

# **BiciMad Bike-Sharing System Dataset creation and Preliminary Hourly Demand Prediction using Machine and Deep Learning Approaches**

Erika Gutierrez  
Javier Roset

# INTRODUCTION



Barcelona School of Economics

# Bikesharing Systems



1) Bikes station



2) Bikes restocking

# Objective

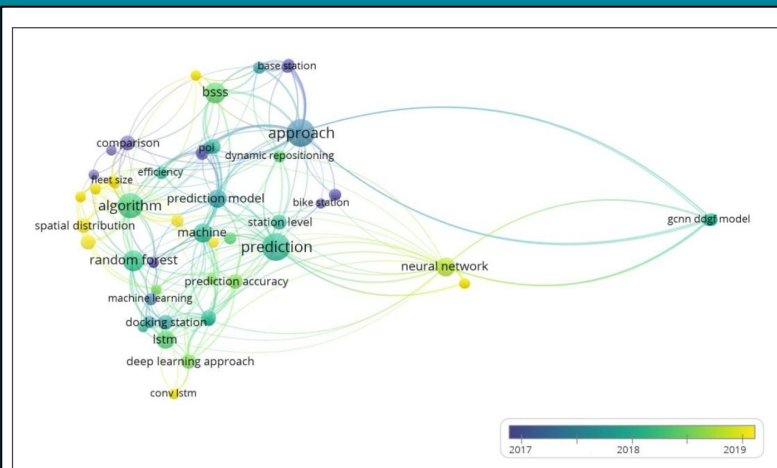
**Creation of a dataset** of BiciMAD bike-sharing system data **oriented to bike demand prediction** that covers from 2019 until the end of 2022, as well as a **preliminary prediction approach** using state-of-the-art techniques such as **XGBOOST** and **GCNN**

## Scope

- Creation of a coherent BiciMAD dataset
- Preliminary prediction study using state-of-the-art techniques (XGBOOST and GCNN)
  - Prediction of plugs and unplugs in each station at each hour
- Hyperparameter optimization of an XGBOOST
- Analysis of the results obtained

- data analysis to determine the quality of the final dataset
- Complete feature engineering to exploit all the potential of the data in terms of predicting power
- Research to find the best model to predict demand in BiciMAD bike-sharing system

# Literature



Papers used to guide our project:

- TS Kim et al. (2019) “Graph convolutional network approach applied to predict hourly bike-sharing demands considering spatial, temporal, and global effects.”
- R. Guo et al. (2019) "BikeNet: Accurate Bike Demand Prediction Using Graph Neural Networks for Station Rebalancing,"
- Lei Lin et al. (2019) “Predicting station-level hourly demand in a large-scale bike-sharing network: A graph convolutional neural network approach”

# DATA

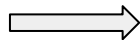


Barcelona School of Economics

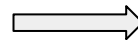
# Data Pipeline

DatosAbiertos  
Madrid

calendario.csv

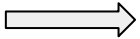


events.csv

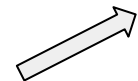


final\_dataset.csv

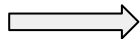
weather.csv



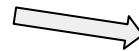
weather.csv



movements.csv  
movements.json



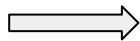
trips.csv



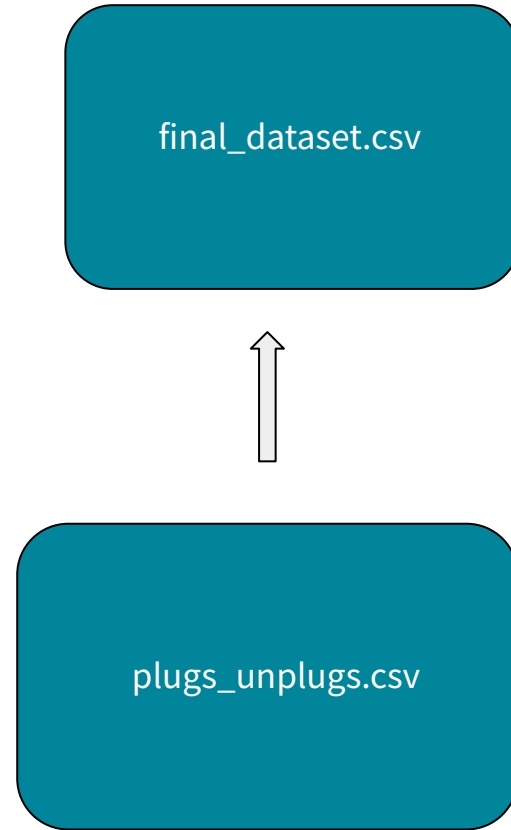
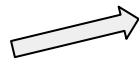
plugs\_unplugs.csv

