

# Taxi Trajectory Data ML Applications

Nikolaos Chiotis | MTN2221

Feb '23



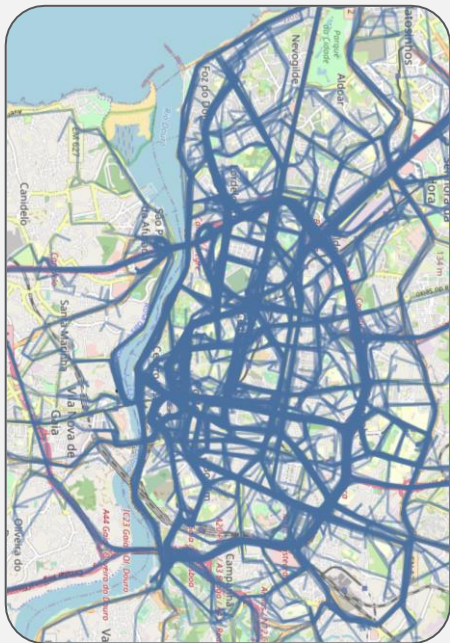
## Agenda

- 1 Dataset Overview
- 2 Data Preprocessing
- 3 Travel Duration Prediction
- 4 Travel Value Prediction

# 1. Dataset Overview



# 1. Dataset Overview



**0.4K**  
Taxis

- ✓ Trajectory data for **442 Taxis**
- ✓ City of **Porto, Portugal**
- ✓ **1 year of data** from July 2013 to June 2014

**1.7M**  
Rows

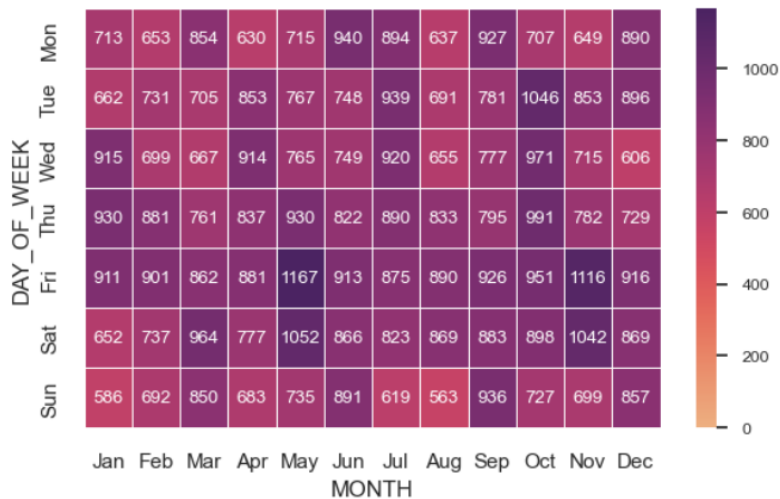
- ✓ Each row contains a **POLYLINE** for each taxi trip
- ✓ Trajectory has taxi's **coordinates every 15 seconds**
- ✓ **Undersampling to 100K rows**

