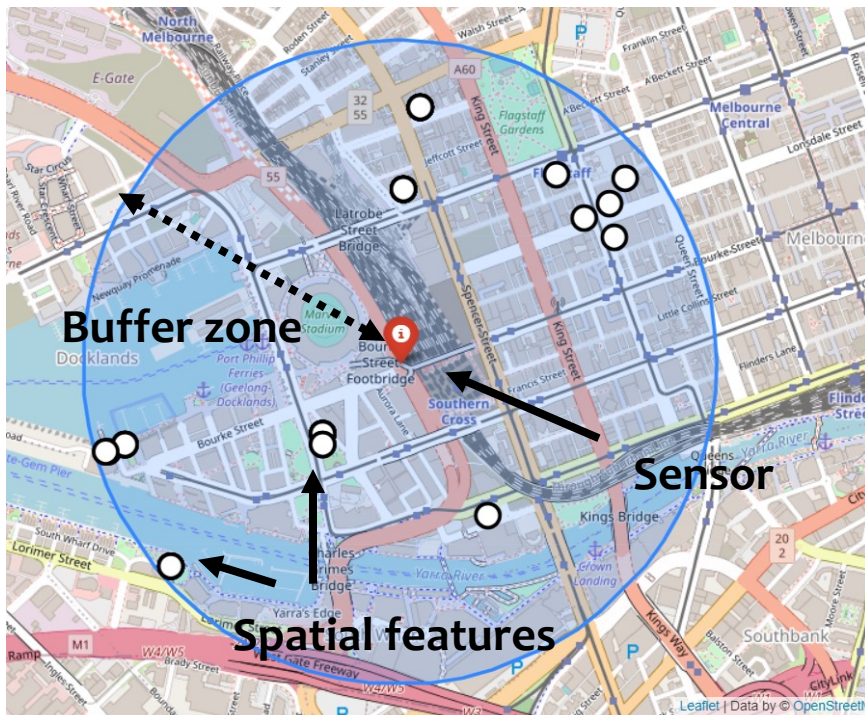
Locations of sensors in the City of Melbourne<sup>1</sup>



<sup>1</sup> <https://www.data.vic.gov.au/discover-city-melbourne-open-data>



# Melbourne Open Data



Example buffer zone within which spatial features are linked to sensors

- Numerous additional open data sets, including:
  - Weather
  - Street furniture (benches, bins etc)
  - Buildings
  - Landmarks
- Buffer zones drawn around sensors to link them to features describing urban environment in vicinity



**UNIVERSITY OF LEEDS**

# Modelling Overview

- Dependent variable:
  - number of pedestrians per sensor per hour
- Explanatory variables:
  - time of day (hour, day, month, year)
  - weather conditions (temperature, humidity, wind speed)
  - road betweenness (a measure of how well integrated the nearest road is to the rest of the network)
  - local built environment variables (number of trees, benches, buildings, public transport, etc., etc.)
- Trained on available sensor data (4 million rows)
- Later used to predict at locations without sensors





# Model selection

- Candidate models evaluated using 10-fold cross-validation
- Error metrics (RMSE and MAE) calculated on the predicted counts-per-hour of pedestrians from 10-fold cross-validation of each model against actual values from the sensor data

## Error metrics

Model	MAE	RMSE
Linear regression	268.40	370.54
Random Forest regression	89.88	179.62
XGBoost	121.35	207.40



# Model selection

- Candidate models evaluated using 10-fold cross-validation
- Error metrics (RMSE and MAE) calculated on the predicted counts-per-hour of pedestrians from 10-fold cross-validation of each model against actual values from the sensor data

Error metrics

Model	MAE	RMSE
Linear regression	268.40	370.54
<b>Random Forest regression</b>	<b>89.88</b>	<b>179.62</b>
XGBoost	121.35	207.40

