

The background of the slide is a detailed architectural rendering of a modern urban environment. It features a mix of tall, glass-fronted office buildings and lower-rise structures with more organic, textured facades. In the foreground, there is a lush green park area with a winding wooden boardwalk. Numerous people are depicted in various activities: walking, sitting on benches, playing in a small garden, and using a stroller. A large, leafy tree stands on the left, and a yellow patio umbrella is visible on the right. The sky is a clear, bright blue with some light clouds.

Semester project

Data-driven Analysis of City Scale Human-building Interaction

Yihan Wang

MSc Civil Engineering

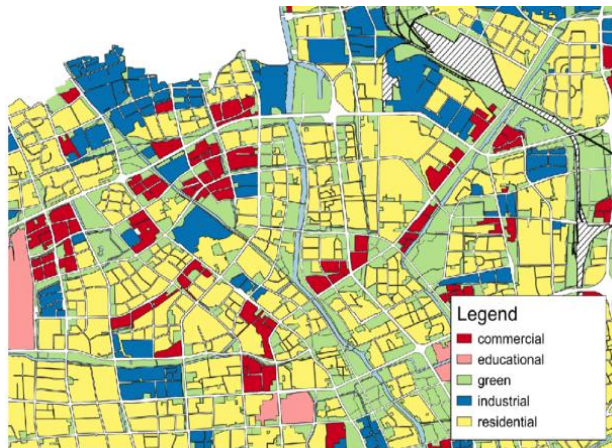
Supervisor:

Andrew Sonta, Kanaha Shoji

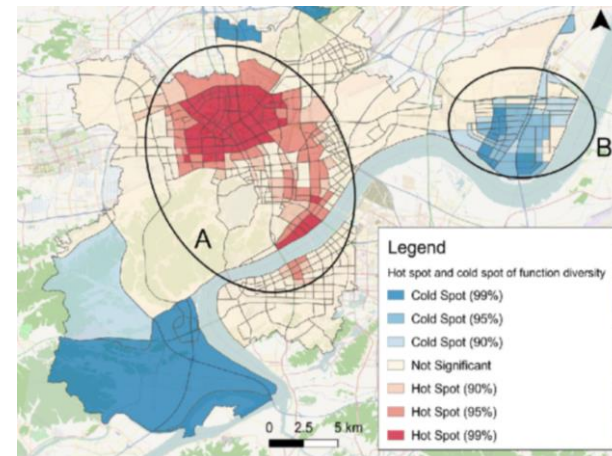
24th Jan 2024



West Lake, Hangzhou, China



Hangzhou urban planning map^[1]

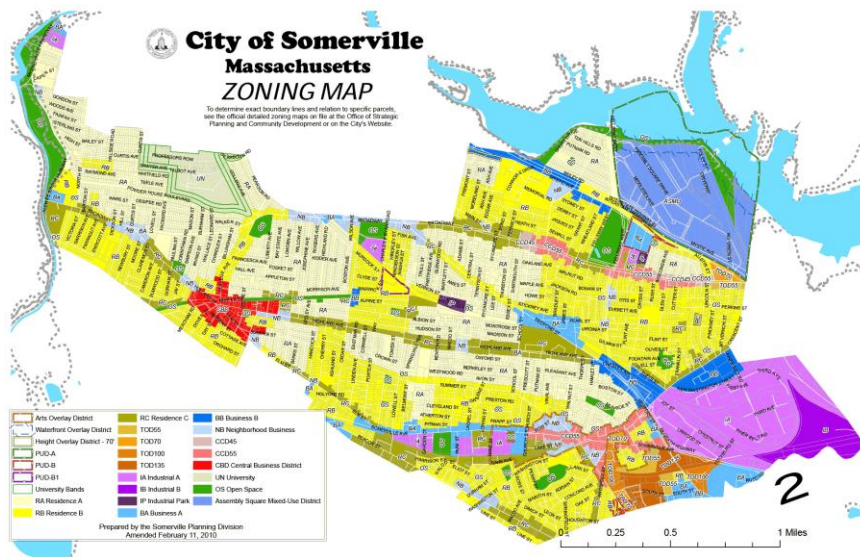


Hangzhou function diversity
hot and cold spot map^[1]

Extraordinary view existence \neq Good living experience !

