

## FINAL PROJECT

EECS3401

INTRODUCTION TO ARTIFICIAL INTELLIGENCE AND LOGIC PROGRAMMING

---

# Predicting Uber Fares with the Uber Fares Dataset

---

*Authors:*

Group 36

Kanwar Partap Pannu

Gabriel Rincon

Abdullah Sajjad

*Student Numbers:*

218918938

218624734

217250945

Wednesday 22<sup>nd</sup> November, 2023

# Contents

<b>1</b>	<b>Framing the Problem and Looking at the Big Picture</b>	<b>3</b>
<b>2</b>	<b>Exploratory Data Analysis</b>	<b>3</b>
2.1	Description of the Data Set . . . . .	3
2.2	Graph Analysis . . . . .	4
<b>3</b>	<b>Data Preparation</b>	<b>6</b>
3.1	Data Cleaning . . . . .	6
3.2	Pipe lining and Preparation for Regression . . . . .	6
<b>4</b>	<b>Training and Comparison</b>	<b>6</b>
<b>5</b>	<b>Best Performing Algorithm</b>	<b>7</b>
<b>6</b>	<b>Limitations</b>	<b>7</b>
<b>A</b>		
	<b>Source Code</b>	<b>9</b>
<b>B</b>		
	<b>Links</b>	<b>15</b>
B.1	Video . . . . .	15
B.2	Dataset . . . . .	15
B.3	GitHub . . . . .	15

# 1 Framing the Problem and Looking at the Big Picture

Uber took the world by storm with it's founding in 2009 and quickly became one of the most used forms of transportation across North America and much of rest of the world. It appeals to many thanks to its convenience and accessibility. When calculating the price of an Uber ride many factors are considered such as time, distance etc. but no clear formula can tell you your fare ahead of time. The ultimate goal of the Uberfares data set is to more accurately predict and understand the Uber pricing model. For Uber it is very important to understand how they're pricing models work to help Uber coordinate any initiatives relating to the pricing model and to ensure better compliance with local regulations. Better fare predictors can also allow drivers to better asses their potential earnings from a certain ride and ensure they can make the best decisions. Most importantly however the user experience would vastly improve, with a more accurate fare the user can plan for trips ahead of time and perform cost-benefit analyses with other options. Framing the problem from the dataset is quite clear and straight forward. This dataset is a supervised learning problem because the data is labelled with the target variable being the fare amount. Since the target variable is a continuous numerical value which we are trying to predict this is a regression task. This dataset is not updated live with a continuous flow of data thus we can train the the model on the whole dataset at once with batch learning.

## 2 Exploratory Data Analysis

### 2.1 Description of the Data Set

The Uber Fares data set includes the fares for many different unique trips mostly centered around New York although some data points are in very far away regions. It contains 200000 entries before data cleaning. Analysis done in this section is done on the cleaned dataset, the details of this cleanup are in section 3. The following fields are included in the orginal data set:

- key - A unique identifier for each trip
- fare\_amount - fare cost in USD
- pickup\_datetime - Date and time at the start of the trip
- passenger\_count - Number of passengers inside the vehicle according to the driver
- pickup\_longitude - Longitude at the start of the trip
- pickup\_latitude - Latitude at the start of the trip
- dropoff\_longitude - Longitude at the end of the trip
- dropoff\_latitude - Latitude at the end of the trip

Below is some information on the cleaned data set where the key field has been removed.

Field	Non-Null Count	Data Type
fare_amount	199259	Float64
pickup_datetime	199259	Int64
pickup_longitude	199259	Float64
pickup_latitude	199259	Float64
dropoff_longitude	199259	Float64
dropoff_latitude	199259	Float64
passenger_count	199259	Int64

Using the head function we can observe the first five entries of the dataset shown below to understand the dataset better.