EMT Madrid

stations.json



stations.csv



# Final Dataset Specifications

- Hourly level observations from January 1st 2019 to December 31st 2022 for 266 bike stations
- Dimensions:
  - 8, 185, 250 rows
  - 33 columns
- Transformed temporal data to sine and cosine expressions to capture cyclic behavior
- 2 target variables: bike **plugs** and **unplugs**
- Inclusion of weekly lag of target variables to account for weekly seasonality
- Standard scaling of all the features

## Example weather variables

Relative humidity  
Temperature  
Precipitation

## Example time variables

Week of year  
Hour  
Month

## Example bike variables

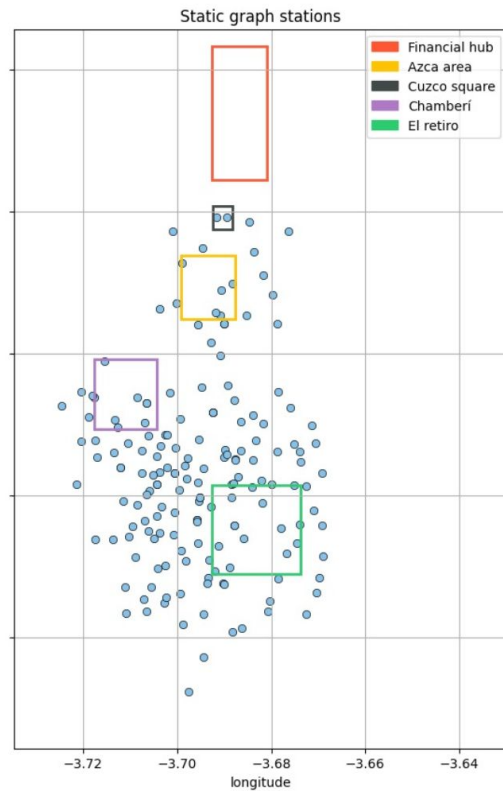
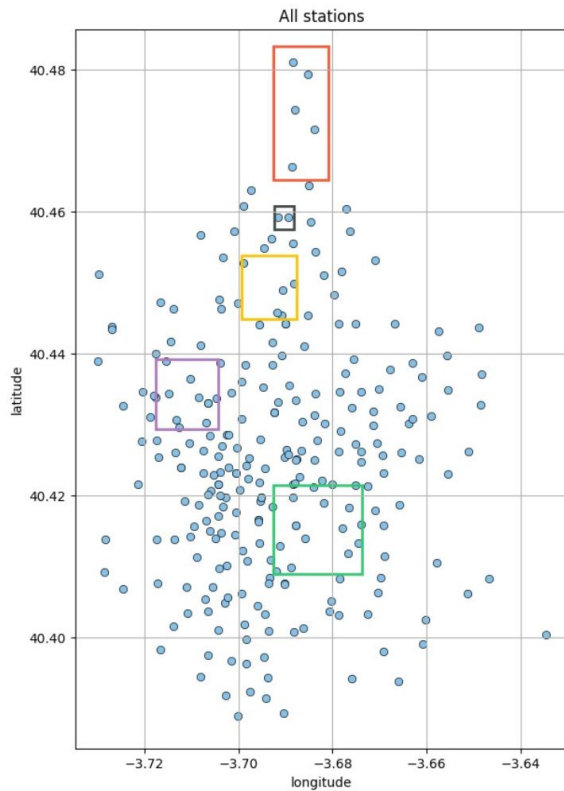
Number of reserved bikes  
Station availability  
Level of occupation

## Example events variables

Workday indicator  
COVID-19 indicator



# Spatial Component



# METHODOLOGY



Barcelona School of Economics

# Models

## Baseline and XGBOOST

### Baseline

---

- It predicts the value of the previous week at that time (day and hour)
- Will be used to compare our models with a trivial solution