More challenging urban design exists...

- **Euclidean zoning (areas with specific rules for building use)**



Question: What kind of urban design is considered benefiting internal individual residents and society?

Objective: Understanding human's interaction with urban forms

Jeopardizing vibrant and socially resilient communities^[1]

Intro

■ 2. Jacobs, J. (1961). *The Death and Life of Great American Cities*. New York.



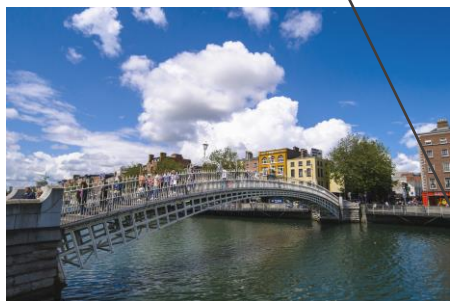
OUTLINE

Data Exploration

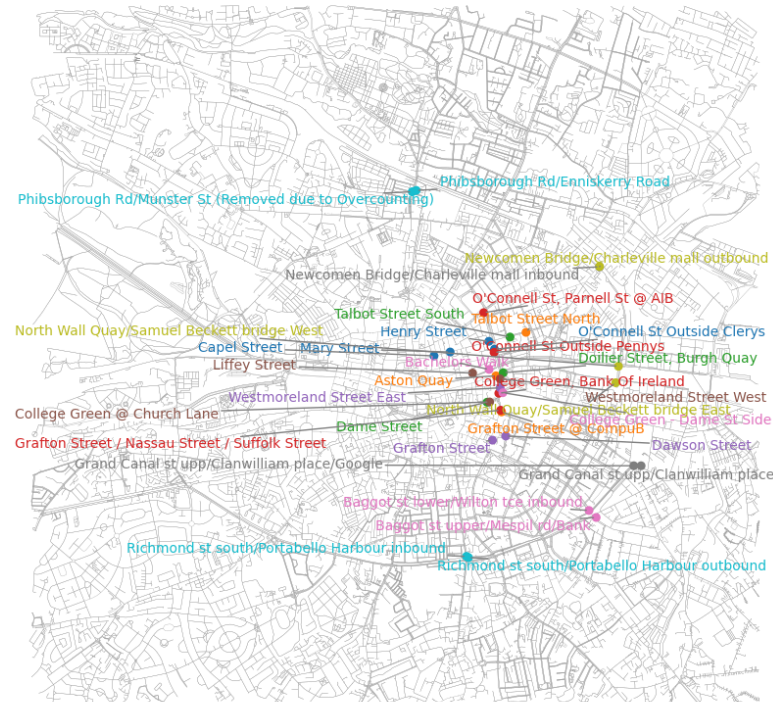
- Pedestrian Counts
- Urban Forms
- Other Features

