

#Predict the price of the Uber ride from a given pickup point to the agreed drop-off location. Perform following tasks.

1. Pre-process the dataset.
2. Identify outliers.
3. Check the correlation.
4. Implement linear regression and random forest regression models.
5. Evaluate the models and compare their respective scores like R2, RMSE, etc. Dataset link: <https://www.kaggle.com/datasets/yassenli/uber-fares-dataset>

```
In [1]: #Importing the required libraries
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt

In [2]: #Importing the dataset
df = pd.read_csv("uber.csv")
```

1. Pre-process the dataset.

```
In [3]: df.head()
```

```
Out[3]:
```

	pickup_datetime	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude	passenger_count
0	2015-06-07 19:52:06-00:00	-73.999817	40.738354	-73.999512	40.723217	1
1	2009-07-17 20:04:56-00:00	-73.994355	40.728225	-73.994710	40.750325	1
2	2009-08-24 21:45:00-00:00	-74.005043	40.740770	-73.962565	40.772647	1
3	2009-06-26 08:22:21-00:00	-73.976124	40.790844	-73.965316	40.803349	3
4	2014-08-28 17:47:00-00:00	-73.925023	40.744085	-73.973082	40.761247	5

```
In [4]: df.info() #To get the required information of the dataset

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 280000 entries, 0 to 199999
Data columns (total 9 columns):
Unnamed: 0      280000 non-null int64
key             280000 non-null object
fare_amount     280000 non-null float64
pickup_datetime 280000 non-null object
pickup_longitude 280000 non-null float64
pickup_latitude 280000 non-null float64
dropoff_longitude 280000 non-null float64
dropoff_latitude 280000 non-null float64
passenger_count 280000 non-null int64
dtypes: float64(5), int64(2), object(2)
memory usage: 13.7+ MB

In [5]: df.columns #To get number of columns in the dataset

Out[5]: Index(['Unnamed: 0', 'key', 'fare_amount', 'pickup_datetime',
        'pickup_longitude', 'pickup_latitude', 'dropoff_longitude',
        'dropoff_latitude', 'passenger_count'],
        dtype='object')

In [6]: df = df.drop(['Unnamed: 0', 'key'], axis=1) #To drop unnamed column as it isn't required

In [7]: df.head()
```

```
Out[7]:
```

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

```
In [8]: df.shape #To get the total (Rows,Columns)

Out[8]: (280000, 7)
```

```
In [9]: df.dtypes #To get the type of each column

Out[9]: fare_amount      float64
        pickup_datetime  object
        pickup_longitude  float64
        pickup_latitude  float64
        dropoff_longitude float64
        dropoff_latitude float64
        passenger_count   int64
        dtype: object
```

Column pickup_datetime is in wrong format (Object). Convert it to DateTime Format

```
In [10]: df.pickup_datetime = pd.to_datetime(df.pickup_datetime)

In [10]: df.dtypes

Out[10]: fare_amount      float64
         pickup_datetime  datetime64[ns, UTC]
         pickup_longitude  float64
         pickup_latitude  float64
         dropoff_longitude float64
         dropoff_latitude float64
         passenger_count   int64
         dtype: object

In [11]: df.isnull().sum()

Out[11]: fare_amount      0
         pickup_datetime  0
         pickup_longitude  0
         pickup_latitude  0
         dropoff_longitude  0
         dropoff_latitude  0
         passenger_count   0
         dtype: int64

In [12]: df["dropoff_latitude"].fillna(value=df["dropoff_latitude"].mean(),inplace = True)
df["dropoff_longitude"].fillna(value=df["dropoff_longitude"].median(),inplace = True)

In [12]: df.isnull().sum()

Out[12]: fare_amount      0
         pickup_datetime  0
         pickup_longitude  0
         pickup_latitude  0
         dropoff_longitude  0
         dropoff_latitude  0
         passenger_count   0
         dtype: int64
```