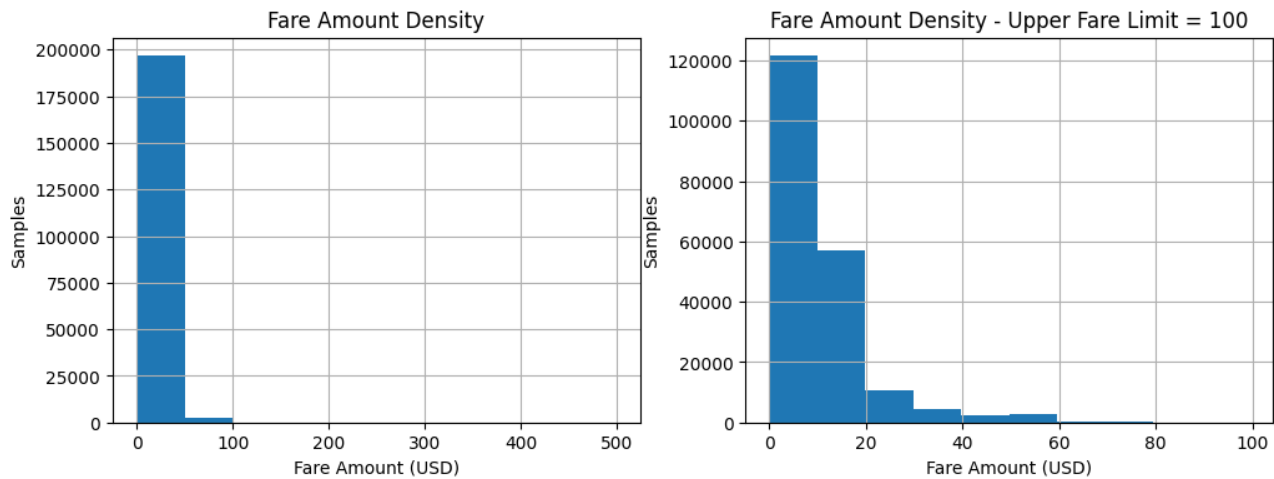
	fare_amount	pickup_datetime	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude	passenger_count
0	7.5	1431028326000000000	-73.999817	40.738354	-73.999512	40.723217	1
1	7.7	1247861096000000000	-73.994355	40.728225	-73.994710	40.750325	1
2	12.9	1251150300000000000	-74.005043	40.740770	-73.962565	40.772647	1
3	5.3	1246004541000000000	-73.976124	40.790844	-73.965316	40.803349	3
4	16.0	1409248020000000000	-73.925023	40.744085	-73.973082	40.761247	5

The describe function gives a numerical description of the data, for example 25% of the Uber fares are below six dollars usd.

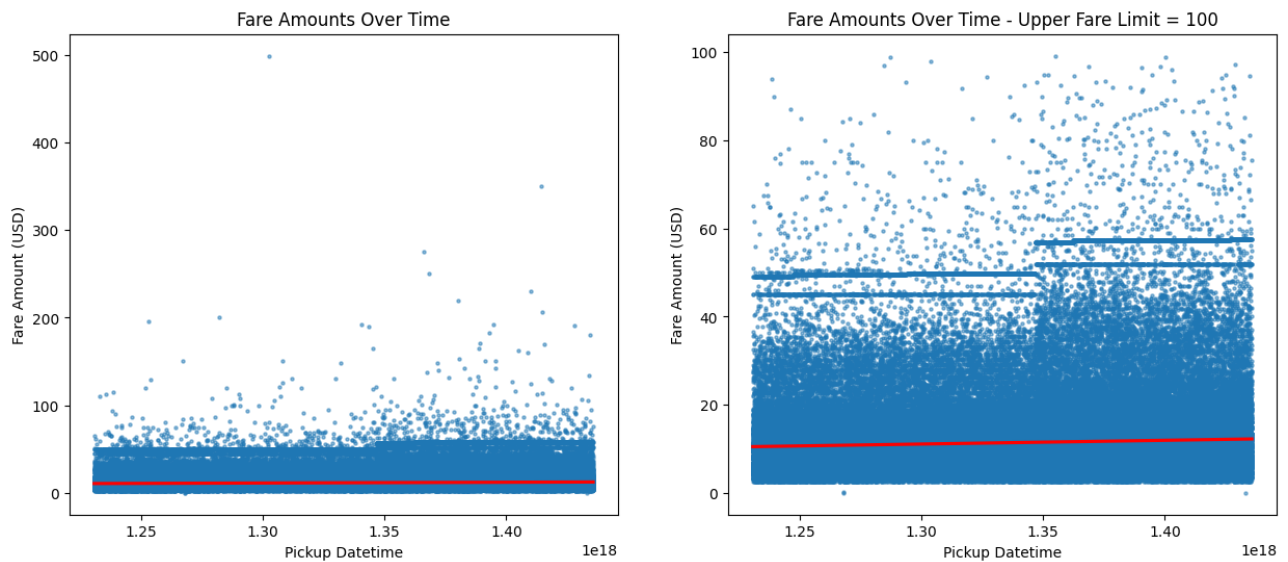
	fare_amount	pickup_datetime	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude	passenger_count
count	199259.000000	1.992590e+05	199259.000000	199259.000000	199259.000000	199259.000000	199259.000000
mean	11.369411	1.332473e+18	-72.504195	39.919184	-72.517762	39.925768	1.690488
std	9.905917	5.830927e+16	10.442166	6.127711	10.509222	6.198131	1.384791
min	0.010000	1.230773e+18	-93.824668	-74.015515	-737.916665	-74.015750	1.000000
25%	6.000000	1.282489e+18	-73.992064	40.734792	-73.991409	40.733828	1.000000
50%	8.500000	1.332573e+18	-73.981825	40.752582	-73.980094	40.753042	1.000000
75%	12.500000	1.382305e+18	-73.967161	40.767155	-73.963663	40.767995	2.000000
max	499.000000	1.435708e+18	40.808425	48.018760	40.831932	493.533332	208.000000