Model Regression & Interpretation

 Dublin



Sensor Locations on Street Network in a region in Dublin

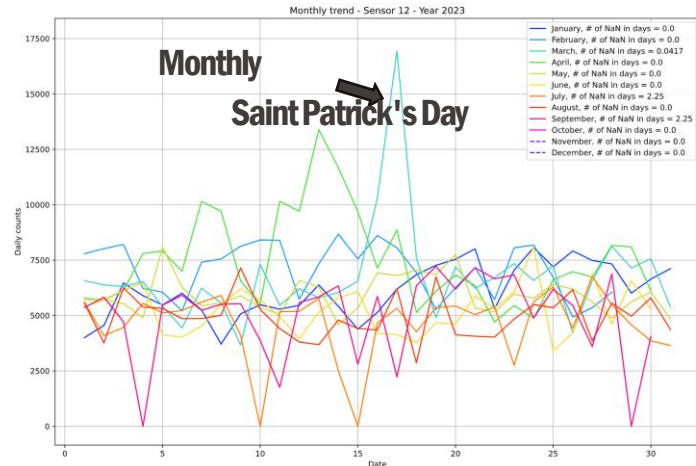
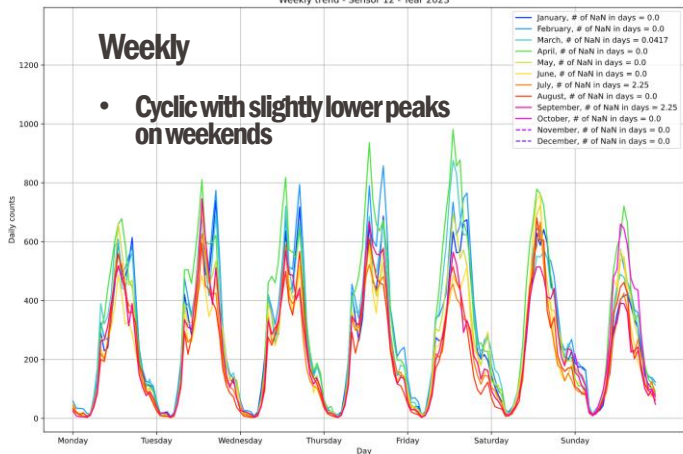
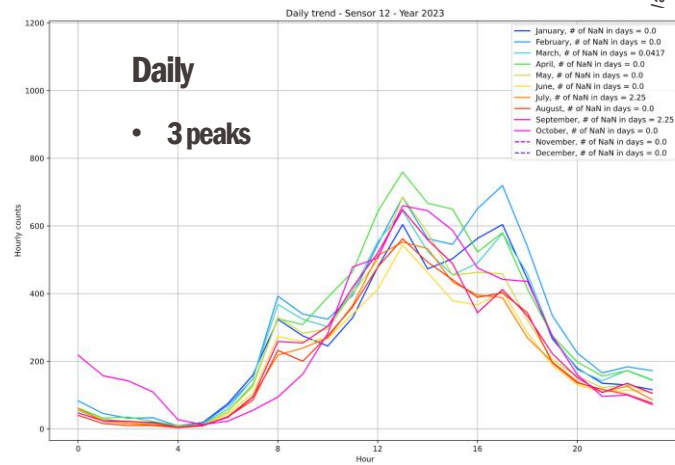
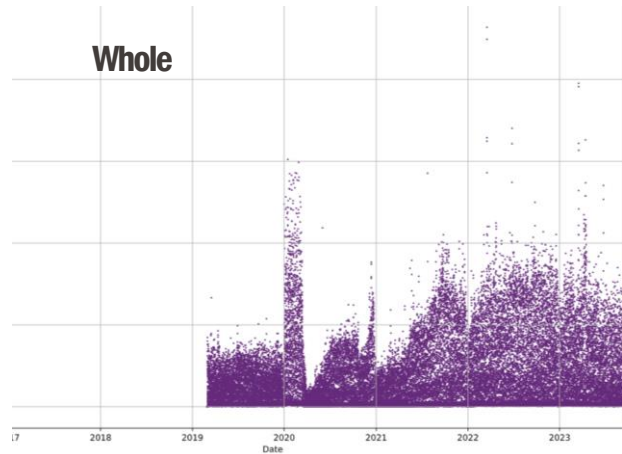