To segregate each time of date and time

```
In [13]: df.assign(hour = df.pickup_datetime.dt.hour,
                 day = df.pickup_datetime.dt.day,
                 month = df.pickup_datetime.dt.month,
                 year = df.pickup_datetime.dt.year,
                 dayofweek = df.pickup_datetime.dt.dayofweek)

In [13]: df.head()
```

```
Out[13]:
```

	fare_amount	pickup_datetime	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude	passenger_count	hour	day	month	year	dayofweek
0	7.5	2015-06-07 19:52:06-00:00	-73.999817	40.738354	-73.999512	40.723217	1	19	7	5	2015	3
1	7.7	2009-07-17 20:04:56-00:00	-73.994355	40.728225	-73.994710	40.750325	1	20	17	7	2009	4
2	12.9	2009-08-24 21:45:00-00:00	-74.005043	40.740770	-73.962565	40.772647	1	21	24	8	2009	0
3	5.3	2009-06-26 08:22:21-00:00	-73.976124	40.790844	-73.965316	40.803349	3	8	26	6	2009	4
4	16.0	2014-08-28 17:47:00-00:00	-73.925023	40.744085	-73.973082	40.761247	5	17	28	8	2014	3

Here we are going to use Haversine formula to calculate the distance between two points and journey, using the longitude and latitude values.

```
Haversine formula hav(θ) = sin²(θ/2).

In [14]: #From math import *
# function to calculate the travel distance from the longitudes and latitudes
def distance_transform(longitude1, latitude1, longitude2, latitude2):
    travel_dist = []

    for pos in range(len(longitude1)):
        long1,lat1, long2,lat2 = map(radians,[longitude1[pos],latitude1[pos],longitude2[pos],latitude2[pos]])
        dist_long = long2 - long1
        dist_lat = lat2 - lat1
        a = sin(dist_lat/2)**2 + cos(lat1) * cos(lat2) * sin(dist_long/2)**2
        c = 2 * asin(sqrt(a))*6371
        travel_dist.append(c)

    return travel_dist

In [14]: df["dist_travel_km"] = distance_transform(df["pickup_longitude"].to_numpy(),
                                                df["pickup_latitude"].to_numpy(),
                                                df["dropoff_longitude"].to_numpy(),
                                                df["dropoff_latitude"].to_numpy())

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

```
Out[14]:
```

	fare_amount	pickup_datetime	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude	passenger_count	hour	day	month	year	dayofweek	dist_travel_km
0	7.5	2015-06-07 19:52:06-00:00	-73.999817	40.738354	-73.999512	40.723217	1	19	7	5	2015	3	1.683233
1	7.7	2009-07-17 20:04:56-00:00	-73.994355	40.728225	-73.994710	40.750325	1	20	17	7	2009	4	2.457990
2	12.9	2009-08-24 21:45:00-00:00	-74.005043	40.740770	-73.962565	40.772647	1	21	24	8	2009	0	5.038377
3	5.3	2009-06-26 08:22:21-00:00	-73.976124	40.790844	-73.965316	40.803349	3	8	26	6	2009	4	1.661883
4	16.0	2014-08-28 17:47:00-00:00	-73.925023	40.744085	-73.973082	40.761247	5	17	28	8	2014	3	4.475450

```
In [15]: # drop the column 'pickup_datetime' using drop()
# 'axis = 1' drops the specified column
df = df.drop('pickup_datetime',axis=1)

In [15]: df.head()
```

```
Out[15]:
```

	fare_amount	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude	passenger_count	hour	day	month	year	dayofweek	dist_travel_km
0	7.5	-73.999817	40.738354	-73.999512	40.723217	1	19	7	5	2015	3	1.683233
1	7.7	-73.994355	40.728225	-73.994710	40.750325	1	20	17	7	2009	4	2.457990
2	12.9	-74.005043	40.740770	-73.962565	40.772647	1	21	24	8	2009	0	5.038377
3	5.3	-73.976124	40.790844	-73.965316	40.803349	3	8	26	6	2009	4	1.661883
4	16.0	-73.925023	40.744085	-73.973082	40.761247	5	17	28	8	2014	3	4.475450

Checking outliers and filling them

```
In [16]: df.plot(kind = "box",subplots = True,layout = (7,2),figsize=(15,20)) #Boxplot to check the outliers

Out[16]:
```

```
In [17]: #Using the InterQuartile Range to fill the values
def remove_outlier(df, col):
    Q1 = df[col].quantile(0.25)
    Q3 = df[col].quantile(0.75)
    IQR = Q3 - Q1
    lower_whisker = Q1-1.5*IQR
    upper_whisker = Q3+1.5*IQR
    df[col] = np.clip(df[col], lower_whisker, upper_whisker)
    return df

def treat_outliers_all(df, col_list):
    for c in col_list:
        df1 = remove_outlier(df, c)
    return df1