## 2.2 Graph Analysis

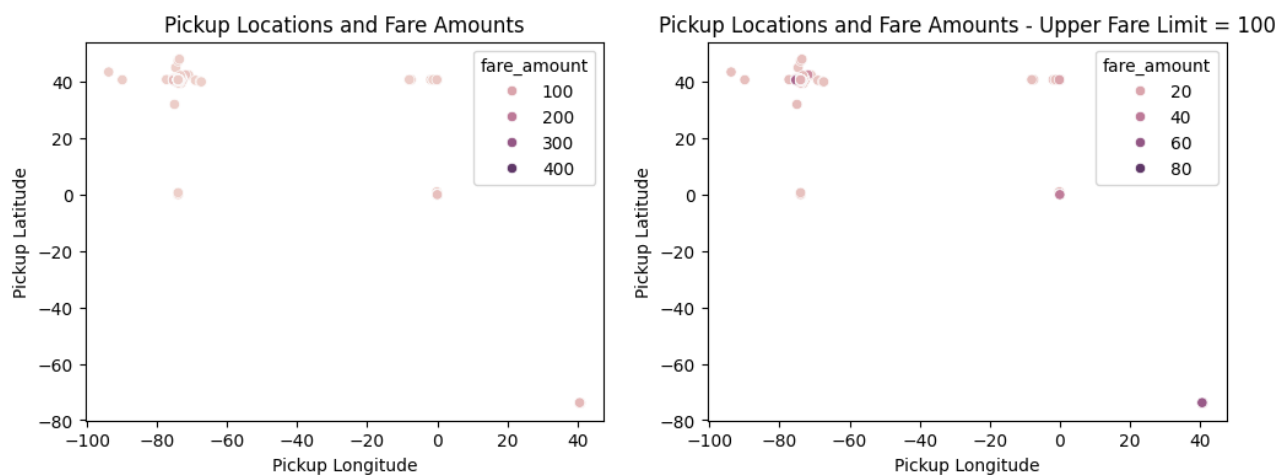
While these tables give us an understanding of the dataset structure and values included graphs will help us better understand the correlations and interesting relationships between the data points. As we can see from the resultant graph displaying all the fare amounts as a histogram, majority of the fares are distributed in the range \$0 - \$25.