Random forest regressor selected as best performing model

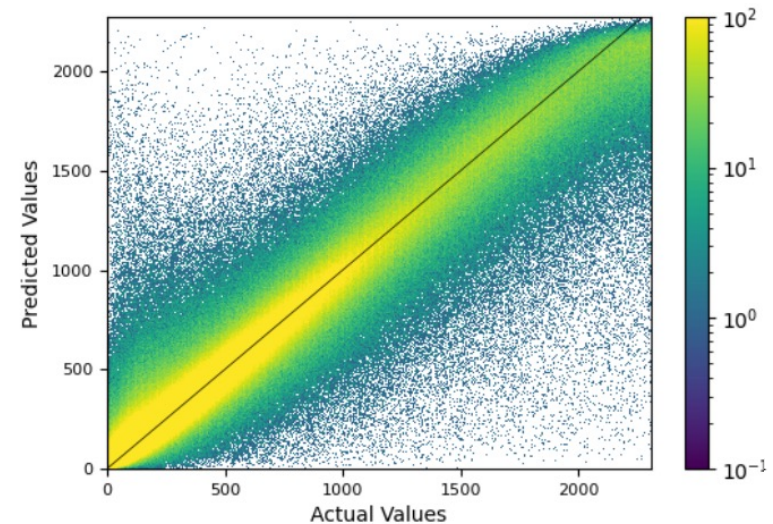




# Model evaluation

- Predicted counts-per-hour of pedestrians plotted against actual values from the sensor data
- Most predictions fall around the diagonal ( $x=y$ ), giving confidence that model is not biased towards smaller or larger counts

## Random forest regressor



MAE

89.88

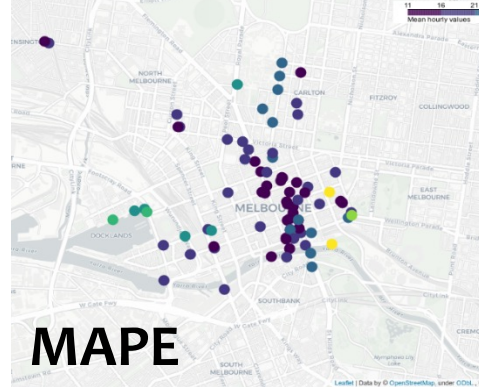
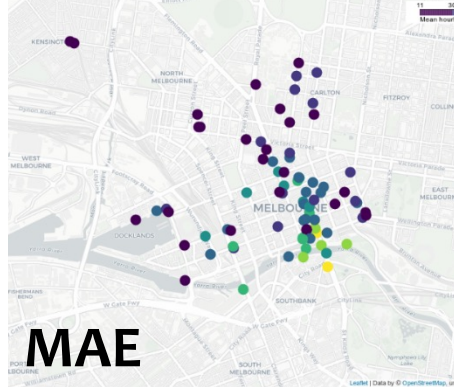
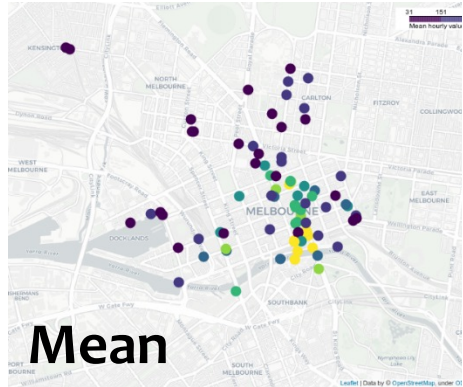
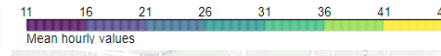
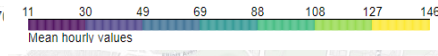
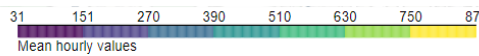
RMSE

179.62



UNIVERSITY OF LEEDS

# Spatial variation



Central and southern sensors capture highest footfall



Patterns of absolute error follow (roughly) the mean, with some deviations



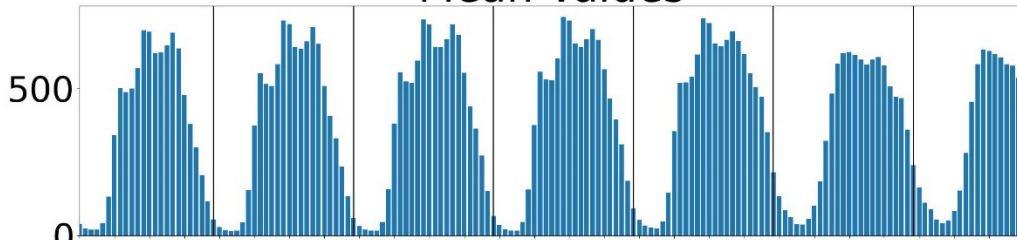
Several sensors with much larger percentage error