In [17]: df = treat_outliers_all(df, df.columns[0 : 11])

In [17]: df.plot(kind = "box",subplots = True,layout = (7,2),figsize=(15,20)) #Boxplot shows that dataset is free from outliers

Out[17]:
```

```
In [18]: #Uber doesn't travel over 130 kms so minimize the distance
df = df.loc[(df.dist_travel_km >= 5)] # (df.dist_travel_km <= 130)]
print("Remaining observations in the dataset:", df.shape)

Remaining observations in the dataset: (280000, 12)

In [18]: #Finding incorrect latitude (less than or greater than 90) and longitude (greater than or less than 180)
incorrect_coordinates = df.loc[(df.pickup_latitude > 90) | (df.pickup_latitude < -90) |
                               (df.dropoff_latitude > 90) | (df.dropoff_latitude < -90) |
                               (df.pickup_longitude > 180) | (df.pickup_longitude < -180) |
                               (df.dropoff_longitude > 180) | (df.dropoff_longitude < -180)]

In [18]: df.drop(incorrect_coordinates,inplace = True, errors = 'ignore')

In [18]: df.head()
```

```
Out[18]:
```

	fare_amount	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude	passenger_count	hour	day	month	year	dayofweek	dist_travel_km
0	7.5	-73.999817	40.738354	-73.999512	40.723217	1	19	7	5	2015	3	1.683233
1	7.7	-73.994355	40.728225	-73.994710	40.750325	1	20	17	7	2009	4	2.457990
2	12.9	-74.005043	40.740770	-73.962565	40.772647	1	21	24	8	2009	0	5.038377
3	5.3	-73.976124	40.790844	-73.965316	40.803349	3	8	26	6	2009	4	1.661883
4	16.0	-73.925023	40.744085	-73.973082	40.761247	5	17	28	8	2014	3	4.475450

```
In [19]: df.isnull().sum()

Out[19]: fare_amount      0
         pickup_longitude  0
         pickup_latitude  0
         dropoff_longitude  0
         dropoff_latitude  0
         passenger_count   0
         hour              0
         day              0
         month            0
         year             0
         dayofweek        0
         dist_travel_km    0
         dtype: int64

In [20]: sns.heatmap(df.isnull()) #Free for null values

Out[20]:
```

```
In [21]: df.corr() #Function to find the correlation

Out[21]:
```

```
corr
```

	fare_amount	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude	passenger_count	hour	day	month	year	dayofweek	dist_travel_km
fare_amount	1.000000	0.154969	-0.110842	0.218675	-0.129888	0.035779	0.023623	0.004563	0.000017	0.341277	0.013562	0.844374
pickup_longitude	0.154969	1.000000	0.259407	0.425619	0.073290	-0.013213	0.011579	-0.003204	0.001169	0.010198	0.034562	0.089894
pickup_latitude	-0.110842	0.259407	1.000000	0.484889	0.515714	-0.012889	0.003681	-0.001553	0.001562	-0.014243	-0.042310	-0.046812
dropoff_longitude	0.218675	0.425619	0.484889	1.000000	0.245667	-0.000308	0.019783	-0.003478	-0.002391	0.011346	-0.003336	0.186531
dropoff_latitude	-0.129888	0.073290	0.515714	0.245667	1.000000	-0.000308	0.019783	-0.003478	-0.001193	-0.009603	-0.033119	-0.038900
passenger_count	0.035779	-0.013213	-0.012889	-0.000308	-0.000308	1.000000	0.002074	0.002712	0.010591	-0.009749	0.048950	0.009709
hour	-0.023623	0.011579	0.003681	-0.000308	-0.000308	0.002074	1.000000	0.004677	-0.003926	0.002156	-0.089647	-0.038306
day	0.004563	-0.003204	-0.001553	-0.000308	-0.000308	0.010591	0.004677	1.000000	-0.017360	-0.021210	0.009617	0.003002
month	0.000017	0.001169	0.001562	0.000407	-0.001193	-0.009603	-0.003926	-0.017360	1.000000	-0.118859	-0.008786	0.011628
year	0.341277	0.010198	-0.014243	0.011346	-0.009603	-0.009749	0.002156	-0.012170	-0.118859	1.000000	0.006113	0.024278
dayofweek	0.013562	-0.024562	-0.042310	-0.003336	-0.033119	0.048950	-0.089647	0.005617	-0.008786	0.006113	1.000000	0.027053
dist_travel_km	0.844374	0.089894	-0.046812	0.186531	-0.038900	0.009709	-0.038306	0.003002	0.011628	0.024278	0.027053	1.000000

