

Table of Contents

- 1. Objectives
- 2. Data
- 3. Modeling
- 4. Conclusions
- 5. Appendix

1. Objectives

The report describes dataset, preparation and modeling done to predict taxi demands for Manhattan and the airports.

2. Data

The dataset used in yellow taxi trip dataset for January 2015. Due to huge size that cannot be processed by current computational power, I opt for smaller dataset as a sample to represent the whole 2015 dataset.

The data is checked for missing values and duplicates. I had taken out 6 regions from the Jan 2015 dataset that concentrate these areas which the client requested.

I have analysed by grouping into dates according to each zone to explore duration, fare amount and passenger counts.

I have dropped unwanted features that are useless, created new features like duration, is Holiday, Day of Pickup and Demand as target.

Statistical tests ANOVA and Chi-Square tests are performed to find out the p-values.

As for correlation, total amount is highly correlates with fare amount and tip amount.

Test train split is performed with 80% training and 20% testing.

Modeling

The model will be multi-classification model since there are demand results are Low, Medium and High.

Two scenarios created:

- 1. Baseline Model
- 2. Tweaked Model to reduce false negatives in Low Demand

I used Logistic Regression and XGBoost for creating models in this project. Hyperparameter tuning is also done for XGBoost. There is no balancing of target feature done for demand. Crossvalidation is also done to ensure accuracy as a whole.

Metrics used are accuracy, recall, precision and F1-score.

Conclusion

Both models gave perfect score for classifying taxi demand, hence logistic regression model is recommended for simplicity.

Appendix: Code Section

Modeling Tasks

Good work so far preparing the data to classify taxi demand and evaluating predictors. Now you can get started with the modeling!

There are a couple scenarios to address below. As usual for classification modeling, they now you need to recall the last following for each of your best/final models:

- Accuracy
- Confusion matrix
- cMetrics output (or equivalent)

Scenario 1:

Start off with a baseline model that emphasizes overall accuracy. Use your test data set to verify your modeling results. This model will be useful for analysis and comparison later. It may be also a good way to investigate model types and hyperparameters, class

imbalance, and which features are most useful.

Scenario 2 Main Points:

- Emphasize reducing false negatives for Low demand, especially Low demand classified as High demand.
- Also, try to reduce false positives for Low, especially missed High.
- Some increase in false positives for Low demand is an acceptable trade-off to reduce false negatives for Low demand.
- Medium/High demand misclassification is not a priority to reduce.
- Use reasonable modifications to the cost matrix to achieve your goals. You'll need to include a discussion of your customizations in your report.
- Use your test data set to verify your modeling results.

Import Libraries

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import sklearn
import shap

from xgboost import XGBClassifier, XGBRegressor
from xgboost import to_graphviz, plot_importance

from sklearn.linear_model import ElasticNet, Lasso, LinearRegression, LogisticRegression, Ridge
from sklearn.tree import DecisionTreeClassifier, DecisionTreeRegressor, ExtraTreeClassifier, ExtraTreeRegressor
from sklearn.ensemble import RandomForestClassifier, RandomForestRegressor, GradientBoostingClassifier, GradientBoostingRegressor
from sklearn.svm import SVC, LinearSVC, NuSVC, LinearSVR

#matplotlib inline
sns.set_style('dark')
sns.set(font_scale=1.2)

from sklearn.pipeline import Pipeline
from sklearn.model_selection import GridSearchCV, RepeatedStratifiedFold
from sklearn.feature_selection import RFECV, SelectKBest, f_classif, f_regression, chi2

from sklearn.inspection import permutation_importance
from sklearn.model_selection import cross_val_score, train_test_split, GridSearchCV, RandomizedSearchCV
from sklearn.preprocessing import LabelEncoder, StandardScaler, MinMaxScaler, OneHotEncoder
from sklearn.pipeline import Pipeline
from sklearn.metrics import confusion_matrix, classification_report, mean_absolute_error, mean_squared_error
from sklearn.metrics import plot_confusion_matrix, plot_precision_recall_curve, plot_roc_curve, accuracy_score
from sklearn.metrics import auc, f1_score, precision_score, recall_score, roc_auc_score