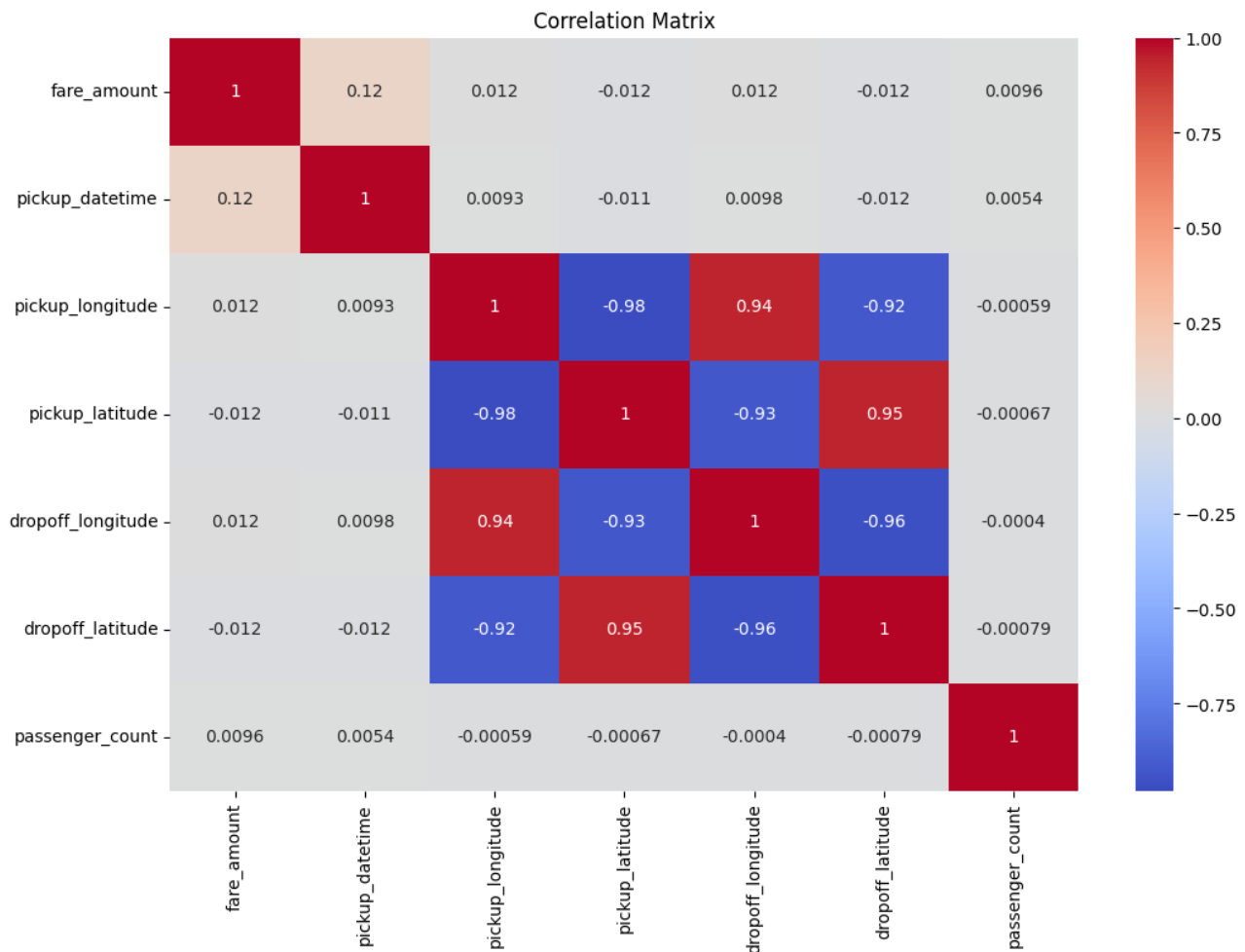
The next following graphs show the fare amount in respect to date-time with the original dataset and one restricted to fare amounts lower than 100 dollars for better visualisation. From these graphs, it's a little hard to visually see any overall correlation between date-time and fare amount. The main observation that can be drawn is the increase in fare amounts near data-time 1.35e18 which corresponds to some time in 2012 where there was likely a mark up in service costs. Other observations will likely need to be left to be seen in the correlation matrix.



The last graphs compare fare amount and pickup location with the original dataset and one restricted to fare amounts lower than 100 dollars for better visualisation. It helps to identify any geographical patterns or areas with higher fares. There doesn't seem to be a clear correlation between geographical location and price. It should be noted however that in the second graph, there is a significantly higher cost fare in the bottom right however those coordinates are in Antarctica so it's legitimacy is questionable.



Finally from the correlation matrix we can conclude that the fare seems to be reliant more on pickup date and time as compared to other variables, with pickup longitude as the runner-up.



### 3 Data Preparation

#### 3.1 Data Cleaning

To make the data suitable for preparation, several data points were removed. Firstly all NA values were removed, then the key field was removed for all entries as it is an unrelated feature in predicting the price of the ride. Additionally all entries with negative fares and zero passenger counts were removed as these signified data entry errors or other discrepancies. Then the data set was restricted to include only valid coordinates as many coordinates went beyond the limit of 180 or -180 thus these points needed to be removed, this was done by limiting the longitude and latitude values. Finally all dates were converted to integer representations.

#### 3.2 Pipe lining and Preparation for Regression

Before pipe lining the target feature being the fare amount is separated. For pipe lining a modified version of the end to end ML project sample on the course website is used. First the missing numerical values are filled with the mean using a SimpleImputer then the numerical columns scaled using StandardScaler without scaling the target. Then the data set is split into 80% training set and 20% test.