**9**  
Columns

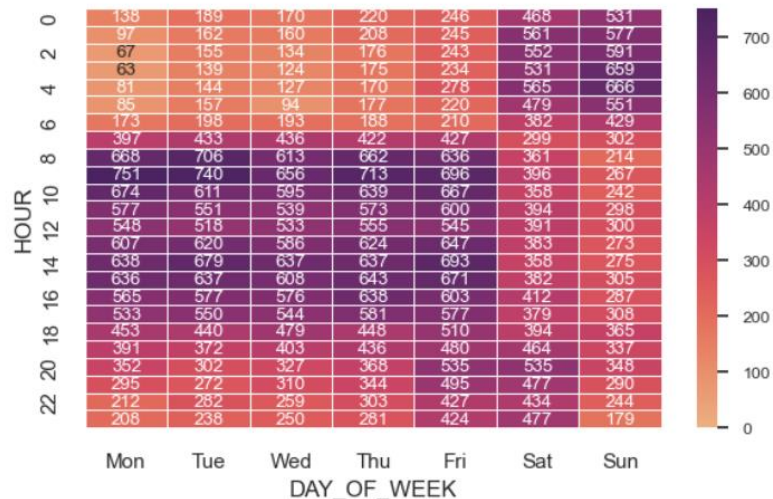
- |                      |                       |                       |
|----------------------|-----------------------|-----------------------|
| ✓ <b>TRIP_ID</b>     | ✓ <b>ORIGIN_STAND</b> | ✓ <b>DAYTYPE</b>      |
| ✓ <b>CALL_TYPE</b>   | ✓ <b>TAXI_ID</b>      | ✓ <b>MISSING_DATA</b> |
| ✓ <b>ORIGIN_CALL</b> | ✓ <b>TIMESTAMP</b>    | ✓ <b>POLYLINE</b>     |