```
In [22]: fig,axis = plt.subplots(figsize = (10,6))
sns.heatmap(df.corr(),annot = True) #Correlation Heatmap (Light values means highly correlated)

Out[22]:
```

```
In [23]: df.corr()

Out[23]:
```

```
corr
```

	fare_amount	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude	passenger_count	hour	day	month	year	dayofweek	dist_travel_km
fare_amount	1	0.15	0.11	0.22	0.13	0.016	0.024	0.005	0.011	0.34	0.014	0.84
pickup_longitude	0.15	1	0.26	0.43	0.073	-0.013	0.012	-0.003	0.012	0.01	0.021	0.09
pickup_latitude	0.11	0.26	1	0.49	0.52	-0.013	0.013	0.016	0.016	0.014	0.042	0.048
dropoff_longitude	0.22	0.43	0.49	1	0.25	0.003	0.047	0.004	0.004	0.011	0.013	0.18
dropoff_latitude	0.13	0.073	0.52	0.25	1	0.003	0.02	0.005	0.016	0.006	0.032	0.039
passenger_count	0.016	-0.013	-0.013	-0.003	-0.003	1	0.02	0.027	0.1	0.007	0.049	0.008
hour	0.024	0.012	0.013	0.047	0.02	0.02	1	0.047	0.039	0.022	0.087	0.038
day	0.0045	0.0032	0.016	0.004	0.005	0.027	0.047	1	0.017	0.012	0.005	0.011
month	0.000017	0.0012	0.016	0.0024	0.002	0.01	0.039	0.017	1	0.012	0.008	0.012
year	0.34	0.01	-0.014	0.011	0.0096	0.0097	0.0022	0.012	0.012	1	0.061	0.024
dayofweek	0.014	0.025	-0.042	-0.0033	-0.034	0.049	0.087	0.0054	0.008	0.006	1	0.027
dist_travel_km	0.84	0.089	-0.047	0.187	-0.039	0.0097	0.038	0.0031	0.012	0.024	0.027	1

Dividing the dataset into training and target values

```
In [24]: x = df["pickup_longitude","pickup_latitude","dropoff_longitude","dropoff_latitude","passenger_count","hour","day","month","year","dayofweek","dist_travel_km"]

In [24]: y = df["fare_amount"]
```

Dividing the dataset into training and testing dataset

```
In [25]: from sklearn.model_selection import train_test_split
X_train,X_test,y_train,y_test = train_test_split(x,y,test_size = 0.33)

Linear Regression
```

```
In [26]: from sklearn.linear_model import LinearRegression
regression = LinearRegression()

In [26]: regression.fit(X_train,y_train)

Out[26]: LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)

In [27]: regression.intercept_ #To find the linear intercept

Out[27]: 2899.192377415925

In [28]: regression.coef_ #To find the linear coefficient

Out[28]: array([ 1.76333884e+01, -9.83170739e+00, 1.54611009e+01, -1.59707278e+01,
        5.48456386e-02, -9.46950748e-03, 1.66728920e-03, 5.48917608e-02,
        3.61743634e-01, -3.69474342e-02, 2.806177959e+00])

In [29]: prediction = regression.predict(X_test) #to predict the target values

Out[29]: [19.80422002  4.74707896  9.95283165 ... 5.89597937 17.80144322
        5.38487972]
```

```
In [30]: y_test

Out[30]:
```

	fare_amount
16850	8.50
181076	4.10
70789	9.30
87421	12.90
109443	22.25
19976	11.00
58924	13.70
199504	14.50
125215	5.30
67523	6.50
85217	22.25
150993	21.50
118795	4.10
132179	16.00
124459	3.70
173299	22.25
51448	13.70
95002	22.25
174407	16.00
78890	20.50
26798	22.25
38504	6.50
63861	12.90
171207	