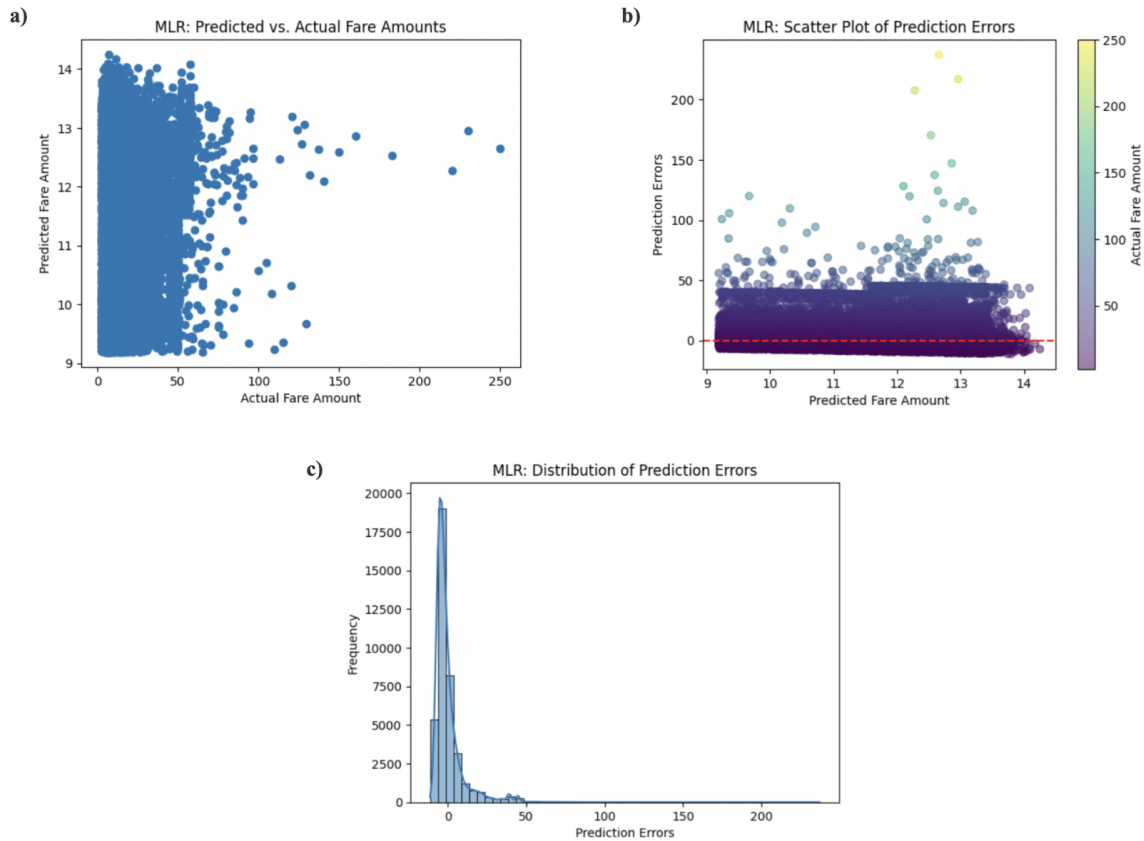
### 4 Training and Comparison

As mentioned in section 1 this task is a regression task so three different models were used which are best suited for regression tasks. The models used were multiple linear regression, ridge regression and elastic net regression. For each model the mean squared error and the mean absolute error were calculated and evaluated based on the scoring from sklearn. Below is a table comparing these results of the different regression models.

	Multiple Linear	Ridge	Elastic Net
Mean Squared Error	9.999959878142283	9.999959958084101	10.032650021093067
Mean Absolute Error	6.033518385262153	6.033518453483443	6.061790508287913
Sklearn Score	0.0158666835743817	0.015866667839637305	0.009421849019126127

## 5 Best Performing Algorithm

Based on the scores provided in the previous section, the multiple linear regression(MLR) model is the best performing algorithm with ridge regression being a close second. Below there are three graphs displaying the results from the MLR model. Graph a displays the predicted values in respect to the actual values, the extreme outliers were not accurately predicted by the model. It is unclear whether this is due to the model not being able to generalize or if these points are data errors that seem legitimate and were not cleaned out during the data cleaning step. Graph b plots predicted values versus error values and uses a colour gradient for the actual fare amount and shows similar trends to graph A. Graph C plots the distribution in errors for the model. Based on the graphs shown below there does not seem to be a clear correlation between the features and the prediction errors which is good in the sense that this means that the model isn't over or under valuing any feature but bad in the sense that this means it is difficult to associate any problematic features in terms of consistency.



## 6 Limitations

Even the best model of the three trained did not achieve the best results possible which may be due to many reasons. Firstly the dataset includes many points which are valid points able to be passed through to the model but make no sense to the overall dataset or big picture. For example the datasets are mostly centered around New York however there are points in Antarctica, Europe, Africa and other points far away from New York which can be passed to the network but are hard to determine as being legitimate. While we decided to include all data points whose coordinates lied within the bounds of the globe, there were other factors that maybe

should have been considered but were not to prevent input bias. For example, since the extreme majority of data points were within New York, should we have restricted the data to New York since for any prediction outside of New York, the accuracy will fall and if we are predicting within the bounds of the city, then restricting the training to the city should improve performance. There were some coordinates as well that while they were valid coordinates, they were in the ocean however these would prove very complicated to clean without the help of an external library so these were left as is. The last consideration was in the fare cost. The extreme majority of the fare costs didn't exceed 25\$ and looked believable up to 100\$. The fare costs that went into the hundreds, even as high as \$500 seem suspicious and if there was more evidence of these being faulty data points, they would have been removed. Since there was no firm evidence of that, they were kept and as a result, are partly to blame for the fact that the best performing model's MAE to be as high as 6 when 75% of trips are under \$12.5.