from imblearn.under_sampling import RandomUnderSampler
from imblearn.over_sampling import RandomOverSampler
from imblearn.over_sampling import SMOTE

import warnings
warnings.filterwarnings('ignore')

import pickle
from pickle import dump, load

np.random.seed(0)

from pycoar.classification import *
```

pd.set_option('display.max_columns',100)
%pd.set_option('display.max_rows',100)
pd.set_option('display.width', 1000)
np.set_printoptions(suppress=True)

Data Exploration and Analysis

```
In [2]: df = pd.read_csv('Finaltrain2.csv')

In [3]: df

Out[3]:
```

	passenger_count	trip_distance	pickup_location	fare_amount	extra	tip_amount	total_amount	duration	pickupday	demand	ishol
0	2	12.74	JFK Airport	47.5	0.5	9.60	58.40	24.0	29	low	
1	1	17.50	JFK Airport	52.0	0.0	7.00	65.13	32.0	1	low	
2	2	18.50	JFK Airport	51.0	0.5	0.00	57.63	30.0	18	low	
3	1	14.80	JFK Airport	41.0	0.0	0.00	41.80	29.0	13	low	
4	1	14.30	JFK Airport	39.5	0.5	0.00	40.80	22.0	29	low	
...
68563	1	1.69	Upper East Side	8.0	0.0	0.00	8.80	8.0	10	low	
68564	1	1.80	Upper East Side	12.5	0.0	2.00	15.30	18.0	3	low	
68565	5	1.54	Upper East Side	7.5	1.0	1.70	11.00	7.0	15	high	
68566	1	0.78	Upper East Side	5.5	1.0	1.25	8.55	5.0	15	low	
68567	1	1.01	Upper East Side	7.0	1.0	1.60	10.40	8.0	15	low	

68568 rows x 12 columns

```
In [4]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 68568 entries, 0 to 68567
Data columns (total 12 columns):
#   Column                Non-Null Count  Dtype
---  --
0   passenger_count        68568 non-null    int64
1   trip_distance          68568 non-null    float64
2   pickup_location        68568 non-null    object
3   fare_amount            68568 non-null    float64
4   extra                  68568 non-null    float64
5   tip_amount             68568 non-null    float64
6   total_amount           68568 non-null    float64
7   duration                68568 non-null    float64
8   pickupday              68568 non-null    int64
9   demand                 68568 non-null    object
10  isholiday              68568 non-null    int64
11  label                   68568 non-null    int64
dtypes: float64(6), int64(4), object(2)
memory usage: 6.3+ MB

In [5]: df.describe(include='all')

Out[5]:
```

	passenger_count	trip_distance	pickup_location	fare_amount	extra	tip_amount	total_amount	duration	pickupday	demand	ishol
count	68568.000000	6856800e+04	68568	68568.000000	68568.000000	68568.000000	68568.000000	68568.000000	68568.000000	68568.000000	
unique	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
top	NaN	NaN	Midtown	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
freq	1.694044	5.194054e+01	39799	13.839213	0.285746	1.767108	17.174227	14.549921	15.838000		
std	1.346843	1.267113e+04	NaN	13.059612	0.365959	2.775889	16.001486	36.983077	8.716981		
min	0.000000	0.000000e+00	NaN	-52.000000	-1.000000	0.000000	-52.800000	0.000000	1.000000		
25%	1.000000	1.000000e+00	NaN	6.500000	0.000000	0.000000	8.160000	6.000000	8.000000		
50%	1.000000	1.200000e+00	NaN	9.000000	0.000000	1.000000	11.300000	10.000000	15.000000		
75%	2.000000	3.400000e+00	NaN	14.500000	0.500000	2.160000	17.800000	17.000000	23.000000		
max	6.000000	3.318000e+06	NaN	570.080000	1.000000	94.510000	590.380000	1440.000000	31.000000		