# 1. Dataset Overview | Dates

## Day of week/Month Distribution



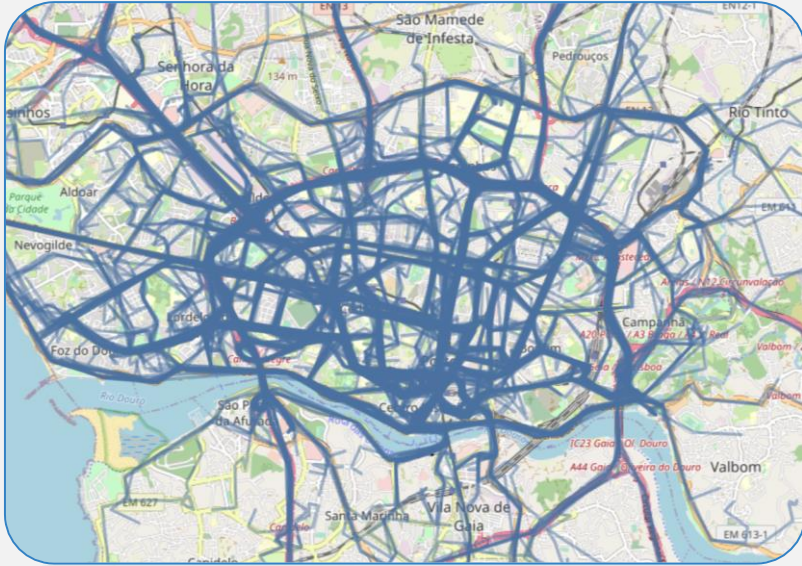
## Hour/Day of week Distribution



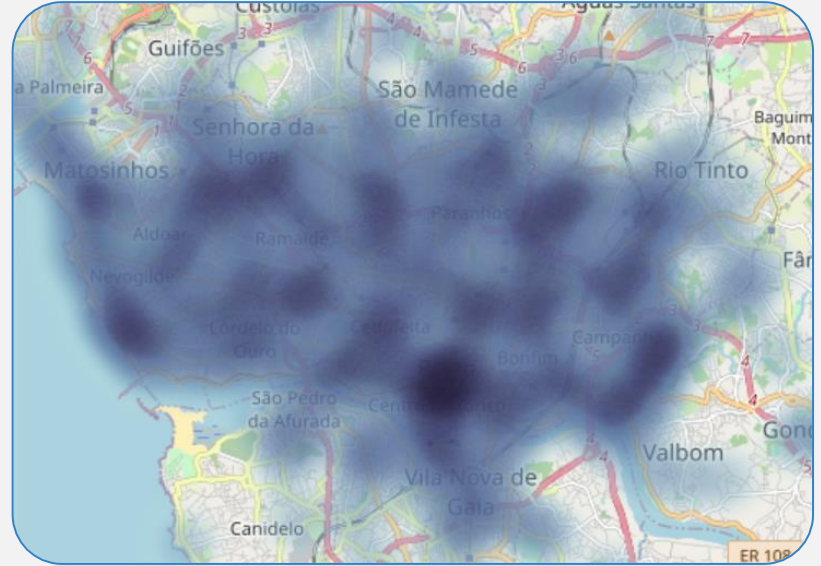


# 1. Dataset Overview | Routes & Starting Points

Routes Sample



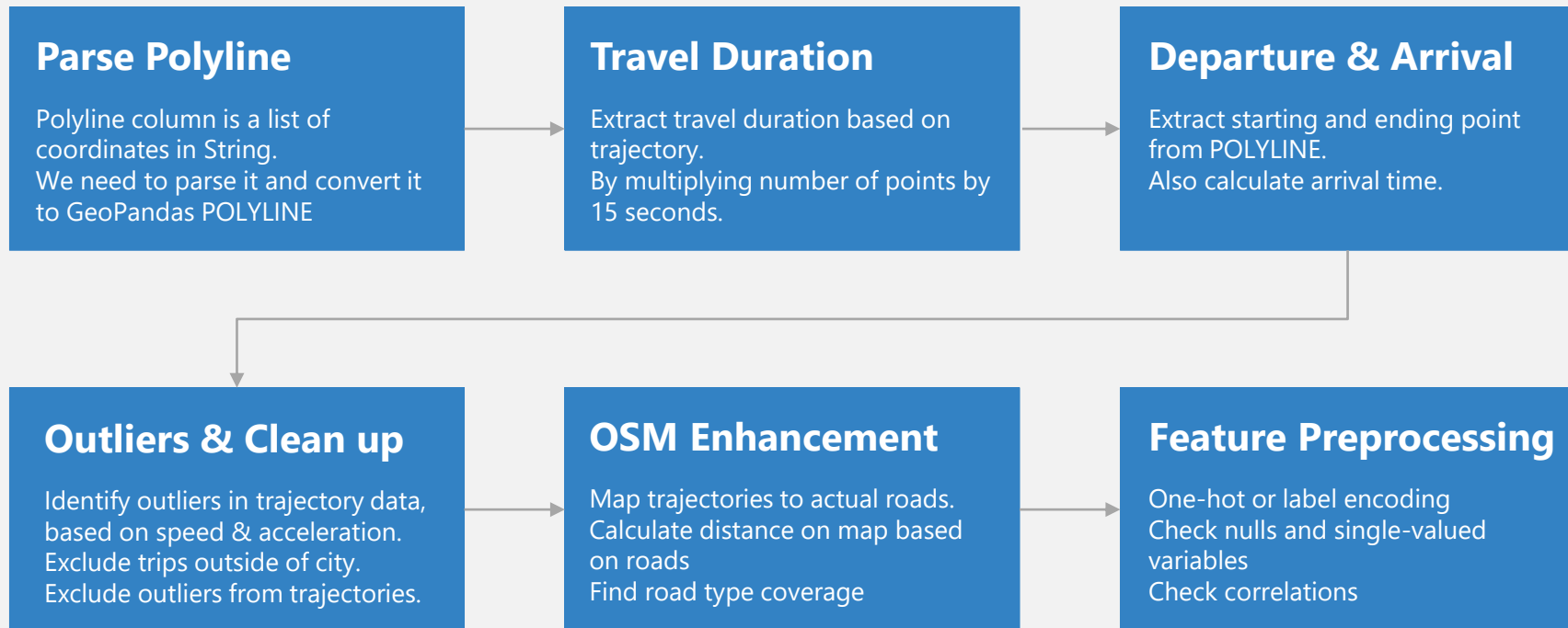
Starting Point Heatmap



## 2. Data Preprocessing



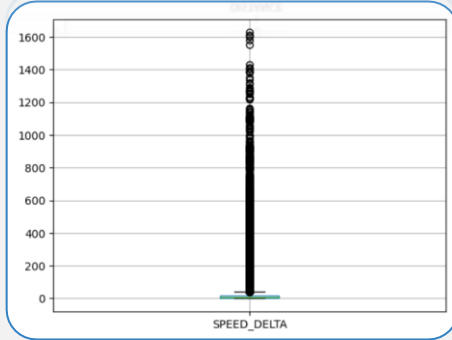
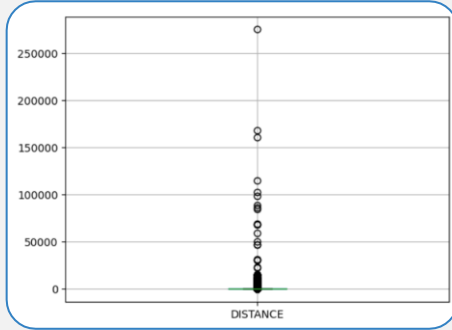
## 2. Data Preprocessing | Geolocation Preprocessing





## 2. Data Preprocessing | Outliers & Clean up

Boxplot



Outlier Example



## 2. Data Preprocessing | Open Street Map

### Load Open Street Map – Road Network



### Actions

- ✓ Load **road network**
- ✓ **Map trajectory** starting and ending points on **actual streets**
- ✓ Calculate **shortest paths**
- ✓ Extract extra feature:
  - **Route length**
  - **Type of road and number of lanes** at starting and ending point
  - **Percentage** of route coverage for each **road type** (roadway, residential etc)