## A

### Source Code

```
1  # %% [markdown]
2  # # Uber Fares Predictions
3
4  # %% [markdown]
5  # ## Preparing The Data
6
7  # %% [markdown]
8  # ### Import Relevant Libraries
9
10 # %%
11 import sklearn
12 import matplotlib.pyplot as plt
13 import numpy as np
14 import pandas as pd
15 import seaborn as sns
16 from IPython.core.interactiveshell import InteractiveShell
17 from IPython.display import Markdown, display
18 from sklearn.compose import ColumnTransformer
19 from sklearn.impute import SimpleImputer
20 from sklearn.metrics import mean_squared_error, mean_absolute_error
21 from sklearn.model_selection import train_test_split
22 from sklearn.pipeline import make_pipeline
23 from sklearn.preprocessing import OneHotEncoder, StandardScaler
24 from sklearn.linear_model import LinearRegression, Ridge, ElasticNet
25
26 # %% [markdown]
27 # ### Configuring Libraries
28
29 # %%
30 pd.options.mode.chained_assignment = None
31
32 # Use a fixed seed to make results reproducible
33 np.random.seed(873409)
34
35 # %% [markdown]
36 # ### Loading The Dataset
37 #
38 # Note: Kaggle does not support importing directly so the .csv must be in the same directory as
39 ↪ the notebook
40
41 # %%
42 df = pd.read_csv("uber.csv")
43
44 # %% [markdown]
45 # ### Cleaning The Dataset
46
47 # %%
48 # Define coordinate bounds
49 longitude_min = -180
50 longitude_max = 180
51 latitude_min = -90
52 latitude_max = 90
53
54 # Dropping unrelated features
55 df.drop(['Unnamed: 0', 'key'], inplace=True, axis=1)
56
57 # Remove na values
58 df.dropna(inplace=True)
59
60 # Remove faulty data
```

```

60 df = df[df.fare_amount > 0]
61
62 df = df[
63     (df.pickup_longitude >= longitude_min) &
64     (df.pickup_longitude <= longitude_max) &
65     (df.pickup_latitude >= latitude_min) &
66     (df.pickup_latitude <= latitude_max)
67 ]
68
69 df = df[df.passenger_count > 0]
70
71 # Convert date to int
72 df.pickup_datetime = pd.to_numeric(pd.to_datetime(df.pickup_datetime))
73
74 # %% [markdown]
75 # ##### Creating A Slightly Stricter Dataframe For EDA Graph Pairs
76
77 # %%
78 filtered_df = df[df.fare_amount < 100]
79
80 # %% [markdown]
81 # ## Exploratory Data Analysis
82 #
83 # Note: For all graphs, a second one is generated with an upper limit of £100USD so that the
84 ↪ patterns are more easily observable at values closer to the mean.
85
86 # %% [markdown]
87 # ### Exploring The Data
88
89 # %%
90 display(df.info())
91
92 # %%
93 display(df.head())
94
95 # %%
96 display(df.describe())
97
98 # %%
99 display(Markdown("> Data Shape: {}".format(df.shape)))
100
101 # %% [markdown]
102 # ### Histogram of Fare Amounts
103 #
104 # As we can see from the resultant graph, majority of the fares are distributed in the range £0 -
105 ↪ £25.
106
107 # %%
108 plt.figure(figsize=(12,4))
109 plt.subplot(121)
110 df.fare_amount.hist()
111 plt.title('Fare Amount Density')
112 plt.xlabel('Fare Amount (USD)')
113 plt.ylabel('Samples')
114 plt.subplot(122)
115 plt.title('Fare Amount Density - Upper Fare Limit = 100')
116 plt.xlabel('Fare Amount (USD)')
117 plt.ylabel('Samples')
118 filtered_df.fare_amount.hist()
119 plt.show()
120
121 # %% [markdown]
122 # ### Scatter Plot With Line of Best Fit of Fare Amounts Over Time