■ Data

- 34 pedestrian counters
- Densely distributed around central region

SensorID: 12

College Green, Bank of Ireland





OUTLINE

Data Exploration

- Pedestrian Counts
- Urban Forms
- Other Features

Model Regression & Interpretation

SensorID: 12

College Green, Bank of Ireland



Determine the closest node → Retrieve 15 mins walk polygon → Inside urban forms

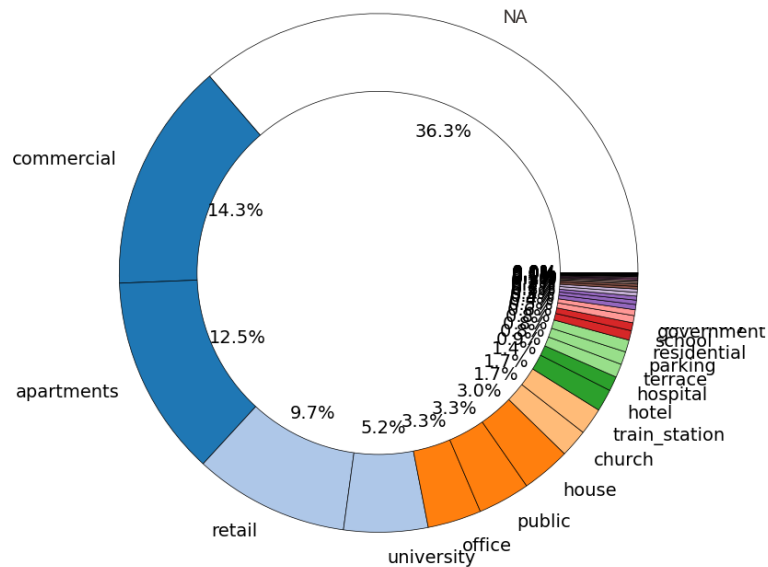
SensorID: 12

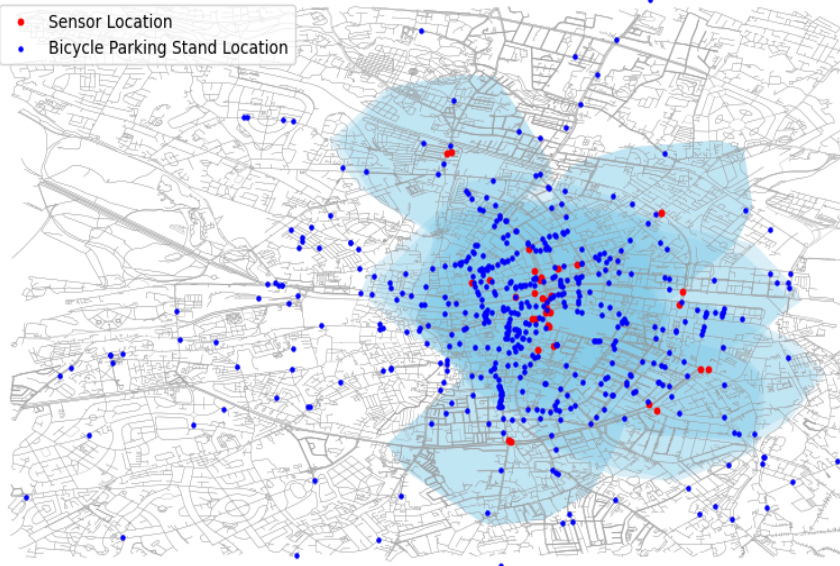
building

College Green, Bank of Ireland

NA Polygon	- 1128081 m2
commercial	- 445261 m2
apartments	- 389587 m2
retail	- 300199 m2
university	- 162703 m2
office	- 104055 m2
public	- 101022 m2
house	- 94768 m2
church	- 54094 m2
train_station	- 51916 m2
hotel	- 43019 m2
hospital	- 28158 m2
terrace	- 26130 m2
parking	- 25108 m2
residential	- 23461 m2
school	- 18477 m2
government	- 15382 m2
dormitory	- 13910 m2
theatre	- 11329 m2
industrial	- 10048 m2
roof	- 8492 m2
cathedral	- 7644 m2
transportation	- 7276 m2
college	- 6253 m2
construction	- 5127 m2
warehouse	- 4362 m2
civic	- 2851 m2
sports_hall	- 2820 m2
library	- 2317 m2
greenhouse	- 2217 m2
service	- 1950 m2
shed	- 1926 m2
chapel	- 1265 m2
presbytery	- 891 m2
ship	- 675 m2
bank	- 528 m2
garage	- 518 m2
guardhouse	- 507 m2
hostel	- 500 m2
monastery	- 424 m2
kindergarten	- 412 m2
restaurant	- 374 m2
hut	- 301 m2
detached	- 276 m2
bridge	- 254 m2
mixed_use	- 233 m2
mews	- 161 m2
arch	- 63 m2