  - **Number of turns**

### 3. Travel Duration Prediction



### 3. Travel Duration Prediction | Regression

Evaluation metric

**Mean square error**  
**Mean root square error**

\_\_\_\_\_ of \_\_\_\_\_

**Travel duration**

Algorithms

**K-Neighbors Regressor**  
**RandomForest Regressor**

\_\_\_\_\_ with \_\_\_\_\_

**GridSearchCV**

Results

**RandomForest Regressor**

with MRSE

\_\_\_\_\_ of \_\_\_\_\_

**6.37 minutes**

Outcomes

**Important**

**Features**

\_\_\_\_\_ were \_\_\_\_\_

**Route Distance & Hour**

## 4. Travel Value Classification





## 4. Travel Value Prediction | Classification

Target Variable

**Taxi Route Value Factor**

**HIGH VALUE: 0/1**

---

Based on the distance and time  
combination

Algorithms

**AdaBoost Classifier**

**GradientBoosting Classifier**

\_\_\_\_\_ in \_\_\_\_\_

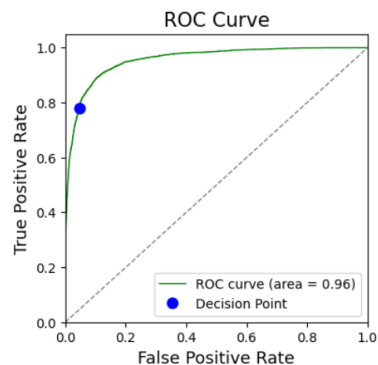
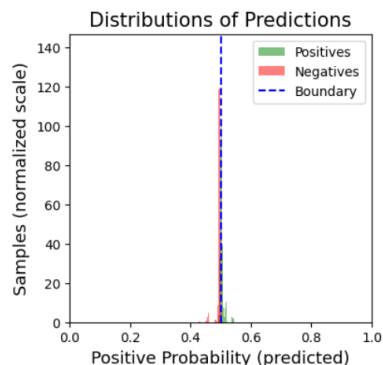
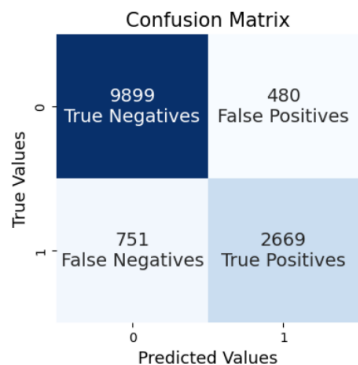
**Voting Classifier**