```

```

121 #
122 # From these graphs, it's a little hard to visually see any overall correlation between Datetime
    ↪ and fare amount. The main observation that can be drawn is the increase in fare amounts near
    ↪ datetime 1.35e18 which corresponds to some time in 2012 where there was likely a mark up in
    ↪ service costs. Other observations will likely need to be left to be seen in the correlation
    ↪ matrix
123
124 # %%
125 plt.figure(figsize=(15, 6))
126 plt.subplot(121)
127 sns.regplot(x='pickup_datetime', y='fare_amount', data=df, scatter_kws={'s': 5, 'alpha': 0.5},
    ↪ line_kws={'color': 'red'})
128 plt.title('Fare Amounts Over Time')
129 plt.xlabel('Pickup Datetime')
130 plt.ylabel('Fare Amount (USD)')
131 plt.subplot(122)
132 sns.regplot(x='pickup_datetime', y='fare_amount', data=filtered_df, scatter_kws={'s': 5, 'alpha':
    ↪ 0.5}, line_kws={'color': 'red'})
133 plt.title('Fare Amounts Over Time - Upper Fare Limit = 100')
134 plt.xlabel('Pickup Datetime')
135 plt.ylabel('Fare Amount (USD)')
136 plt.show()
137
138 # %% [markdown]
139 # ### Scatter Plot of Fare Amounts Against Pickup Locations
140 #
141 # The following graph displays the distribution of fares with respect to location. It helps to
    ↪ identify any geographical patterns or areas with higher fares. There doesn't seem to be a
    ↪ clear correlation between geographical location and price. It should be noted however that in
    ↪ the second graph, there is a significantly higher cost fare in the bottom right however
    ↪ those coordinates are in Antarctica so it's legitimacy is questionable.
142
143 # %%
144 plt.figure(figsize=(12, 4))
145 plt.subplot(121)
146 sns.scatterplot(x='pickup_longitude', y='pickup_latitude', hue='fare_amount', data=df)
147 plt.title('Pickup Locations and Fare Amounts')
148 plt.xlabel('Pickup Longitude')
149 plt.ylabel('Pickup Latitude')
150 plt.subplot(122)
151 sns.scatterplot(x='pickup_longitude', y='pickup_latitude', hue='fare_amount', data=filtered_df)
152 plt.title('Pickup Locations and Fare Amounts - Upper Fare Limit = 100')
153 plt.xlabel('Pickup Longitude')
154 plt.ylabel('Pickup Latitude')
155 plt.show()
156
157 # %% [markdown]
158 # ### Correlation
159
160 # %%
161 # Calculate the correlation matrix
162 correlation_matrix = df.corr()
163
164 # Visualize the correlation matrix as a heatmap
165 plt.figure(figsize=(12, 8))
166 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm')
167 plt.title('Correlation Matrix')
168 plt.show()
169
170 # %% [markdown]
171 # From this matrix we can conclude that the fare seems to be reliant more on pickup date and time
    ↪ as compared to other variables, with pickup longitude as the runner-up.
172

```

```

173 # %% [markdown]
174 # ## Preparing Data For Regression
175
176 # %% [markdown]
177 # ### Separate Target Feature
178
179 # %%
180 dft = df.fare_amount
181 df.drop(['fare_amount'], inplace=True, axis=1)
182 display(Markdown("> Feature Set Shape: {}".format(df.shape)))
183 display(Markdown("> Target Set Shape: {}".format(dft.shape)))
184
185 # %% [markdown]
186 # ### Apply Pipeline
187 #
188 # Note: This code is an adapted version of the code in the end-to-end ML samples on the course
189 ↪ website
190
191 # %%
192 num_cols = df.select_dtypes(include='number').columns.to_list()
193
194 # Create pipelines for numeric and categorical columns
195 num_pipeline = make_pipeline(SimpleImputer(strategy='mean'), StandardScaler())
196
197 # Use ColumnTransformer to set the estimators and transformations
198 preprocessing = ColumnTransformer([('num', num_pipeline, num_cols)],
199                                   remainder='passthrough'
200                                   )
201
202 display(preprocessing)
203
204 df_prepared = preprocessing.fit_transform(df)
205
206 feature_names = preprocessing.get_feature_names_out()
207 df_prepared = pd.DataFrame(data=df_prepared, columns=feature_names)
208
209 display(df_prepared.shape)
210
211 # %% [markdown]
212 # ### Split the datasets
213
214 # %%
215 df_train, df_test, dft_train, dft_test = train_test_split(df_prepared, dft, test_size=0.2,
216 ↪ random_state=42)
217 display(Markdown("> %s %s %s %s" % (df_train.shape, dft_train.shape, df_test.shape,
218 ↪ dft_test.shape)))
219
220 # %% [markdown]
221 # ## Training 3 Different Regression Models
222
223 # %% [markdown]
224 # ### Multiple Linear Regression
225
226 # %%
227 mlr = LinearRegression()
228 mlr.fit(df_train, dft_train)
229
230 # %% [markdown]
231 # ### Ridge Regression
232
233 # %%
234 rr = Ridge()

