





```
(49): hour
AveSpeed 0.01
TotalCharge -0.01
Distance -0.00
trip_duration 0.02
day 0.02
Passengers 0.02
dayofweek 0.65
TripCount 1.00
month 0.00
Name: TripCount, dtype: float64

In (50):
plt.figure(figsize=(16,9))
sns.heatmap(df5.corr(),cmap="coolwarm",annot=True,fmt=".2f",linewidths=2)
plt.title("", fontsize=20)
plt.show()

Passengers Distance TotalCharge AveSpeed trip_duration month day dayofweek hour TripCount
Passengers 1.00 0.02 0.02 0.01 0.01 0.02 0.02 0.03 0.01 0.02
Distance 0.02 1.00 0.91 0.83 0.16 0.07 0.04 0.04 0.00
TotalCharge 0.02 0.91 1.00 0.63 0.24 0.01 0.01 0.08 0.01
AveSpeed 0.01 0.83 0.63 1.00 0.02 0.13 0.08 0.01 0.01
trip_duration 0.01 0.16 0.24 0.02 1.00 0.03 0.01 0.03 0.01
month 0.02 0.07 0.01 0.13 0.03 1.00 0.04 0.01 0.02
day 0.03 0.04 0.01 0.06 0.01 0.04 1.00 0.02 0.65
dayofweek 0.01 0.04 0.09 0.01 0.03 0.01 0.02 1.00 0.06
hour 0.02 0.00 0.01 0.01 0.01 0.02 0.65 0.06 1.00
TripCount 0.02 0.00 0.01 0.01 0.01 0.02 0.06 0.01 0.00

Data Preprocessing

Drop unwanted features

In (51):
df5.columns

Out(51):
Index(['PickupTime', 'Passengers', 'Distance', 'TotalCharge', 'AveSpeed', 'trip_duration', 'month', 'day', 'dayofweek', 'hour', 'Location', 'TripCount'], dtype='object')

In (52):
df5.drop(['PickupTime', 'month', 'Location'], axis=1, inplace=True)

In (53):
df5

Out(53):
Passengers Distance TotalCharge AveSpeed trip_duration day dayofweek hour TripCount
0 1 6.51 21.80 19.91 29 1 1 12 14
1 1 1 100 8.80 5.89 56 1 1 14 14
2 1 9.09 42.41 21.54 71 1 1 14 14
3 2 0.99 14.04 3.78 61 1 1 14 14
4 2 3.20 24.36 6.49 73 1 1 14 14
... ..
16118 4 0.82 8.80 4.45 43 31 3 12 37
16119 3 3.05 15.38 14.52 34 31 3 12 37
16120 1 1.87 8.80 15.37 27 31 3 12 37
16121 1 8.70 39.99 25.53 70 31 3 10 37
16122 2 2.80 15.80 15.87 55 31 3 23 37

16123 rows x 9 columns

Treat Missing Values

In (54):
df5.isnull().sum()

Out(54):
Passengers 0
Distance 0
TotalCharge 0
AveSpeed 0
trip_duration 0
day 0
dayofweek 0
hour 0
TripCount 0
dtype: int64

Treat Duplicate Values

In (55):
df5.duplicated(keep="first").sum()

Out(55):
0

Treat Data Types

In (56):
df5.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 16123 entries, 0 to 16122
Data columns (total 9 columns):
 # Column Non-Null Count Dtype
---
 0 Passengers 16123 non-null int64
 1 Distance 16123 non-null float64
 2 TotalCharge 16123 non-null float64
 3 AveSpeed 16123 non-null float64
 4 trip_duration 16123 non-null int32
 5 day 16123 non-null int64
 6 dayofweek 16123 non-null int64
 7 hour 16123 non-null int64
 8 TripCount 16123 non-null int64
dtypes: float64(3), int32(1), int64(5)
memory usage: 1.2 MB

Create and save processed dataset

In (57):
df5.to_csv("train2.csv",index=False)

In (58):
df5.shape

Out(58):
(16123, 9)

Train Test Split