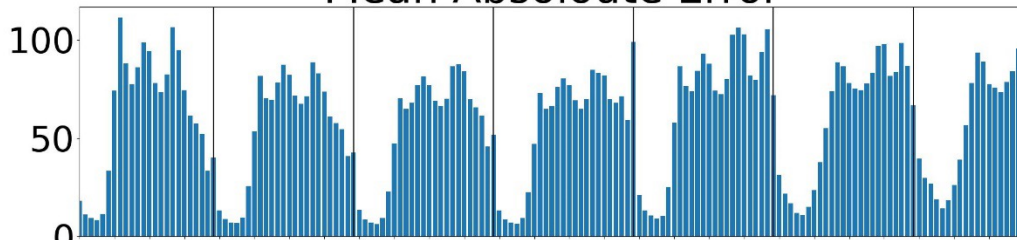
# Temporal variation

Mean Values



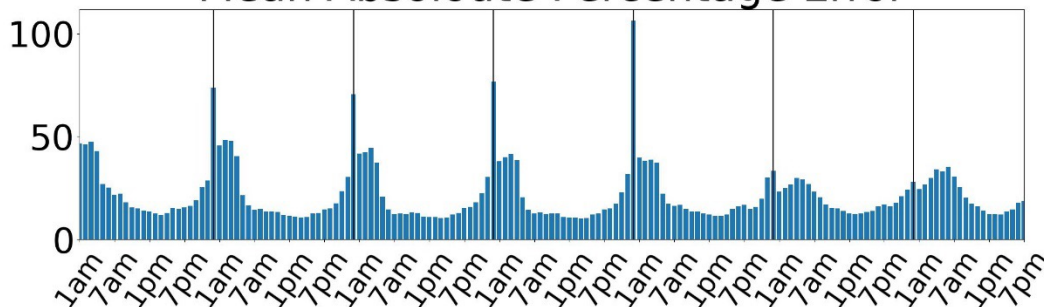
Reflect typical city centre patterns

Mean Absolute Error



Similar patterns to mean (larger counts = larger errors)

