

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
from haversine import haversine
import warnings
warnings.filterwarnings("ignore")
import seaborn as sns

In [2]: data = pd.read_csv('nyc_taxi_trip_duration Dataset.csv')

In [3]: #checking the dataset
data.head()
```

	id	vendor_id	pickup_datetime	dropoff_datetime	passenger_count	pickup_longitude	dropoff_longitude	dropoff_latitude	store_and_fwd_flag	trip_duration		
Out[3]:	0	4003070	1	2016-02-20 16:40:21	2016-02-20 16:47:02	1	-73.963028	-40.778873	-73.963075	40.77184	N	400
	1	4009085	1	2016-02-11 23:36:37	2016-02-11 23:53:57	2	-73.98332	-40.731743	-73.984763	40.694931	N	1100
	2	40057912	2	2016-02-21 17:59:33	2016-02-21 18:26:48	2	-73.997314	-40.721468	-40.744029	40.774958	N	1635
	3	40744273	2	2016-01-05 09:44:31	2016-01-05 10:03:32	6	-73.961670	-40.759720	-73.956779	40.780658	N	1141
	4	40232939	1	2016-02-17 06:42:23	2016-02-17 06:56:31	1	-74.017120	-40.708469	-73.988182	40.740031	N	848

```
In [4]: data.shape

Out[4]: (729322, 13)
```

```
In [5]: #checking for null values in the dataset
data.isnull().sum()
```

```
Out[5]: id                0
vendor_id              0
pickup_datetime       0
dropoff_datetime      0
passenger_count       0
pickup_longitude      0
dropoff_longitude     0
dropoff_latitude     0
store_and_fwd_flag    0
trip_duration         0
dtype: int64

• There are no null values in this dataset
```

```
In [6]: #converting datetime variable
data['pickup_datetime'] = pd.to_datetime(data['pickup_datetime'])
data['dropoff_datetime'] = pd.to_datetime(data['dropoff_datetime'])

In [7]: #generating new variables from datetime variable
data['year'] = data.pickup_datetime.dt.year
data['month'] = data.pickup_datetime.dt.month
data['weekno'] = data.pickup_datetime.dt.weekday
data['hour'] = data.pickup_datetime.dt.hour

In [8]: #using haversine for calculation of distance
def dist(df):
    pick = df['pickup_latitude'], df['pickup_longitude']
    drop = df['dropoff_latitude'], df['dropoff_longitude']
    return haversine(pick,drop)

In [9]: #Calculating speed and distance for each trip
data['distance'] = data.apply(lambda x: dist(x), axis = 1)
data['speed'] = (data.distance/data.trip_duration*3600)
```

```
In [10]: #defining encoding the categorical variable 'store_and_fwd_flag'
data = pd.concat([data, pd.get_dummies(data['store_and_fwd_flag'], prefix = 'flag'), axis=1].drop(['store_and_fwd_flag'],axis=1)

In [11]: #dropping unnecessary variables
data = data.drop(['pickup_datetime', 'dropoff_datetime', 'pickup_longitude', 'pickup_latitude', 'dropoff_longitude', 'dropoff_latitude', 'flag_N'], axis= 1)

In [12]: data.head()
```

	id	vendor_id	passenger_count	trip_duration	year	weekno	hour	distance	speed	flag_N		
Out[12]:	0	401080784	2	1	400	2016	2	0	16	1.199074	10.781668	1
	1	40090805	1	2	1100	2016	3	4	23	4.139117	13.513472	1
	2	40057912	2	2	1635	2016	2	6	17	2.90973	15.964960	1
	3	40744273	2	6	1141	2016	1	1	9	2.261101	7.449673	1
	4	40232939	1	1	848	2016	2	2	6	4.320540	18.375877	1

```
In [13]: data.shape

Out[13]: (729322, 13)
```

```
In [14]: data.dtypes
```

```
Out[14]: id                object
vendor_id              int64
passenger_count        int64
trip_duration          int64
year                  int64
month                 int64
weekno                int64
hour                  int64
distance              float64
speed                 float64
flag_N                 uint8
dtype: object
```

Removing outliers and tuning the data

```
In [15]: data.passenger_count.value_counts()
```

```
Out[15]: 1    517415
         2    189971
         3    38026
         4    29962
         6    24397
         8    14958
         9      13
         7         1
         5         1
Name: passenger_count, dtype: int64

In [16]: data.passenger_count.describe()
```

```
Out[16]: count      729322  0.00000
mean         1.662955
std          1.22446
min           0.000000
max           6.000000
25%           1.000000
50%           1.000000
75%           2.000000
max           9.000000
Name: passenger_count, dtype: float64