```

```

233 rr.fit(df_train, dft_train)
234
235 # %% [markdown]
236 # ### Elastic Net Regression
237
238 # %%
239 enr = ElasticNet()
240 enr.fit(df_train, dft_train)
241
242 # %% [markdown]
243 # ### Comparing Performances Between The Models
244
245 # %% [markdown]
246 # ### Multiple Linear Regression
247
248 # %%
249 mlr_dft_predict = mlr.predict(df_test)
250
251 mlr_RMSE = mean_squared_error(dft_test, mlr_dft_predict, squared=False)
252 mlr_MAE = mean_absolute_error(dft_test, mlr_dft_predict)
253
254 display(Markdown("> Multiple Linear Regression RMSE: %s" % mlr_RMSE))
255 display(Markdown("> Multiple Linear Regression MAE: %s" % mlr_MAE))
256
257 # %% [markdown]
258 # ### Ridge Regression
259
260 # %%
261 rr_dft_predict = rr.predict(df_test)
262
263 rr_RMSE = mean_squared_error(dft_test, rr_dft_predict, squared=False)
264 rr_MAE = mean_absolute_error(dft_test, rr_dft_predict)
265
266 display(Markdown("> Ridge Regression RMSE: %s" % rr_RMSE))
267 display(Markdown("> Ridge Regression MAE: %s" % rr_MAE))
268
269 # %% [markdown]
270 # ### Elastic Net Regression
271
272 # %%
273 enr_dft_predict = enr.predict(df_test)
274
275 enr_RMSE = mean_squared_error(dft_test, enr_dft_predict, squared=False)
276 enr_MAE = mean_absolute_error(dft_test, enr_dft_predict)
277
278 display(Markdown("> Elastic Net Regression RMSE: %s" % enr_RMSE))
279 display(Markdown("> Elastic Net Regression MAE: %s" % enr_MAE))
280
281 # %% [markdown]
282 # ### Comparing The Performances using sklearn's scoring
283
284 # %%
285 display(Markdown("> Multiple Linear Regression Score: %s" % (mlr.score(df_test, dft_test))))
286 display(Markdown("> Ridge Regression Score: %s" % (rr.score(df_test, dft_test))))
287 display(Markdown("> Elastic Net Regression: %s" % (enr.score(df_test, dft_test))))
288
289 # %% [markdown]
290 # ### Analysing The Best Performing Model - Multiple Linear Regression
291
292 # %% [markdown]
293 # ### Calculating Prediction Errors
294
295 # %%

```

```

296 prediction_errors = dft_test - mlr_dft_predict
297
298 # %% [markdown]
299 # ### Graphing Predicted vs Actual Values
300
301 # %%
302 plt.scatter(dft_test, mlr_dft_predict)
303 plt.xlabel("Actual Fare Amount")
304 plt.ylabel("Predicted Fare Amount")
305 plt.title("MLR: Predicted vs. Actual Fare Amounts")
306 plt.show()
307
308 # %% [markdown]
309 # ### Graphing Scatter Plot of Prediction Errors
310
311 # %%
312 # Scatter Plot of Prediction Errors
313 plt.scatter(mlr_dft_predict, prediction_errors, c=dft_test, cmap='viridis', alpha=0.5)
314 plt.axhline(y=0, color='r', linestyle='--')
315 plt.colorbar(label="Actual Fare Amount")
316 plt.xlabel("Predicted Fare Amount")
317 plt.ylabel("Prediction Errors")
318 plt.title("MLR: Scatter Plot of Prediction Errors")
319 plt.show()
320
321 # %% [markdown]
322 # ### Graphing Distribution of Prediction Errors
323
324 # %%
325 sns.histplot(prediction_errors, bins=50, kde=True)
326 plt.xlabel("Prediction Errors")
327 plt.ylabel("Frequency")
328 plt.title("MLR: Distribution of Prediction Errors")
329 plt.show()
330
331 # %% [markdown]
332 # There does not seem to be a clear correlation between the features and the prediction errors
333 ↪ which is good in the sense that this means that the model isn't over or under valuing any
334 ↪ feature but bad in the sense that this means it is difficult to associate any problematic
335 ↪ features in terms of consistency.

```

## **B**

### **Links**

#### **B.1 Video**

<https://youtu.be/nywGmpOWxcU>

#### **B.2 Dataset**

<https://www.kaggle.com/datasets/yasserh/uber-fares-dataset/>

#### **B.3 GitHub**

<https://github.com/GabrielERincon/IntrotoAI-UberFares>