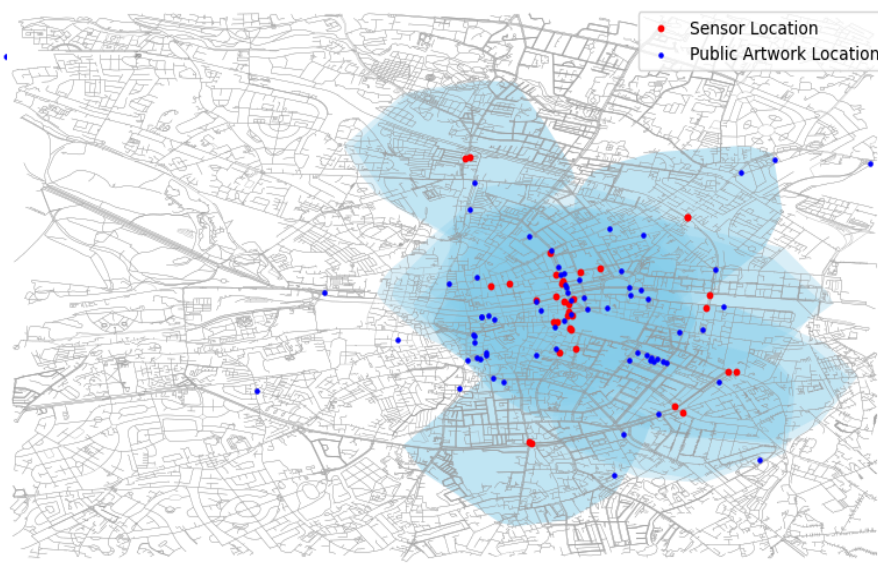
Total Area of Each Building Type (Ordered by Area) within the 15 min walk distance from College Green Counter



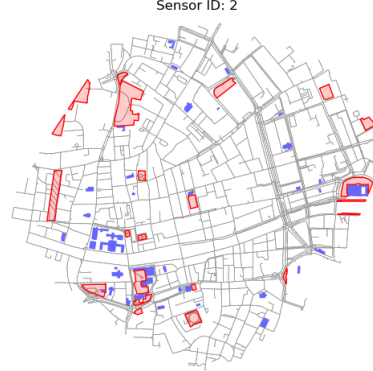
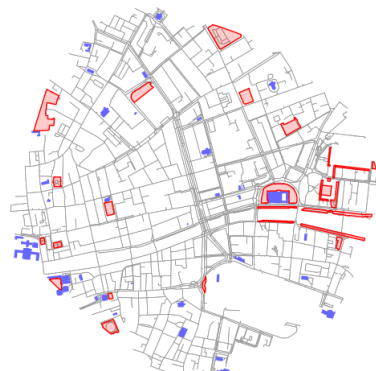
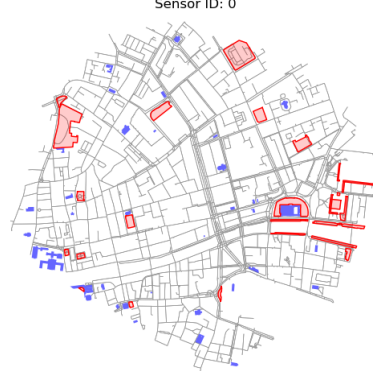


Sensor ID: 0

Sensor ID: 1



Sensor ID: 2



■ From OSMnx (public area:
['public', 'church', 'cathedral', 'chapel', 'civic'])

■ From open dataset (park)



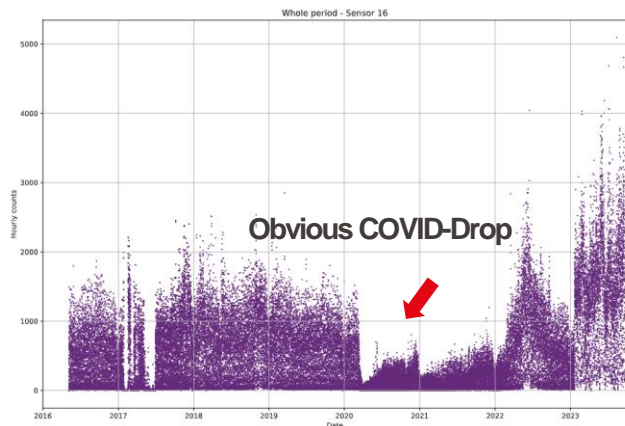
OUTLINE

Data Exploration

- Pedestrian Counts
- Urban Forms
- Other Features

Model Regression & Interpretation

- Filtering out data with COVID-19 restrictions: ($\rightarrow 4115348 \times 135$ dataframe)



Restrictions including:

- School Closing
- Workplace Closing
- Cancel Public Events
- Restrictions on Gatherings
- Close Public Transports
- Protection of elderly people
- Stay at Home Requirements
- Restrictions on Internal Movement
- International travel controls
- Restrictions on Gatherings
- Contact Tracing
- Facial Coverings