### XGBOOST

---

- Provided best results of Innova
- Can handle seasonality (by adding lagged variables)
- Allow to obtain feature importance of the model

# Models

## Graphic Convolutional Neural Network (GCNN)

### GCNN

---

- Very popular in the literature
- Well suited for capturing spatial relationships and leveraging with graph structured data
- Can handle dynamic graphs (graphs with changing number of nodes over time)

### Considerations

---

- We will use static graphs

# Training process

## Training sets creation and hyperparameter search

### Splits analyzed

---

#### Innova split:

- Training: from 2019 to June 19th of 2021
- Test: from June 19th to July 3rd of 2021

#### 2022 split:

- Training: from 2019 to January 2nd of 2022
- Test: From January 2nd of 2022 to December 31st of 2022

### Hyperparameter search optimization



- Smart hyperparameter search, always looks to improve results
- Pruning stops useless trainings

# Training process

## Specificities of our approach

### XGBOOST

---

2 approaches for training:

- Use all data to train in one go
- Use a rolling window

### Considerations

---

- All stations were used
- Rolling window approach did not give very good results
- 1 separated model per target

### GCNN

---

1 approach for training:

- Use a rolling window

### Considerations

---

- Subsample of the data due to static graph approach
- Iteration lagging weather variables from next week to simulate weather forecast
- 1 separated model per target

# RESULTS



Barcelona School of Economics

# Results

## Metrics

### Test set Innova

MAE comparison for the Innova-TSN test set

| Model                       | Plugs | Unplugs |
|-----------------------------|-------|---------|
| Innova-TSN XGBOOST          | 1.106 | 1.106   |
| Our XGBOOST                 | 2.362 | 2.371   |
| GCNN without lagged weather | 2.844 | 2.84    |

### Test set 2022

Metrics of the models for plugs on the 2022 test set

| Model                       | RMSE  | MAE   | $R^2$ |
|-----------------------------|-------|-------|-------|
| Baseline                    | 2.831 | 1.573 | .077  |
| XGBOOST                     | 2.053 | 1.395 | .331  |
| GCNN with lagged weather    | 2.326 | 1.541 | .2    |
| GCNN without lagged weather | 2.364 | 1.584 | .168  |

Metrics of the models for unplugs on the 2022 test set

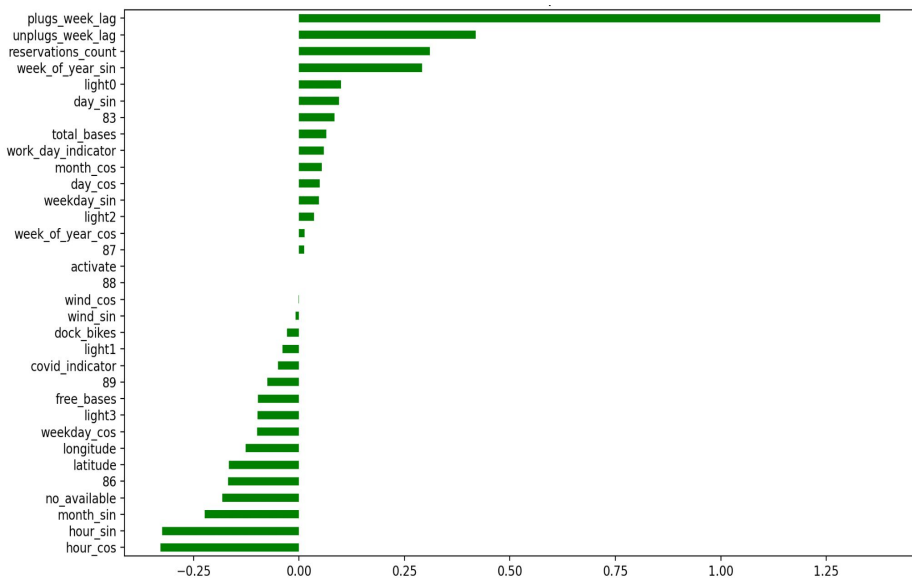
| Model                       | RMSE  | MAE   | $R^2$ |
|-----------------------------|-------|-------|-------|
| Basic Baseline              | 2.831 | 1.573 | .077  |
| XGBOOST                     | 2.219 | 1.513 | .273  |
| GCNN with lagged weather    | 2.464 | 1.635 | .151  |
| GCNN without lagged weather | 2.46  | 1.647 | .147  |