# 4. Travel Value Prediction | AdaBoost

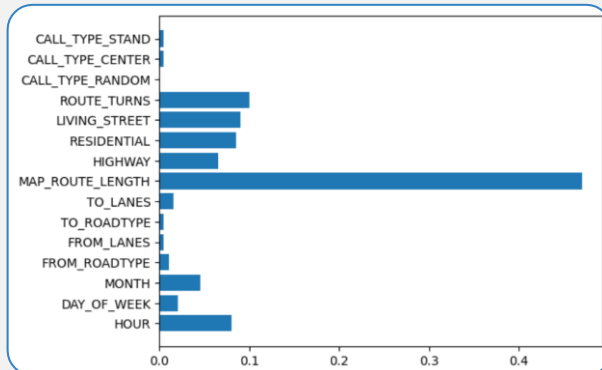
## AdaBoost Classification Results

```
evalBinaryClassifier(ada_clf, X_test, y_test)
```

	precision	recall	f1-score	support
0	0.93	0.95	0.94	10379
1	0.85	0.78	0.81	3420
accuracy			0.91	13799
macro avg	0.89	0.87	0.88	13799
weighted avg	0.91	0.91	0.91	13799



## Feature Importance

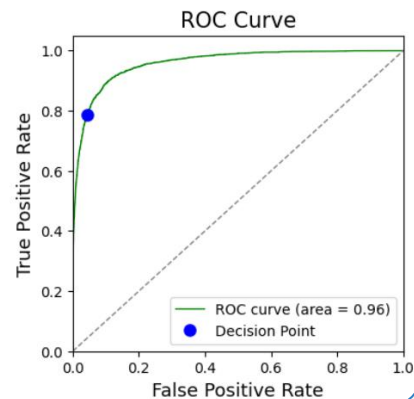
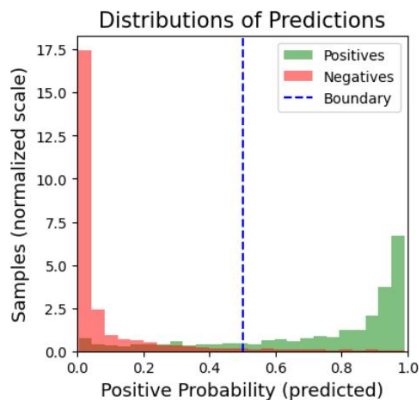
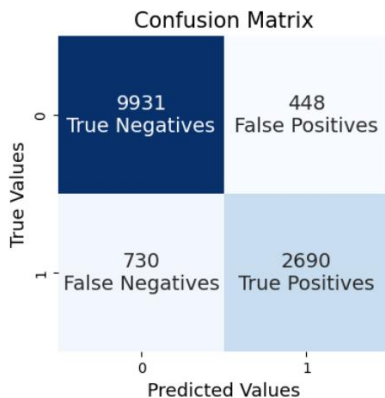


## 4. Travel Value Prediction | GradientBoosting

### GradientBoosting Classification Results

```
evalBinaryClassifier(hgb, X_test, y_test)
```

	precision	recall	f1-score	support
0	0.93	0.96	0.94	10379
1	0.86	0.79	0.82	3420
accuracy			0.91	13799
macro avg	0.89	0.87	0.88	13799
weighted avg	0.91	0.91	0.91	13799

