

Project Title

Uber Fare Prediction

TEAM MEMBERS:

Sailesh Pilla (11593815)

Yaswanth Sai Satish Sreerama (11601948)

Kishore Kumar Paila (11600316)

Project Proposal Link: [GitHub](#)

Demo Video Link: [demo](#)

Idea Description:

Our Project is about Uber, the largest taxi company in the world, serving thousands of consumers every day. I noticed some open-source data from Kaggle. Due to the service's global reach, it is necessary to precisely compute fare pricing considering several factors, including distance, weather, time, and demand.

Goals and Objectives:

Goals of this project is to predict fare accordingly in different places with respect to different variables by training the model first and then test it on test dataset.

Objectives:

1. Thoroughly understand the dataset and determine whether the data is correct or not, and if not, perform a clean-up.
2. Create a regression model over the training data that has been separated.
3. Evaluate the model and compare R² and Root mean squared error scores to find the projected model error so that you may understand how accurate your predictions were.

Motivation:

The reason I wanted to know what type of algorithm or model they were using to see what all elements play a role in forecasting price and see which dimension plays a significant function in accurately predicting the fare is that we utilize uber a lot in our daily lives to get to our destinations.

Significance:

This is accomplished by predicting fare quickly whenever a user enters his pickup and destination; it calculates it automatically considering all the factors using a machine learning algorithm. It attracts users when he can see the amount that he needs to pay before booking; it will give user to allocate or plan their budget accordingly and they can choose the type of vehicle accordingly. Here, utilizing ML models increases the likelihood that the prediction will be accurate. If it is incorrect, we can compute the error and improve our model's training by using more accurate data or by omitting irrelevant data. A better prediction can be made with the right training data.

Literature Survey:

We are attempting to comprehend how Uber operates, how the uber fare is computed, and what additional criteria are taken into account to estimate the fare using ML models. Given that there are numerous variables that affect how fare is predicted, we will attempt to use a multi-linear regression.

1. Data Exploration
2. EDA (Explanatory Data Analysis)
3. Data Pre-Processing
4. Feature Selection
5. Splitting of Data
6. Applying Models and Performing Prediction

Features:

- key - a unique identifier for each trip
- fare_amount - the cost of each trip in USD
- pickup_datetime - date and time when the meter was engaged
- passenger_count - the number of passengers in the vehicle
- pickup_longitude - the longitude where the meter was engaged
- pickup_latitude - the latitude where the meter was engaged
- dropoff_longitude - the longitude where the meter was disengaged
- dropoff_latitude - the latitude where the meter was disengaged

Expected Outcome:

Using the data's longitude, latitude, and passenger count, predict the fare amount. To improve our ability to predict the fare amount variable from the given dataset, we can extract new features from the dataset.

Increment 1

Related Work (Background):

Attained data set contains over two lakh travels and spans the years 2009 to 2015.

Dataset: Uber.csv

Detail Design of features:

Given Features:

- key - a unique identifier for each trip
- fare_amount - the cost of each trip in usd
- pickup_datetime - date and time when the meter was engaged
- passenger_count - the number of passengers in the vehicle
- pickup_longitude - the longitude where the meter was engaged
- pickup_latitude - the latitude where the meter was engaged
- dropoff_longitude - the longitude where the meter was disengaged
- dropoff_latitude - the latitude where the meter was disengaged

Created Features:

- Distance (including pickup latitude and longitude. Latitude and longitude of drop off)

Analysis:

- After inspecting the dataset, I discovered that it is in csv format with comma separated values, which I must read in. As I'm doing it in Google Colab, I'm uploading my dataset using Google Colab library and using browse function to upload dataset and save it in a dataframe. In order to calculate the fare from an Uber dataset, where each feature denotes a distinct representation, I intended to use a multiple linear regression model.
- The key represents the transaction number for that particular trip, which will be distinct. If we want to know how a particular trip went, we can check it with the key to find out how long it took, when the pickup time was, how many passengers traveled, what kind of car was used more frequently, and whether long or short distances were traveled.
- Before removing null data or filling it with mode values, all necessary analysis must be completed, starting with checking the data's head and null values.
- Once the necessary preprocessing processes have been completed, we can build new features, such as distance calculations utilizing pickup and dropoff latitude and longitude with all four features, and then drop the other four features to reduce the number of columns.
- In order to see the surge in fare depending on time by pickup time column, we can create detailed analysis by creating features like type of car using passenger count and another feature like time of travel and creating bin's in the feature with morning, noon, evening, and night. Then we can drop this column and use the newly created feature as it will give more insights for the prediction.
- Once all necessary fields have been made available, dummification or one-hot encoding will be used to turn continuous data into categorical variables. As soon as the data is prepared, we divide the X and Y variables into independent and dependent ones in a ratio of 70:30 for the train and test runs.
- Once the data has been divided, the model will need to be fitted using a multiple linear regression model once the data has been trained. Now, in order to determine residual errors, we predict the data for the train fare amount for which we already have a number. We then compare it to the projected value and note any differences.

- Then use the same model on a test dataset to examine how it will predict, how accurate it will be with observed and predicted values, and how the train and test data are fitted. At this point, you should calculate the MSE, RMSE, and R Square metrics to determine how well your model predicts. Additionally, you may utilize any fresh data to predict, estimate its accuracy, and visually assess the performance of your model by drawing a best fit line.

Implementation: Work Completed:

- **Description:** After obtaining a dataset, we completed the necessary preprocessing procedures, followed by the required data analysis, visualization, and train and test split data.

Responsibility:

- **Sailesh** conducted the necessary research to identify the best dataset from which to make the prediction. He then uploaded the dataset and performed a few preprocessing operations, such as computing null values, reading the dataset, and cleaning the dataset.
- **Kishore** completed the feature engineering portion by adding new features, deleting unnecessary columns, performing the required categorical encoding, dummifying data using a single hot encoding operation or label encoder, and completing the required visualizations to gain insights into the removal of unnecessary features.
- **Yaswanth** completed the necessary visualizations, the train-test split, the multiple linear regression model, and determined the model's accuracy to apply the model to any additional dataset.

Data Overview and Documentation:

2. Checking head of the datasets

1uber_data.head()

	Unnamed: 0	key	fare_amount	pickup_datetime	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude	passenger_count
0	24238194	2015-05-07 19:52:06.0000003	7.5	2015-05-07 19:52:06 UTC	-73.999817	40.738354	-73.999512	40.723217	1
1	27835199	2009-07-17 20:04:56.0000002	7.7	2009-07-17 20:04:56 UTC	-73.994355	40.728225	-73.994710	40.750325	1
2	44984355	2009-08-24 21:45:00.00000061	12.9	2009-08-24 21:45:00 UTC	-74.005043	40.740770	-73.962565	40.772647	1
3	25894730	2009-06-26 08:22:21.0000001	5.3	2009-06-26 08:22:21 UTC	-73.976124	40.790844	-73.965316	40.803349	3
4	17610152	2014-08-28 17:47:00.000000188	16.0	2014-08-28 17:47:00 UTC	-73.925023	40.744085	-73.973082	40.761247	5

Info of the data:

: 1uber_data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 200000 entries, 0 to 199999
Data columns (total 9 columns):
Column Non-Null Count Dtype
--- ---
0 Unnamed: 0 200000 non-null int64
1 key 200000 non-null object
2 fare_amount 200000 non-null float64
3 pickup_datetime 200000 non-null object
4 pickup_longitude 200000 non-null float64
5 pickup_latitude 200000 non-null float64
6 dropoff_longitude 199999 non-null float64
7 dropoff_latitude 199999 non-null float64
8 passenger_count 200000 non-null int64
dtypes: float64(5), int64(2), object(2)
memory usage: 13.7+ MB