• There are 33 trips with 0 passenger counts. Since the mean is approx 1 we will convert 0 passenger counts to 1.
• We will remove passenger count more than 6 as we consider it an outlier.
```

```
In [17]: data['passenger_count'] = data.passenger_count.apply(lambda x : 1 if x==0 else x)
data = data[data.passenger_count < 7]
```

```
In [18]: #checking speed and trip_duration
plt.scatter(data['speed'],data['trip_duration'])
plt.show()
```

```
In [19]: #plotting trip duration to check for outliers
plt.figure(figsize=(15,5))
sns.boxplot(data.trip_duration)
plt.show()
```

- the outliers we see in this plot will cause inconsistency in our predictions.

```
In [19]: #plotting trip duration to check for outliers
plt.figure(figsize=(15,5))
sns.boxplot(data.trip_duration)
plt.show()
```

- We will focus on the data with less than 24 hours of trip duration to avoid inconsistency due to outliers.

```
In [20]: data = data[data.trip_duration <= 57600]
```

```
In [21]: plt.figure(figsize=(15,5))
sns.boxplot(data.trip_duration)
plt.show()
```

```
In [22]: #checking for outliers in speed
plt.figure(figsize=(15,5))
sns.boxplot(data.speed)
plt.show()
```

- the speed limit in NYC is 105km/h so we remove the data with speed more than 105km/h.

```
In [23]: data = data[data.speed< 105]
plt.figure(figsize=(15,5))
sns.boxplot(data.speed)
plt.show()
```

```
In [24]: #hypocasting the categorical variable and dropping the 'id' variable.
data = data.drop(['id'], axis = 1)
data['vendor_id'] = data['vendor_id'].astype('category')
```

```
In [25]: data.shape

Out[25]: (728287, 10)
```

```
In [26]: # Looking at data of less than 3 hours and less than 50 km of distance travelled.
bi_dist = data.loc[(data.trip_duration<3600) & (data.distance<50)]
plt.scatter(bi_dist.trip_duration, bi_dist.distance, s=1, alpha=0.5)
plt.xlabel('duration(seconds)')
plt.ylabel('distance')
plt.show()
```

- there are a lot of trips with more than 30 and 60 minutes which cover 0km distance. it rarely happens that passenger keeps sitting in the taxi for more than an hour and travels nowhere.

```
In [27]: #focusing on distances with only low and more than 60 minutes of trip duration.
data.loc[data['distance'] <= 1] & [data['trip_duration'] >= 3600], ['distance', 'trip_duration']].reset_index(drop=True)
```

```
In [28]: sns.regplot(min.trip_duration, min.distance)
plt.show()
```

- Although having some linear relationship, still it is not a linear distribution and shows no correlation between the two.
- these values must be removed for better consistency in the dataset.

```
In [29]: data = data[~((data['trip_duration'] <=1) & (data['distance'] <=3600))]

In [30]: data.shape

Out[30]: (728276, 10)
```

1. Build a K-Nearest neighbours model for the given dataset and find the best value of K.

```
In [31]: #separating the independent and dependent variables.
x = data.drop(['trip_duration'], axis = 1)
y = data['trip_duration']

In [32]: from sklearn.preprocessing import StandardScaler
#scaling the dataset
scaled = StandardScaler()
x_scaled = scaler.fit_transform(x)
#splitting the array into dataframe
x = pd.DataFrame(x_scaled)

In [33]: from sklearn.model_selection import train_test_split as tts

• we will use the same split data for all the models in this notebook

In [34]: #splitting the data
train_x, test_x ,train_y, test_y = tts(x,y, test_size = 0.25, random_state= 12)
train_x.shape, test_x.shape, train_y.shape , test_y.shape

Out[34]: ((546287, 9), (182089, 9), (546287,1), (182089,1))

In [35]: from sklearn.neighbors import KNeighborsRegressor as KNN

In [36]: #Creating instance of KNN
kreg = KNN()

In [37]: #reporting mean squared error and r square
from sklearn.metrics import mean_squared_error as mse
from sklearn.metrics import r2_score

In [38]: def Elbow(k):
#creating empty list
test_error = []

#training model for every value of K
for i in K:
    #instance of KNN
    kreg = KNN(n_neighbors = i, n_jobs = -1)
    #reg fit(train_x, train_y)
    kreg.fit(train_x, train_y)
    #predicting over the Train Set and calculating MSE
    pred_knn = kreg.predict(test_x)
    k = r2_score(test_y, pred_knn)
    #mse = mse(pred(test_y),
    #mse = mse(test_y,
    test_error.append(test_y)
    test_error.append(test_y)

    return test_error

In [39]: #setting k range and assigning it to a variable
k = range(1,40,2)
test = Elbow(k)