- One-hot encoding for categorical features

- Z - score Standardization: $x'_{i,j} = \frac{x_{i,j} - \min(x_j)}{\max(x_j) - \min(x_j)}$

- Normalization: $x'_{i,j} = \frac{x_{i,j} - \text{mean}(x_j)}{\text{std}(x_j)}$

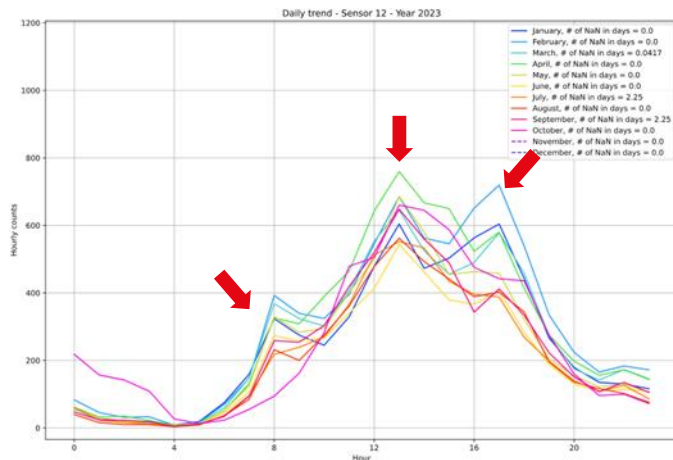
- Deal with NaN values: fill with feature means

Interpretable model: linear regression

Dependent Variable		Hourly pedestrian counts		Perimeter Mean	Recreational	PopSum	distance	B15_Comm ercial_per	B5_others_ per	S15_street_le ngth_avg
Method		Ordinary Least Squares		Perimeter Stdev	Green View Mean	Men	feature_x	B15_Reli gious_per	S15_n	S15_circuitry_ avg
Predictors		Degree	Footprint Proportion	Complexity Mean	Sky View Mean	Women	feature_y	O5_Walkab leLength_m	S15_m	S15_self_loop _proportion
cc	ws	Clustering	Footprint Mean	Complexity Stdev	Building View Mean	Elderly	nearest_x	O5_PublicTranspor tAccessCounts	S15_k_avg	S5_n
speed_u	Node Density	Closeness Centrality	Footprint Stdev	Building Count	Building View Stdev	Youth	nearest_y	O5_Amenity_Shan non_Entropy	S15_edge_ length_total	S5_m
speed_v	PageRank	Betweenness Centrality	Commercial	Food	Road View Mean	Children	O15_Walkabl eLength_m	B5_Accomm odation_per	S15_edge_ length_avg	S5_k_avg
rh	Street Length	Civic	Social	Healthcare	Road View Stdev	O15_Green Area_m2	O15_PublicTranspor tAccessCounts	B5_Comme rcial_per	S15_streets_per _node_avg	S5_edge_ length_total
temp_K	day_type	Entertainment	Perimeter Total	Institutional	Visual Complexity	n	O15_Amenity_Sha nnon_ Entropy	B5_Religious _per	S15_street_ length_total	S5_edge_length _avg ...

Features including: weather, urban form, network property, population...

- Plain model: $R^2 = 0.287$
- Baseline model: Categorizing hour in day: $R^2 = 0.523$



- Morning Peak: 7, 8, 9
- Lunch Peak: 12, 13, 14
- Evening Peak: 16, 17, 18

Hour categorization helps explaining the variation a lot!

- Adding a predictor of 'year': $R^2 = 0.538$

Other temporal predictor remains unexplored

Baseline model: $R^2 = 0.523$

• Exploring feature interactions: $R^2 = 0.545$

Intuition: For population groups,
the fractions are more meaningful.