Describing Data:

1uber_data.describe()

	Unnamed: 0	fare_amount	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude	passenger_count
count	2.000000e+05	200000.000000	200000.000000	200000.000000	199999.000000	199999.000000	200000.000000
mean	2.771250e+07	11.359955	-72.527638	39.935885	-72.525292	39.923890	1.684535
std	1.601382e+07	9.901776	11.437787	7.720539	13.117408	6.794829	1.385997
min	1.000000e+00	-52.000000	-1340.648410	-74.015515	-3356.666300	-881.985513	0.000000
25%	1.382535e+07	6.000000	-73.992065	40.734796	-73.991407	40.733823	1.000000
50%	2.774550e+07	8.500000	-73.981823	40.752592	-73.980093	40.753042	1.000000
75%	4.155530e+07	12.500000	-73.967154	40.767158	-73.963658	40.768001	2.000000
max	5.542357e+07	499.000000	57.418457	1644.421482	1153.572603	872.697628	208.000000

Number of columns and rows in the data:

1uber_data.shape

(200000, 9)

Haversine function to return distance by taking pickup and dropoff locations:

```
1 def haversine(lon_1, lon_2, lat_1, lat_2):
2
3     lon_1, lon_2, lat_1, lat_2 = map(np.radians, [lon_1, lon_2, lat_1, lat_2]) #Degrees to Radians
4
5
6     diff_lon = lon_2 - lon_1
7     diff_lat = lat_2 - lat_1
8
9
10    km = 2 * 6371 * np.arcsin(np.sqrt(np.sin(diff_lat/2.0)**2 +
11                                     np.cos(lat_1) * np.cos(lat_2) * np.sin(diff_lon/2.0)**2))
12
13    return km
```

```
1 uber_data['distance'] = haversine(uber_data['pickup_longitude'],uber_data['dropoff_longitude'],
2                                   uber_data['pickup_latitude'],uber_data['dropoff_latitude'])
```

Retrieving year,month weekday and hour by using to_date from pandas:

```
1 uber_data.pickup_datetime=pd.to_datetime(uber_data.pickup_datetime)
```

```
1 uber_data['year'] = uber_data.pickup_datetime.dt.year
2 uber_data['month'] = uber_data.pickup_datetime.dt.month
3 uber_data['weekday'] = uber_data.pickup_datetime.dt.weekday
4 uber_data['hour'] = uber_data.pickup_datetime.dt.hour
```

Dividing the to_date into segments and monthly quarter for better understanding of data:

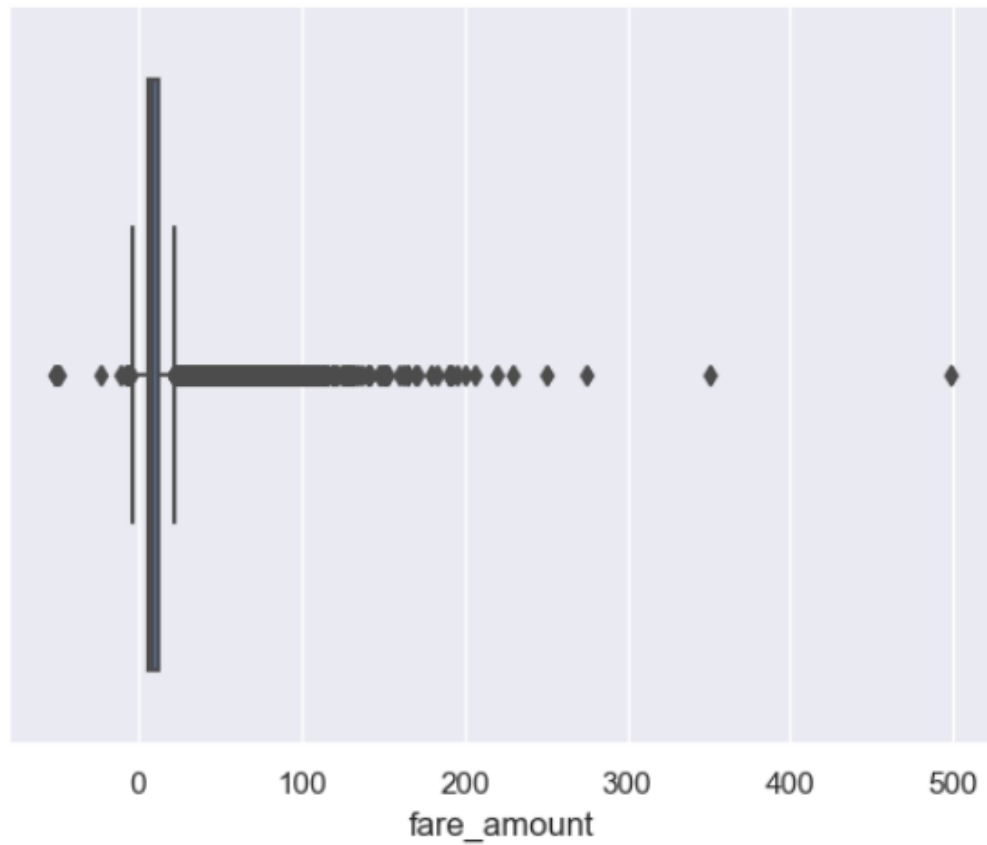
```
1 uber_data['Monthly_Quarter'] = uber_data.month.map({1:'Q1',2:'Q1',3:'Q1',4:'Q2',5:'Q2',6:'Q2',7:'Q3',
2                                                    8:'Q3',9:'Q3',10:'Q4',11:'Q4',12:'Q4'})
3 uber_data['Hourly_Segments'] = uber_data.hour.map({0:'H1',1:'H1',2:'H1',3:'H1',4:'H2',5:'H2',6:'H2',7:'H2',8:'H3',
4                                                    9:'H3',10:'H3',11:'H3',12:'H4',13:'H4',14:'H4',15:'H4',16:'H5',
5                                                    17:'H5',18:'H5',19:'H5',20:'H6',21:'H6',22:'H6',23:'H6'})
```

Visualisations representing the data:

Box Plot:

```
1 sns.boxplot(x=uber_data['fare_amount'])
```

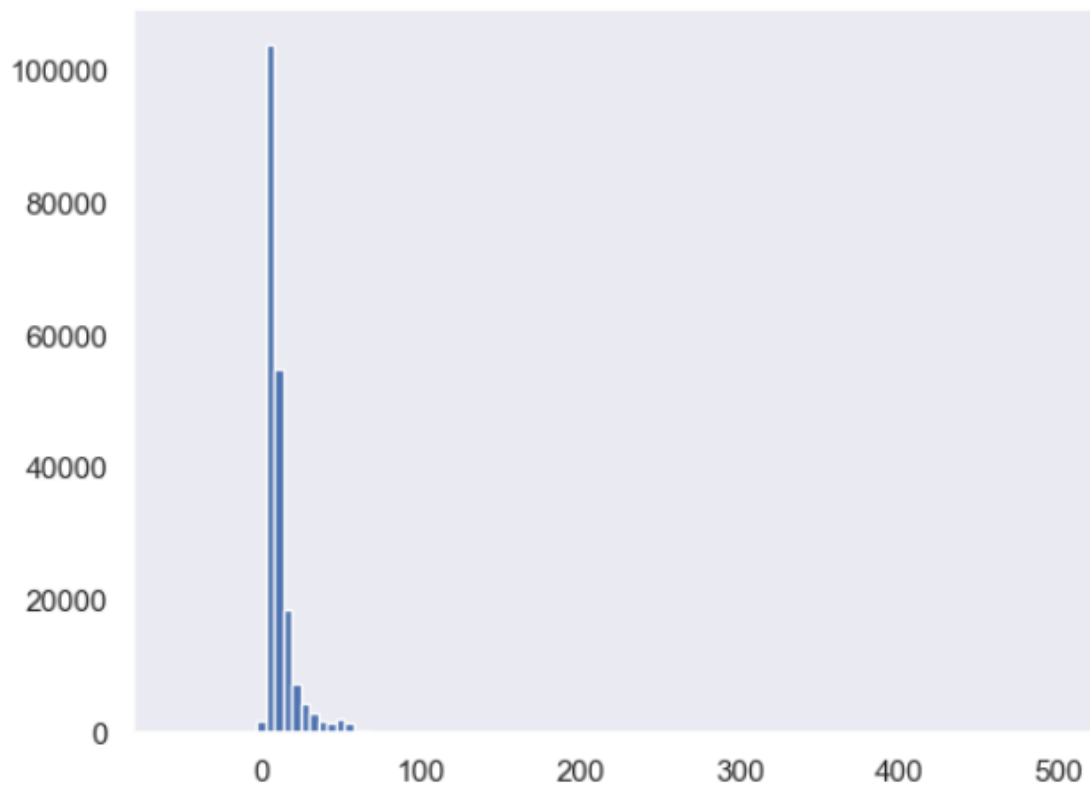
<AxesSubplot: xlabel='fare_amount'>



Histogram:

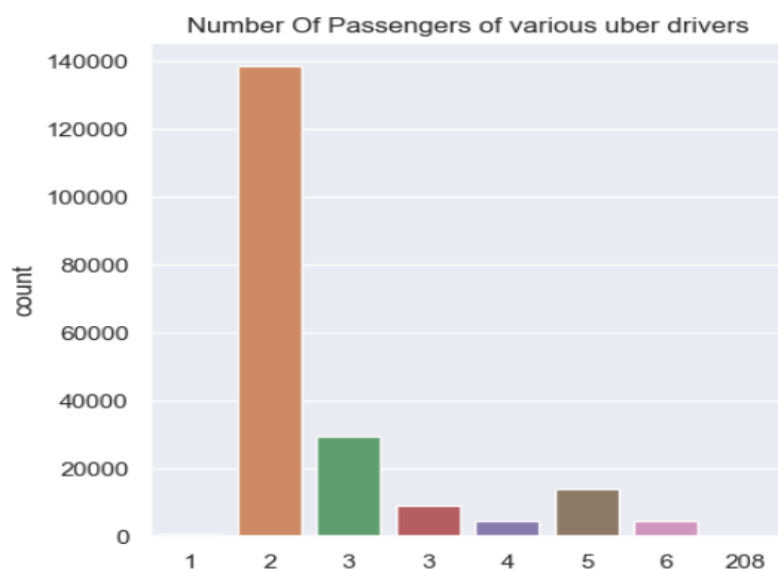
```
1 uber_data['fare_amount'].hist(bins=100,grid=False)
```

<AxesSubplot: >



Count Plot:

```
1 fig, ax = plt.subplots(figsize = (5, 5))
2 sns.countplot(x = uber_data.passenger_count.values, data=uber_data)
3 labels = [item.get_text() for item in ax.get_xticklabels()]
4 labels[0] = '1'
5 labels[1] = '2'
6 labels[2] = '3'
7 ax.set_xticklabels(labels)
8 ax.set_title("Number Of Passengers of various uber drivers")
9 plt.show()
```



Heat Map:

```
1 #heat map
2 plt.figure(figsize=(10,5))
3 hm = sns.heatmap(uber_data.corr(),annot=True)
4 plt.show()
```

C:\Users\yaswa\AppData\Local\Temp\ipykernel_10192\2389066797.py:3: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.