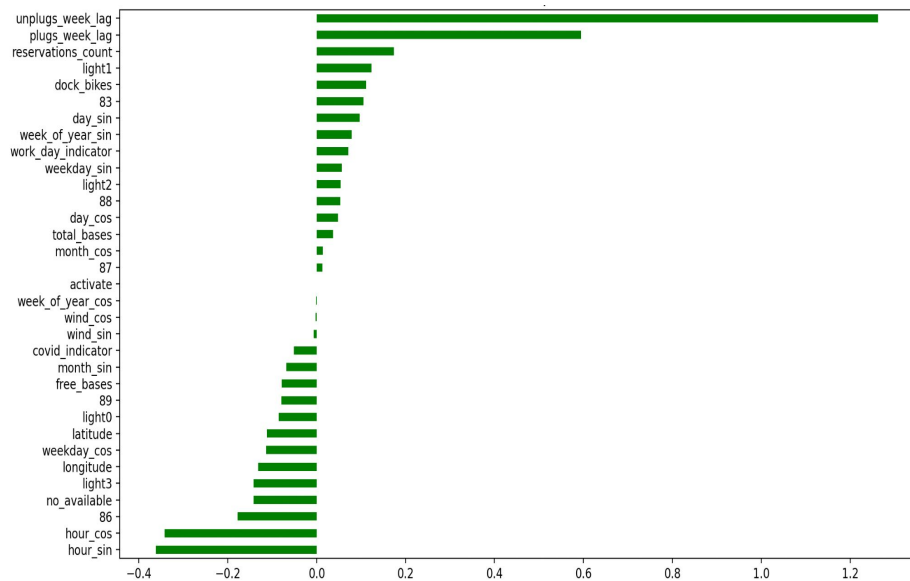
# Results

## Feature importances of XGBOOST

### Plugs



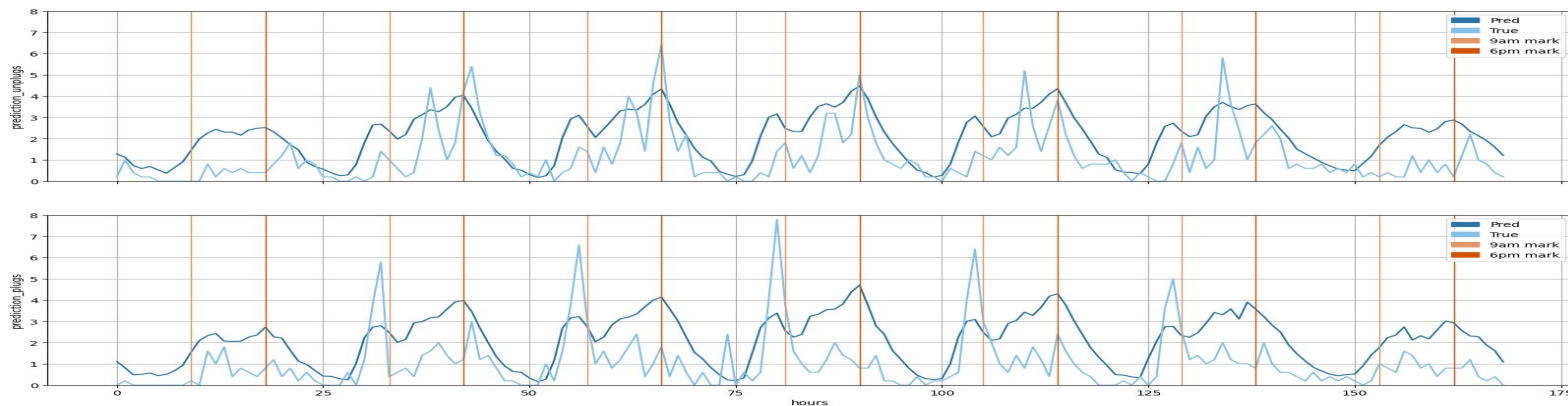
### Unplugs



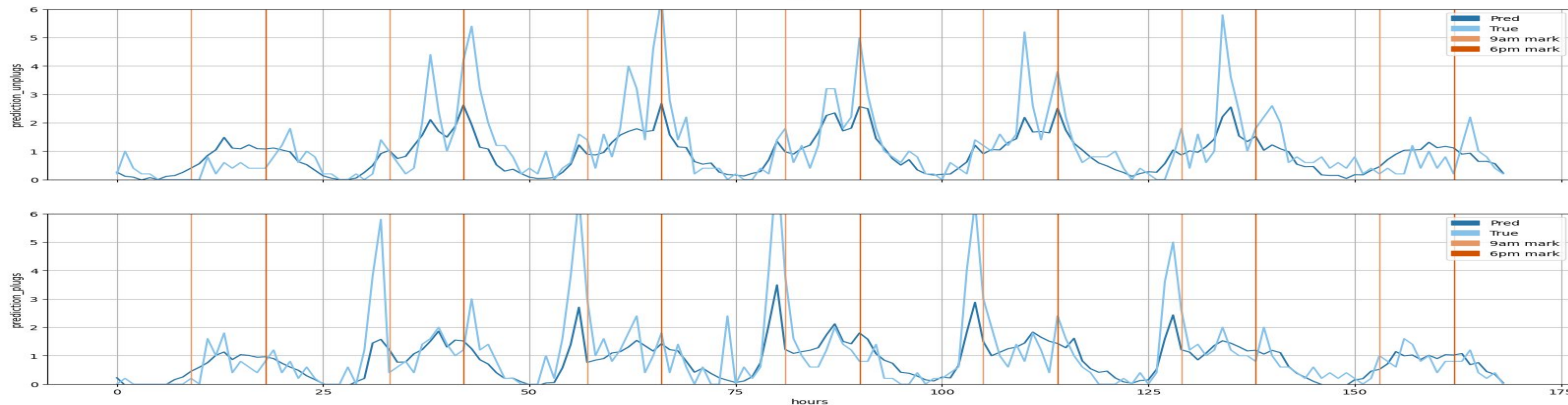
# Results

## Comparisons over Azca area (office district)

GCNN



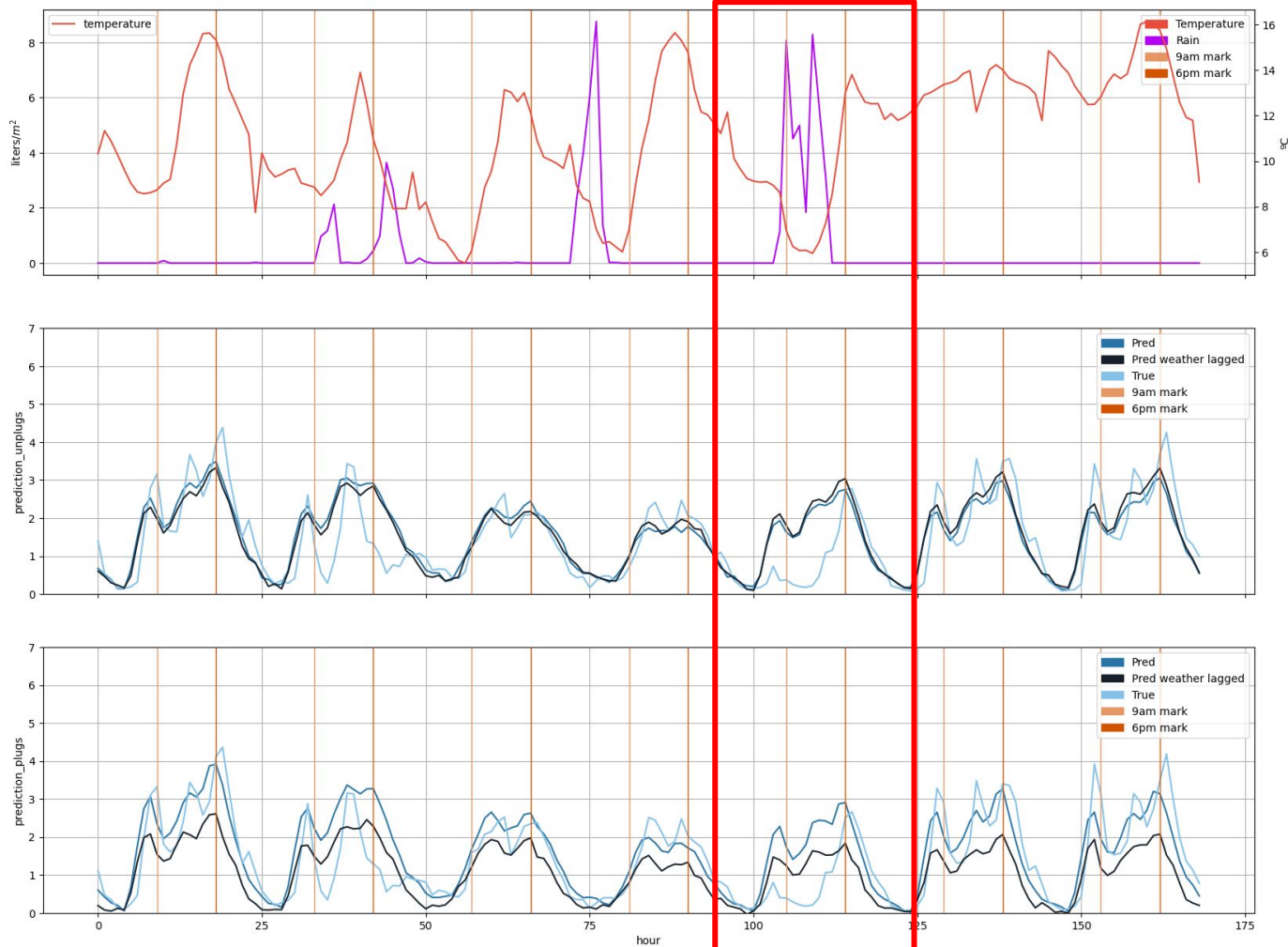