Included interactions

Intuition: Young people may be
influenced by day type more

Men / PopSum	Youth * Commercial	<day_type * Institutional>	< day_type * B15_Commercial_per >
Youth / PopSum	Youth * Entertainment	<day_type * Recreational>	< day_type * B15_ Religious_per >
Children / PopSum	PopSum * Commercial	<day_type * O15_GreenArea_m2>	< day_type * B5_ Religious_per >
day_type * morning_peak	PopSum * Entertainment	<day_type * O15_WalkableLength_m>	<day_type * Youth>
day_type * lunch_peak	PopSum * Building Count	< day_type*O5_WalkableLength_m >	<day_type * Men>
day_type * evening_peak	<day_type * Commercial>	<day_type * B5_Accommodation_per>	<day_type * Women>
< day_type * O15_PublicTransportAccessCounts >	<day_type * Entertainment>	<day_type * B5_Commercial_per>	<day_type * Elderly>
<day_type * O5_PublicTransportAccessCount>	<day_type * Building Count>	< day_type * B15_Accommodation_per >	<day_type* Children>

Intuition: people may have different public
transport use behavior on holidays vs
workdays

<>: $p < 0.05$

• Adding additional features: $R^2 = 0.580$

(yearly average data: gdp, green gas emission from different sections)

Decomposition of the residuals

- Weather features:**

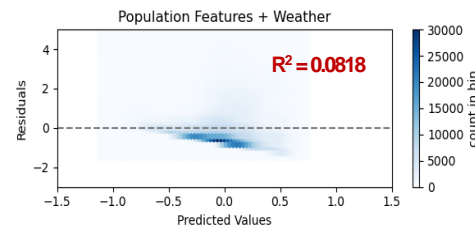
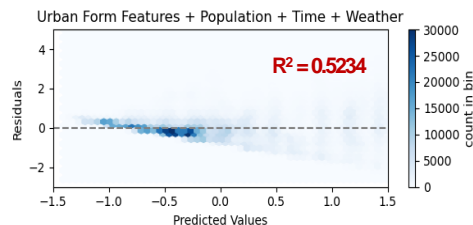
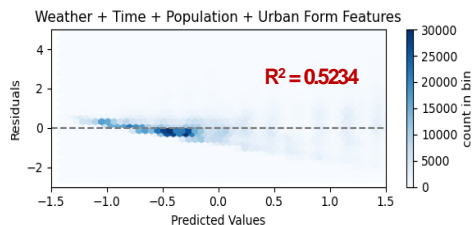
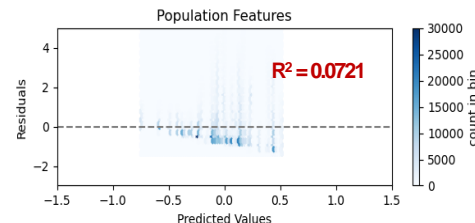
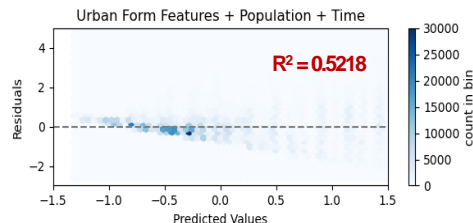
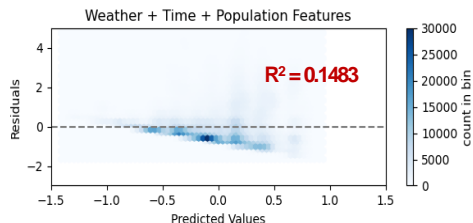
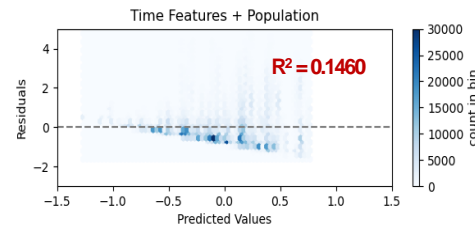
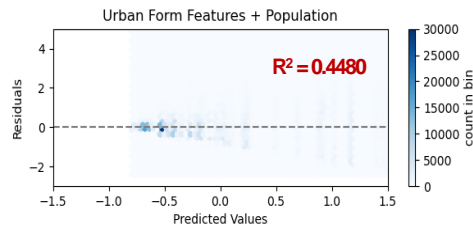
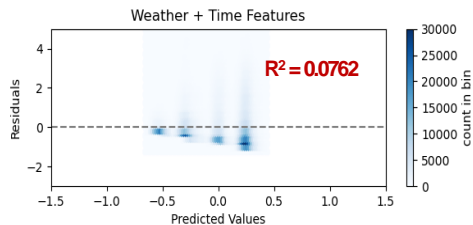
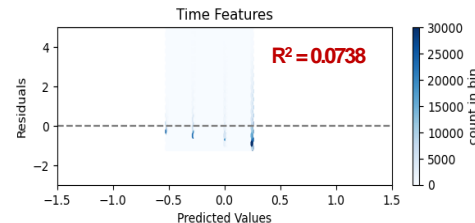
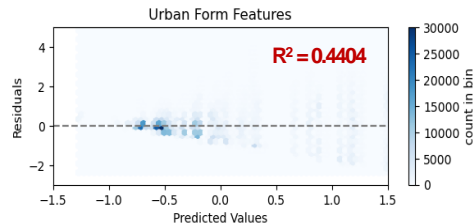
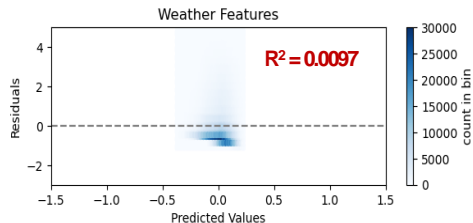
cc, speed_u, speed_v, rh, temp_K, ws