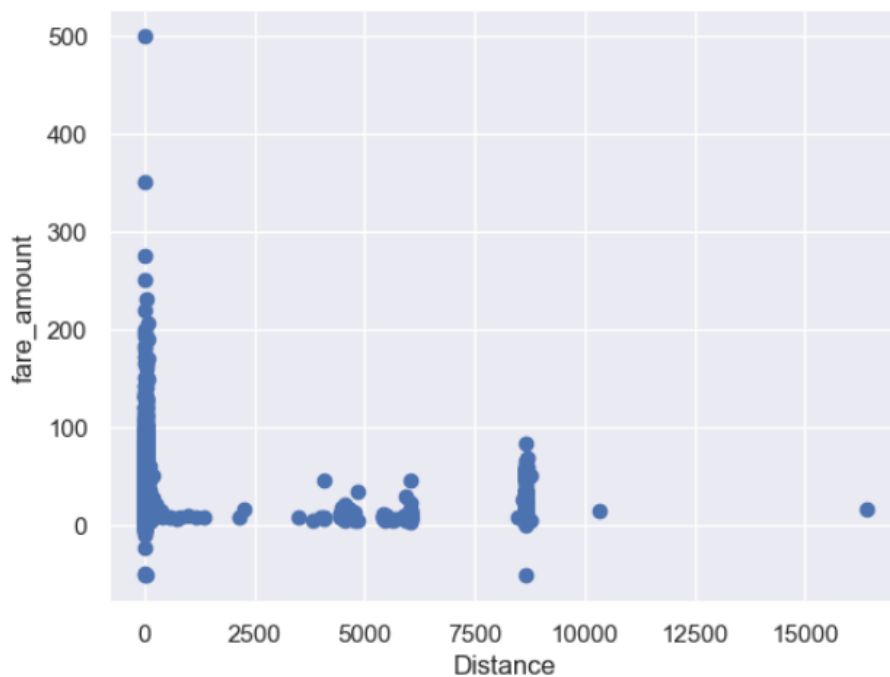
```
hm = sns.heatmap(uber_data.corr(),annot=True)
```



Scatter Plot:

```
1 plt.scatter(uber_data['distance'], uber_data['fare_amount'])
2 plt.xlabel("Distance")
3 plt.ylabel("fare_amount")
```

```
Text(0, 0.5, 'fare_amount')
```



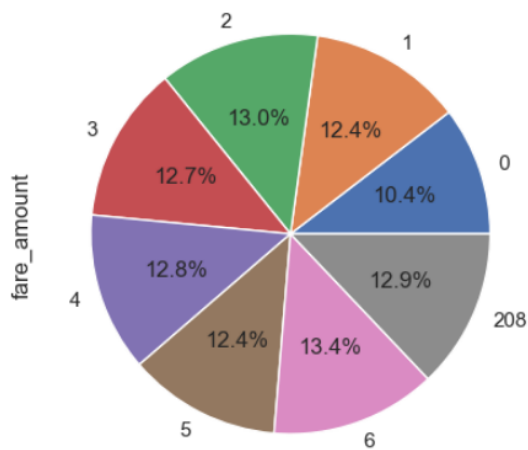
Pie Chart:

```
1 grouped=uber_data.groupby('passenger_count')
2 print(grouped['fare_amount'].mean())
```

```
passenger_count
0      9.439266
1     11.254158
2     11.784452
3     11.486731
4     11.642472
5     11.199698
6     12.158537
208    11.700000
Name: fare_amount, dtype: float64
```

```
1 grouped['fare_amount'].mean().plot(kind='pie',y='fare_amount',autopct='%1.1f%%')
```