XGB



# Results

Comparisons over  
rain

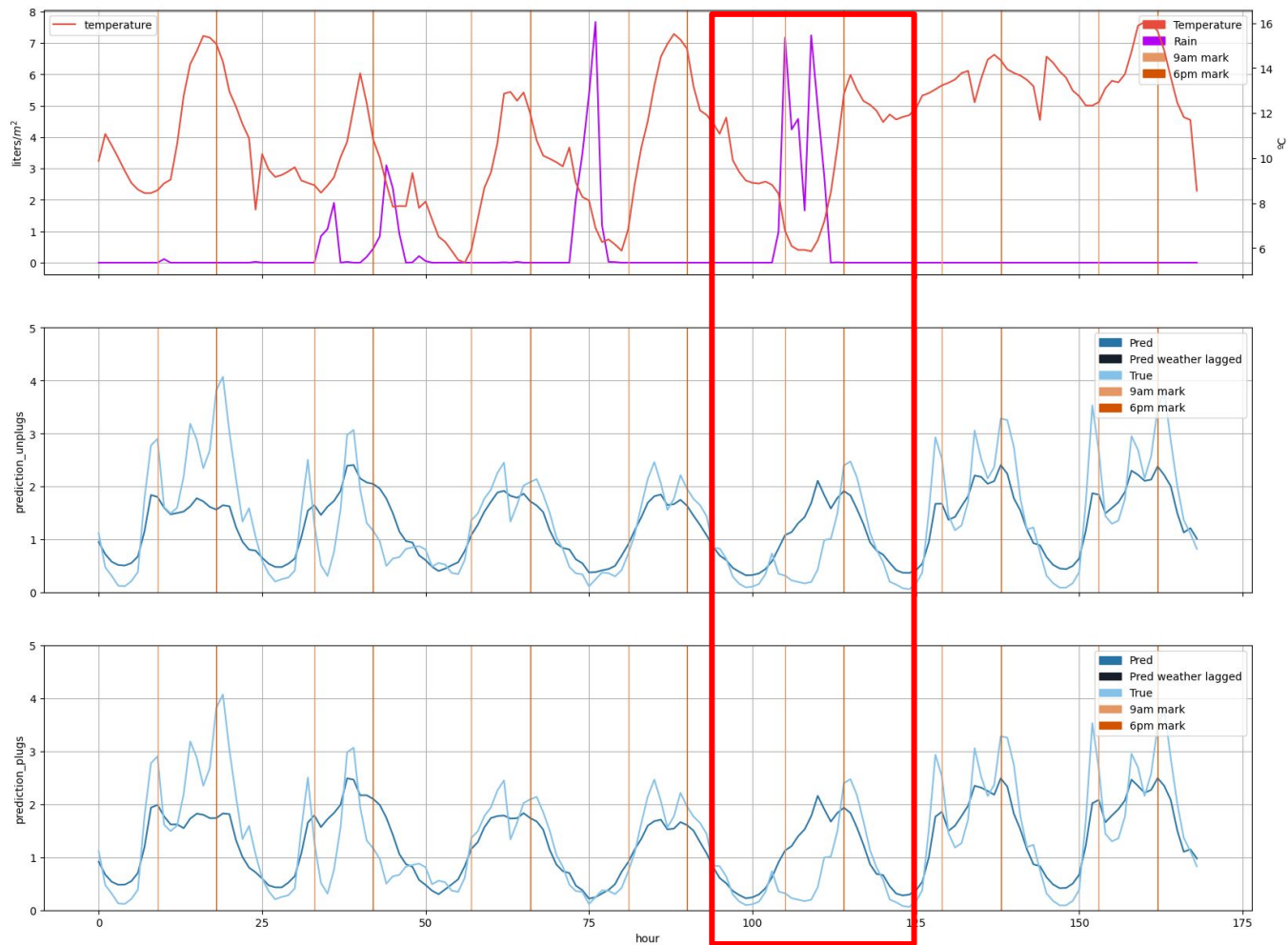
GCNN



# Results

Comparisons over  
rain

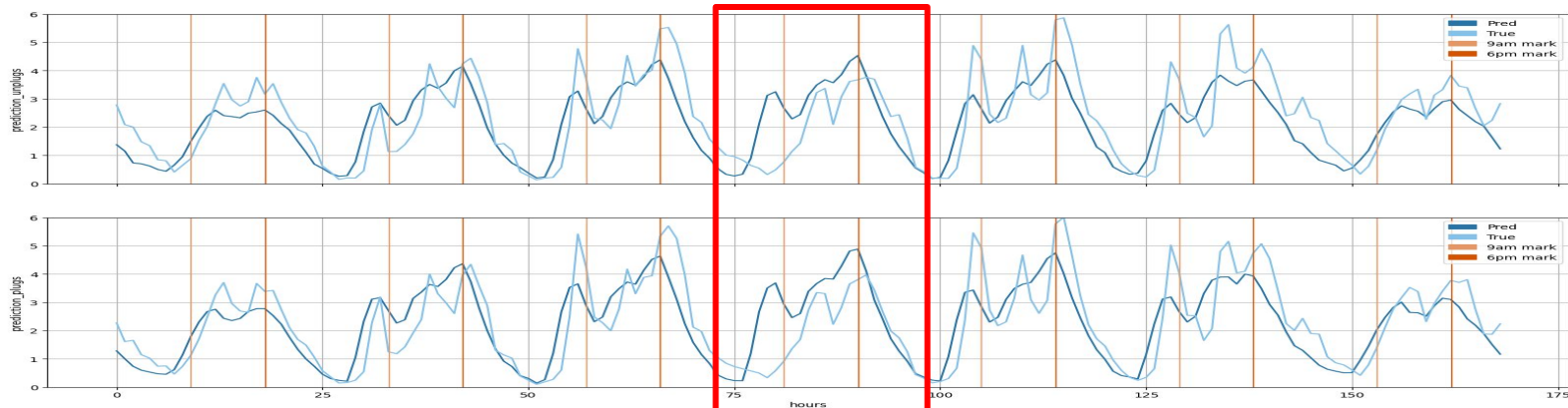
XGB



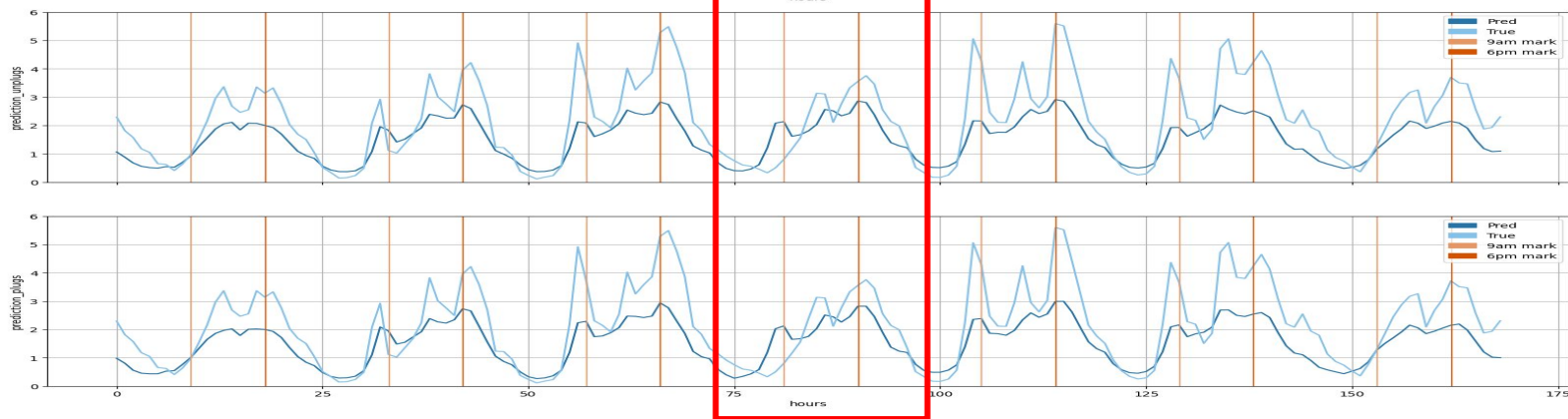
# Results

## Comparisons over “Día de la hispanidad”

GCNN



XGB