```
In [6]: df.shape

Out[6]: (68568, 12)

In [7]: df.columns

Out[7]: Index(['passenger_count', 'trip_distance', 'pickup_location', 'fare_amount', 'extra', 'tip_amount', 'total_amount', 'duration', 'pickupday', 'demand', 'isholiday', 'label'], dtype='object')

In [8]: df[pickup_location == 'Midtown'].head()

Out[8]:
```

	passenger_count	trip_distance	pickup_location	fare_amount	extra	tip_amount	total_amount	duration	pickupday	demand	ishol
0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	0
...
68563	0	0	0	0	0	0	0	0	0	0	0
68564	0	0	0	0	0	1	1	0	0	0	0
68565	0	0	0	0	0	1	1	0	0	0	0
68566	0	0	0	0	0	1	1	0	0	0	0
68567	0	0	0	0	0	1	1	0	0	0	0

68568 rows x 12 columns

```
In [10]: df2 = pd.concat([df, df[pickup_location == 'Midtown']], axis=1)

Out[10]:
```

	passenger_count	trip_distance	pickup_location	fare_amount	extra	tip_amount	total_amount	duration	pickupday	demand	ishol
0	2	12.74	JFK Airport	47.5	0.5	9.60	58.40	24.0	29	low	
1	1	17.50	JFK Airport	52.0	0.0	7.00	65.13	32.0	1	low	
2	2	18.50	JFK Airport	51.0	0.5	0.00	57.63	30.0	18	low	
3	1	14.80	JFK Airport	41.0	0.0	0.00	41.80	29.0	13	low	
4	1	14.30	JFK Airport	39.5	0.5	0.00	40.80	22.0	29	low	
...
68563	1	1.69	Upper East Side	8.0	0.0	0.00	8.80	8.0	10	low	
68564	1	1.80	Upper East Side	12.5	0.0	2.00	15.30	18.0	3	low	
68565	5	1.54	Upper East Side	7.5	1.0	1.70	11.00	7.0	15	high	
68566	1	0.78	Upper East Side	5.5	1.0	1.25	8.55	5.0	15	low	
68567	1	1.01	Upper East Side	7.0	1.0	1.60	10.40	8.0	15	low	

68568 rows x 16 columns

```
In [12]: df2.columns

Out[12]: Index(['passenger_count', 'trip_distance', 'pickup_location', 'fare_amount', 'extra', 'tip_amount', 'total_amount', 'duration', 'pickupday', 'demand', 'isholiday', 'LaGuardia Airport', 'Lower Manhattan', 'Midtown', 'Upper East Side'], dtype='object')

In [13]: df2.drop(['pickup_location', 'label'], axis=1, inplace=True)

Out[13]:
```

	passenger_count	trip_distance	fare_amount	extra	tip_amount	total_amount	duration	pickupday	demand	isholiday	LaGuardia Airport	Lower Manhattan	Midtown	Upper East Side
0	2	12.74	47.5	0.5	9.6	58.40	24.0	29	low	0	0	0	0	
1	1	17.50	52.0	0.0	7.0	65.13	32.0	1	low	1	1	0	0	
2	2	18.50	51.0	0.5	0.0	57.63	30.0	18	low	0	0	0	0	
3	1	14.80	41.0	0.0	0.0	41.80	29.0	13	low	0	0	0	0	
4	1	14.30	39.5	0.5	0.0	40.80	22.0	29	low	0	0	0	0	
...	
68563	0	0	0	0	0	0	0	0	0	0	0	0	0	
68564	0	0	0	0	0	1	1	0	0	0	0	0	0	
68565	0	0	0	0	0	1	1	0	0	0	0	0	0	
68566	0	0	0	0	0	1	1	0	0	0	0	0	0	
68567	0	0	0	0	0	1	1	0	0	0	0	0	0	

68568 rows x 16 columns