Mean Absolute Percentage Error

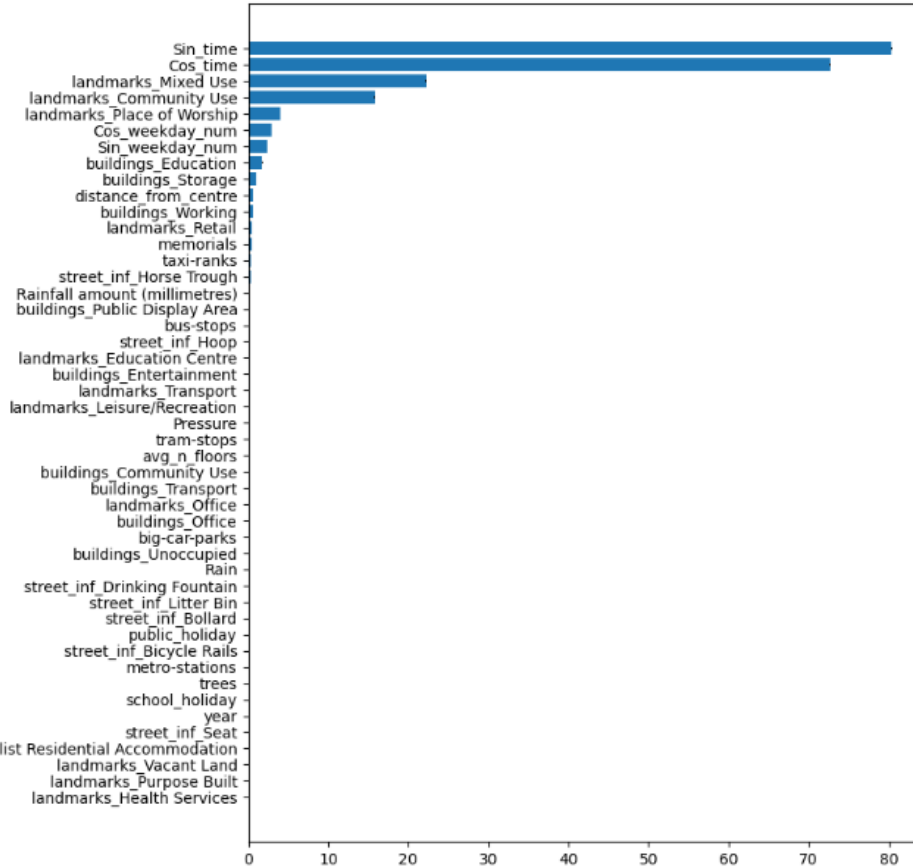


Largest errors at night



# Feature importance

Permutation importance

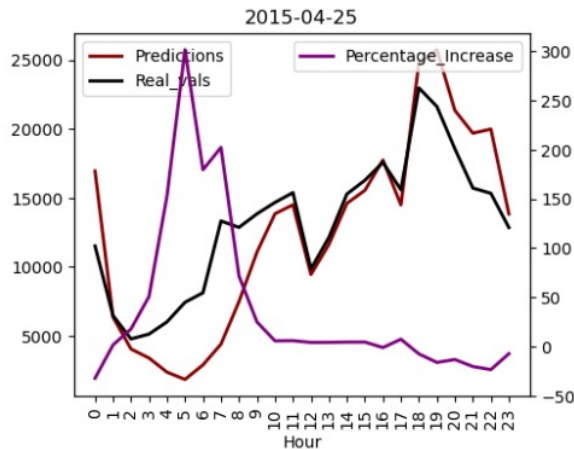


- Ranks features contribution to the model's predictions
- Most important features:
  - Hour of day
  - Landmarks (mixed, community use, places of worship)
  - Weekday
  - Educational buildings
- (surprisingly) lower importance features:
  - Betweenness
  - School/public holiday

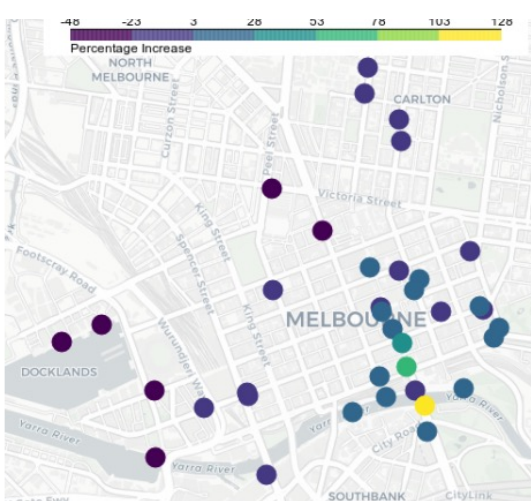




# Evaluating events



- Model can be used as a tool to evaluate success of events
- E.g. Anzac Day Parade:
  - 5% more footfall in whole city over 24h
  - 72% more footfall from 3-10am
  - 128% more footfall at a sensor in south-east near parade location



# Conclusions

- Ongoing work aiming to:
  - accurately predict the number of pedestrians in time and space at un-sampled locations under different conditions
  - better understand the impact of the built environment and other contextual factors on pedestrian counts
  - Evaluate the success of past events
- Model performs reasonably well overall
- Some spatial and temporal variations in prediction error
- Beginning to make inferences about impact of urban environment







# Predicting Pedestrian Counts using Machine Learning

Molly Asher<sup>1</sup>, Yannick Oswald<sup>2</sup>, Nick Malleson<sup>2</sup>

<sup>1</sup> School of Earth and Environment, University of Leeds

<sup>2</sup> School of Geography, University of Leeds

These slides: <https://urban-analytics.github.io/dust/presentations.html>



**UNIVERSITY OF LEEDS**