# CONCLUSIONS AND FUTURE LINES OF WORK



Barcelona School of Economics

# Conclusions

- We have created a coherent BiciMAD bikesharing system dataset oriented to plugs and unplugs demand prediction with hourly weather data per each station
- We have done a preliminary prediction study using SOA techniques (XGBOOST and GCNN)
- We have observed how the models capture the behaviors present in the data

# Future lines of work

- deep quality data assessment
- Enrich the dataframe with more variables (such as interactions with other public transports)
- Improve feature engineering to exploit the potential of the data
- Implement GCNN with dynamic graphs
- Work on AI explainability for GCNN



# THANK YOU FOR YOUR ATTENTION



Barcelona School of Economics

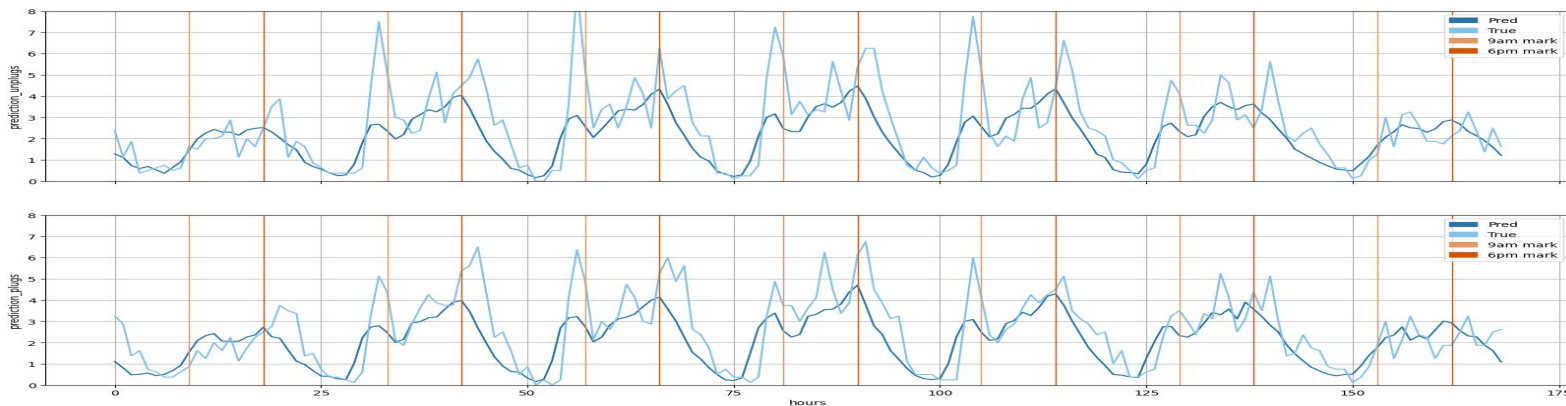
# APPENDIX



# Results

## Comparisons over Chamberi area (residential)

GCNN



XGB