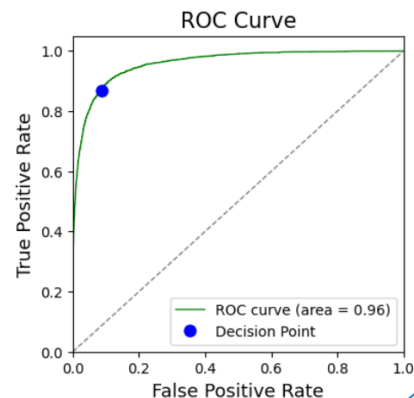
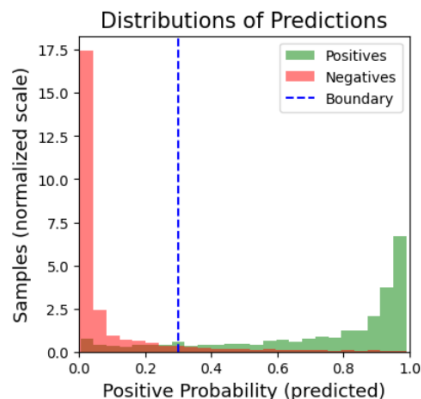
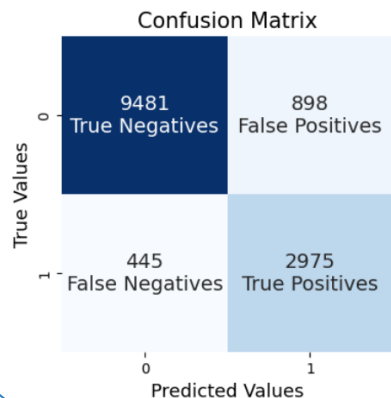


## 4. Travel Value Prediction | GradientBoosting (thres=0.3)

### GradientBoosting Classification Results

```
evalBinaryClassifier(hgb, X_test, y_test, thres=0.30)
```

	precision	recall	f1-score	support
0	0.96	0.91	0.93	10379
1	0.77	0.87	0.82	3420
accuracy			0.90	13799
macro avg	0.86	0.89	0.87	13799
weighted avg	0.91	0.90	0.90	13799

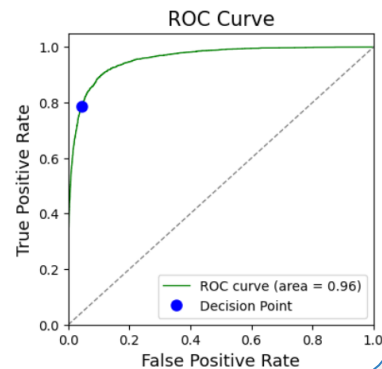
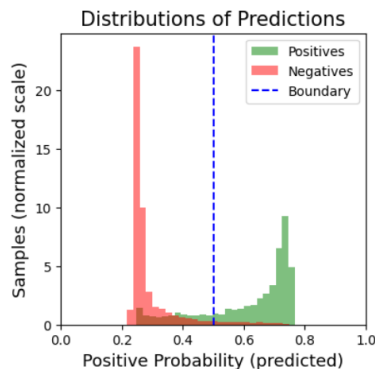
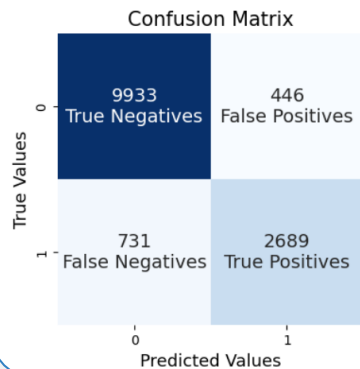


## 4. Travel Value Prediction | Voting Classifier

### Voting Classification Results

```
clf = VotingClassifier(estimators=[('ada_clf', ada_clf), ('hgb', hgb)], voting='soft')
clf.fit(X_train, y_train)
evalBinaryClassifier(clf, X_test, y_test)
```

	precision	recall	f1-score	support
0	0.93	0.96	0.94	10379
1	0.86	0.79	0.82	3420
accuracy			0.91	13799
macro avg	0.89	0.87	0.88	13799
weighted avg	0.91	0.91	0.91	13799





## Possible Enhancements

Graph ML to predict road traffic flow

Outlier detection using ML on trajectories

Apply node embeddings for route predictions

**Thank you!**

**Questions?**