<AxesSubplot: ylabel='fare_amount'>



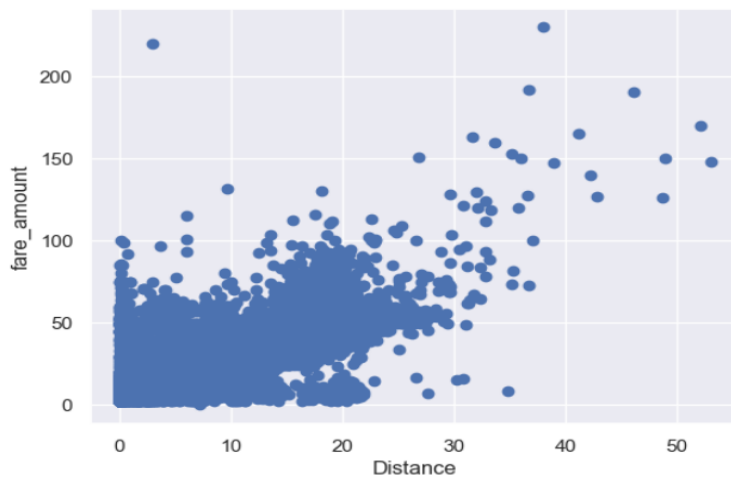
Scatter Plot after removing irrelevant rows :

```
1 uber_data.drop(uber_data[uber_data['distance'] > 60].index, inplace = True)
2 uber_data.drop(uber_data[uber_data['distance'] == 0].index, inplace = True)
3 uber_data.drop(uber_data[uber_data['fare_amount'] == 0].index, inplace = True)
4 uber_data.drop(uber_data[uber_data['fare_amount'] < 0].index, inplace = True)
```

```
1 uber_data.drop(uber_data[(uber_data['fare_amount']>100) & (uber_data['distance']<1)].index, inplace = True )
2 uber_data.drop(uber_data[(uber_data['fare_amount']<100) & (uber_data['distance']>100)].index, inplace = True )
```

```
1 plt.scatter(uber_data['distance'], uber_data['fare_amount'])
2 plt.xlabel("Distance")
3 plt.ylabel("fare_amount")
```

Text(0, 0.5, 'fare_amount')



Splitting the dataset into train and test sets:

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=0)
```

Implementing linear regression:

```
1 from sklearn.linear_model import LinearRegression
2 l_reg = LinearRegression()
3 l_reg.fit(X_train, y_train)
4
5 print("Training set score: {:.2f}".format(l_reg.score(X_train, y_train)))
6 print("Test set score: {:.7f}".format(l_reg.score(X_test, y_test)))
```

Training set score: 0.80
Test set score: 0.8014349

Predicted and actual values after fitting the training data into linear regression:

```
1 y_pred = l_reg.predict(X_test)
2 df = {'Actual': y_test, 'Predicted': y_pred}
3
4 from tabulate import tabulate
5 print(tabulate(df, headers = 'keys', tablefmt = 'psql'))
```

	Actual	Predicted
1	12.5	16.2194
2	6.5	7.89846
3	7.3	7.39028
4	19.7	24.4064
5	4.1	4.51103
6	15	15.3584
7	13.7	17.8741
8	6.1	6.07715
9	5.5	5.28881
10	31.3	21.9703
11	4.9	7.04995
12	8.1	4.29591
13	6	6.69743
14	15	7.17872
15	8.5	11.132
16	10	10.3526
17	8.5	7.38459
18	16	11.5931
19	16.9	17.3298

Checking the accuracy of the data with predicting:

```
1 from sklearn import metrics
2 print('Mean Absolute Error:', metrics.mean_absolute_error(y_test, y_pred))
3 #print('Mean Absolute % Error:', metrics.mean_absolute_percentage_error(y_test, y_pred))
4 print('Mean Squared Error:', metrics.mean_squared_error(y_test, y_pred))
5 print('Root Mean Squared Error:', np.sqrt(metrics.mean_squared_error(y_test, y_pred)))
```

Mean Absolute Error: 2.2843774732419555
Mean Squared Error: 17.95984340408041
Root Mean Squared Error: 4.237905544497235

References:

[1][Fares Dataset](#)

[2][Fare Amount Prediction](#)

[3][Comparing Fares](#)