```
In [14]: df2.head()

Out[14]:
```

	passenger_count	trip_distance	fare_amount	extra	tip_amount	total_amount	duration	pickupday	demand	isholiday	LaGuardia Airport	Lower Manhattan	Midtown	Upper East Side
0	2	12.74	47.5	0.5	9.6	58.40	24.0	29	low	0	0	0	0	
1	1	17.50	52.0	0.0	7.0	65.13	32.0	1	low	1	1	0	0	
2	2	18.50	51.0	0.5	0.0	57.63	30.0	18	low	0	0	0	0	
3	1	14.80	41.0	0.0	0.0	41.80	29.0	13	low	0	0	0	0	
4	1	14.30	39.5	0.5	0.0	40.80	22.0	29	low	0	0	0	0	
...	
68563	0	0	0	0	0	0	0	0	0	0	0	0	0	
68564	0	0	0	0	0	1	1	0	0	0	0	0	0	
68565	0	0	0	0	0	1	1	0	0	0	0	0	0	
68566	0	0	0	0	0	1	1	0	0	0	0	0	0	
68567	0	0	0	0	0	1	1	0	0	0	0	0	0	

68568 rows x 16 columns

```
In [15]: df2.columns

Out[15]: Index(['passenger_count', 'trip_distance', 'fare_amount', 'extra', 'tip_amount', 'total_amount', 'duration', 'pickupday', 'demand', 'isholiday', 'LaGuardia Airport', 'Lower Manhattan', 'Midtown', 'Upper East Side'], dtype='object')

In [16]: df2 = df2[['passenger_count', 'trip_distance', 'fare_amount', 'extra', 'tip_amount', 'total_amount', 'duration', 'pickupday', 'demand', 'isholiday', 'LaGuardia Airport', 'Lower Manhattan', 'Midtown', 'Upper East Side', 'demand']]

Out[16]:
```

	passenger_count	trip_distance	fare_amount	extra	tip_amount	total_amount	duration	pickupday	demand	isholiday	LaGuardia Airport	Lower Manhattan	Midtown	Upper East Side	demand
0	2	12.74	47.5	0.5	9.6	58.40	24.0	29	low	0	0	0	0	0	
1	1	17.50	52.0	0.0	7.0	65.13	32.0	1	low	1	1	0	0	0	
2	2	18.50	51.0	0.5	0.0	57.63	30.0	18	low	0	0	0	0	0	
3	1	14.80	41.0	0.0	0.0	41.80	29.0	13	low	0	0	0	0	0	
4	1	14.30	39.5	0.5	0.0	40.80	22.0	29	low	0	0	0	0	0	
...	
68563	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
68564	0	0	0	0	0	1	1	0	0	0	0	0	0	0	
68565	0	0	0	0	0	1	1	0	0	0	0	0	0	0	
68566	0	0	0	0	0	1	1	0	0	0	0	0	0	0	
68567	0	0	0	0	0	1	1	0	0	0	0	0	0	0	

68568 rows x 16 columns

```
In [17]: df2.head()

Out[17]:
```

	passenger_count	trip_distance	fare_amount	extra	tip_amount	total_amount	duration	pickupday	demand	isholiday	LaGuardia Airport	Lower Manhattan	Midtown	Upper East Side	demand
0	2	12.74	47.5	0.5	9.6	58.40	24.0	29	low	0	0	0	0	0	
1	1	17.50	52.0	0.0	7.0	65.13	32.0	1	low	1	1	0	0	0	
2	2	18.50	51.0	0.5	0.0	57.63	30.0	18	low	0	0	0	0	0	
3	1	14.80	41.0	0.0	0.0	41.80	29.0	13	low	0	0	0	0	0	
4	1	14.30	39.5	0.5	0.0	40.80	22.0	29	low	0	0	0	0	0	
...	
68563	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
68564	0	0	0	0	0	1	1	0	0	0	0	0	0	0	
68565	0	0	0												