- Time features:**

day_type, morning_peak, lunch_peak, evening_peak

- Population features:**

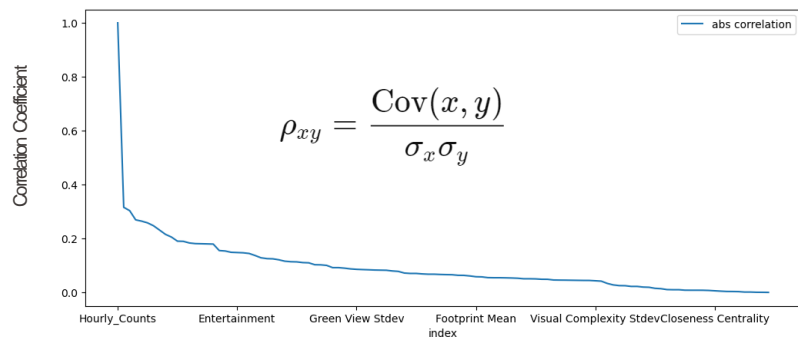
PopSum, Men, Women, Elderly, Youth, Children



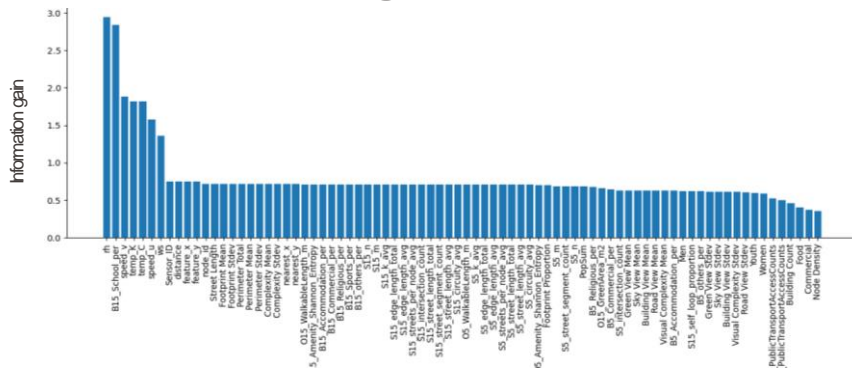
Baseline model: $R^2 = 0.523$

- **Correlation analysis**

Sorted Correlation Coefficients between Features and Pedestrian Count



- **Information gain**



Rule out the features with correlation coefficient < 0.01

➡ $R^2 = 0.509$

cc, S5_edge_length_avg, S5_street_length_avg, S5_circuitry_avg, S15_self_loop_proportion, Closeness Centrality, B5_others_per, B5_Accommodation_per, Institutional, speed_u, distance, Building View Stdev, Clustering (Weighted): 13

Rule out the features with information gain < 0.04

➔ $R^2 = 0.507$

Healthcare, day_type, lunch_peak, Degree, Clustering(Weighted),
Institutional, Civic, evening_peak, Clustering, Eigenvector Centrality,
PageRank, Social, n: 14

- Explored and analyzed urban features from open datasets.
- Developed mathematical models statistically explaining the spatial-temporal relationship between pedestrian flow and different groups of features.
- Implemented traditional machine learning techniques, improving the model R^2 up to **0.580**, with Luis F. Miranda-Moreno et al.^[1]
highest $R^2 = \mathbf{0.60}$.