In [40]: #plotting the Elbow curve
plt.plot(k, test)
plt.xlabel('K Neighbors')
plt.ylabel('Test Mean Squared Error')
plt.title('Elbow Curve For Test')
plt.show()
```

```
In [41]: # Creating instance of KNN
kreg = KNN(n_neighbors = 5)

# fitting the model
kreg.fit(train_x, train_y)

# Predicting over the Train Set and calculating MSE
pred_knn = kreg.predict(test_x)
k = r2_score(test_y, pred_knn)
print('R squared error:', k)

R squared error: 0.7543655349328569
```

2. Build a Linear model for the given dataset with regularisation.

```
In [42]: from numpy import arange
from sklearn.linear_model import Lasso
from sklearn.metrics import r2_score
from sklearn.model_selection import RepeatedKFold
from sklearn.model_selection import RandomizedSearchCV

In [43]: from sklearn.metrics import mean_absolute_error

In [44]: #define model
lassoreg = Lasso()
# define model evaluation method
cv = RepeatedKFold(n_splits=10, n_repeats=5, random_state=1)
# define grid
grid = {'alpha': arange(0.100,0.01)}
# define search
search = RandomizedSearchCV(lassoreg, grid, scoring='neg_mean_absolute_error', cv=cv, n_jobs=-1)
# perform the search
results = search.fit(x, y)

print('mae: %.3f' % results.best_score_)
print('Config: %s' % results.best_params_)

mae = 182.161
Config: {'alpha': 14.96}

In [45]: #training the model
lassoreg = Lasso(alpha = 14.96)
lassoreg.fit(train_x,train_y)

Out[46]: Lasso(alpha=14.96)

In [47]: #calculating the r squared error
pred_lasso = lassoreg.predict(test_x)
r_squared = r2_score(test_y, pred_lasso)
print('Test R squared error is:', r_squared)

Test R squared error is: 0.688052607690843

In [48]: # putting together the coefficients and their corresponding variable names
coeff_df = pd.DataFrame()
coeff_df['columns'] = train_x.columns
coeff_df['coefficient_values'] = pd.Series(lassoreg.coef_)
print(coeff_df.head(5))

Columns coefficient values
0      0      0.000000
1      1      0.000000
2      2      0.000000
3      3      0.299509
4      4      0.000000
5      5      0.000000
6      6      0.78.394727
7      7      -128.395058
8      8      0.000000

In [49]: #plotting the coefficient values and columns.
plt.rcParams['figure.figsize'] = 20,10
plt.bar(coeff_df['columns'],coeff_df['coefficient_values'])

Out[49]: <BarContainer object of 9 artists>
```

```
In [50]: data.shape

Out[50]: (728276, 10)
```

3. Build a Random Forest model for the given dataset.

```
In [51]: #importing random forest regressor
from sklearn.ensemble import RandomForestRegressor

In [52]: Reg = RandomForestRegressor()

In [53]: #importing random search for better hyperparameters
from sklearn.model_selection import RandomizedSearchCV

In [54]: #giving inputs for hyperparameter search
n_estimators = [int(x) for x in np.linspace(start = 10, stop = 200, num = 10)]
max_depth = [int(x) for x in np.linspace(10, 40, num = 11)]
max_depth.append(None)
min_samples_split = [2,4,8]

In [55]: #putting the lists in a dictionary
random_search = {'n_estimators': n_estimators,
                  'max_features': 'max_features',
                  'max_depth': max_depth,
                  'min_samples_split': min_samples_split}
print(random_search)

{'n_estimators': [10, 31, 52, 73, 94, 115, 136, 157, 178, 200], 'max_features': ['auto', 'sqrt'], 'max_depth': [10, 19, 28, 37, 46, 55, 64, 73, 82, 91, 100, None], 'min_samples_split': [2, 4, 8]}

In [56]: random_grid = RandomizedSearchCV(estimator = Reg, param_distributions = random_search, scoring = r2_score, cv = 3, random_state = 21, verbose = 2, n_jobs = -1)

In [57]: random_grid.fit(train_x, train_y)

Fitting 3 folds for each of 10 candidates, totalling 30 fits
RandomizedSearchCV(cv=3, estimator=RandomForestRegressor(), n_jobs=-1,
                  param_distributions={'max_depth': [10, 19, 28, 37, 46, 55, 64, 73, 82, 91, 100, None],
                  'max_features': ['auto', 'sqrt'],
                  'min_samples_split': [2, 4, 8],
                  'n_estimators': [10, 31, 52, 73, 94, 115, 136, 157, 178, 200]},
                  random_state=21,
                  scoring=funcion r2_score at 0x8000023506472820,
                  verbose=2)

In [58]: random_grid.best_params_

{'n_estimators': 136,
 'min_samples_split': 4,
 'max_features': 'sqrt',
 'max_depth': 100}

In [59]: reg = RandomForestRegressor(n_estimators = 136, min_samples_split = 4, max_features = 'sqrt', max_depth = 100, random_state = 21)

In [60]: reg.fit(train_x,train_y)

RandomForestRegressor(max_depth=100, max_features='sqrt', min_samples_split=4,
                      n_estimators=136, random_state=21)

In [61]: #calculating error
pred_rf = reg.predict(test_x)
r_squared = r2_score(test_y, pred_rf)
print(r_squared)

0.89921638878106
```

4. Build a Gradient Boosting model for the given dataset.

```
In [62]: # importing gradient boosting regressor
from sklearn.ensemble import GradientBoostingRegressor as gbr

In [63]: #fitting the model
grad = gbr(random_state = 21)
grad.fit(train_x,train_y)

Out[63]: GradientBoostingRegressor(random_state=21)

In [64]: #calculating error
pred_gb = grad.predict(test_x)
Error_rm = r2_score(test_y, preds)
print(Error_rm)

0.94533065210191

In [65]: from sklearn.model_selection import RandomizedSearchCV

In [66]: #setting inputs for hyperparameter tuning
n_estimators = [int(x) for x in np.linspace(start = 10, stop = 200, num = 10)]
max_depth = [int(x) for x in np.linspace(10,120, num = 11)]
max_depth.append(None)
learning_rate = [0.15,0.1,0.3,0.5,0.7,0.9,0.05]

In [67]: random_search1 = {'n_estimators': n_estimators,
                        'learning_rate': [0.15,0.1,0.3,0.5,0.7,0.9,0.05],
                        'max_depth': max_depth}
print(random_search1)

{'n_estimators': [10, 31, 52, 73, 94, 115, 136, 157, 178, 200], 'learning_rate': [0.15, 0.1, 0.3, 0.5, 0.7, 0.9, 0.05], 'max_depth': [10, 21, 32, 43, 54, 65, 76, 87, 98, 109, 120, None]}

In [68]: random_grid1 = RandomizedSearchCV(estimator = grad, param_distributions = random_search1, scoring = r2_score, cv= 3, random_state = 21, verbose = 2, n_jobs = -1)

In [69]: random_grid1.fit(train_x,train_y)

Fitting 3 folds for each of 10 candidates, totalling 30 fits
RandomizedSearchCV(cv=3, estimator=GradientBoostingRegressor(random_state=21),
                  n_jobs=-1,
                  param_distributions={'learning_rate': [0.15, 0.1, 0.3, 0.5, 0.7, 0.9, 0.05],
                  'max_depth': [10, 21, 32, 43, 54, 65, 76, 87, 98, 109, 120, None],
                  'n_estimators': [10, 31, 52, 73, 94, 115, 136, 157, 178, 200]},
                  random_state=21,
                  scoring=funcion r2_score at 0x8000023506472820,
                  verbose=2)

In [70]: random_grid1.best_params_

{'n_estimators': 84, 'max_depth': None, 'learning_rate': 0.7}

In [71]: gbr = gbr(n_estimators = 94, max_depth = None, learning_rate = 0.7)
gbr.fit(train_x,train_y)

GradientBoostingRegressor(learning_rate=0.7, max_depth=None, n_estimators=94)

In [72]: #calculating error
pred_gb = gbr.predict(test_x)
error = r2_score(test_y, pred_gb)

In [73]: # r squared error
print(error)

0.95912963718018

In [74]: #getting train score
pr = gbr.predict(train_x)
print(r2_score(train_y,pr))

0.99936163002543
```

5. Combine all the models above using the averaging technique to generate the final predictions.

```
In [66]: # averaging all the above models:
w_avg = (pred_lasso) + (pred_knn) + (pred_rf) + (pred_gb))/(4)
r_sq = r2_score(test_y,w_avg)
print('R squared error of weighted average of all models is: ',r_sq)

R squared error of weighted average of all models is: 0.8703041881535512

In [ ] :
```