In (59):
X = df5.iloc[:,1:8]
y = df5.iloc[:,8]

In (60):
X.values, y.values

Out(60):
(array([[ 1., 6.51, 21.8, ..., 1., 1., 12., 14.],
[ 1., 1., 100, ..., 1., 1., 14., 14.],
[ 1., 9.09, 42.41, ..., 1., 1., 14., 14.],
[ 2., 0.99, ..., 1., 1., 14., 14.],
[ 2., 3.2, ..., 1., 1., 14., 14.],
[ 1., 1.87, 8.8, ..., 31., 3., 12., 37.],
[ 1., 8.7, 39.99, ..., 31., 3., 10., 37.],
[ 2., 2.8, 15.8, ..., 31., 3., 23., 37.]],
array([14, 14, 14, ..., 37, 37, 37], dtype=int64))

In (61):
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)

In (62):
X_train.shape, X_test.shape, y_train.shape, y_test.shape

Out(62):
(12898, 8), (3225, 8), (12898,), (3225,1)

Feature Scaling

In (63):
X_train

Out(63):
Passengers Distance TotalCharge AveSpeed trip_duration day dayofweek hour
9228 2 1.78 11.16 13.46 39 17 3 0
12106 1 1.40 12.25 8.26 51 26 5 23
12415 1 3.60 17.16 16.12 46 24 3 23
9091 1 1.18 8.80 6.78 25 17 3 13
... ..
13123 3 3.20 11.80 20.35 9 26 5 10
3264 2 3.92 20.80 8.77 61 7 0 18
9845 2 1.65 12.80 5.93 35 19 5 20
10799 2 10.30 44.19 23.75 85 20 6 17
2732 1 1.29 12.30 5.64 50 5 5 20

12898 rows x 8 columns

In (64):
scaler = StandardScaler()

In (65):
X_train_scaled = scaler.fit_transform(X_train)

In (66):
X_test_scaled = scaler.transform(X_test)

In (67):
X_test_scaled

Out(67):
array([[ -0.53436466, 1.6951792, 2.26092482, ..., 1.60454256,
-0.43213072, 0.9375585, ..., 0.53582371, ..., 0.15876683,
0.21454715, -0.6600744, -0.85358371, ..., 0.15876683,
0.07986836, -2.39972044, ..., 0.15876683,
-1.10386653, -1.43810311],
...,
[0.21454715, -0.70806551, -0.67548768, ..., 0.38119386,
1.10386653, 0.9375585, ..., 0.49240738,
0.21454715, 2.4851859, 2.73331367, ..., 1.38211553,
1.61586561, -0.43696669],
[-0.53436466, -0.84096383, -0.72978528, ..., -1.17579539,
1.10386653, 0.9375585 ]])

In (68):
X_test_scaled

Out(68):
array([[ -0.53436466, -0.72652361, -0.59404133, ..., -0.50851428,
0.59186745, -0.56421894],
[0.21454715, -0.84096383, -1.05557066, ..., 1.04847497,
0.59186745, 1.10442244],
[-0.53436466, -0.695074, -0.45829741, ..., -1.62064947,
-0.9441298, 0.60383063],
...,
[-0.53436466, -0.4681102, -0.67548768, ..., -1.17579539,
1.10386653, -2.06599237],
1.46128226, -0.33150205, -0.72978525, ..., 1.38211553,
-1.45612888, 0.10321881],
[-0.53436466, 0.03395014, -0.51259498, ..., 1.71575608,
0.07986836, 1.10442244]])

Model Training

Using PyCaret

In (69):
from pycaret.regression import *

In (70):
exp_reg = setup(data = df5, target = "TripCount", session_id=0, normalize=True,
train_size = 0.8, numeric_features=['Passengers', 'trip_duration'])

Description Value
0 session_id 0
1 Target TripCount
2 Original Data (16123, 9)
3 Missing Values 7
4 Numeric Features False
5 Categorical Features 1
6 Ordinal Features False
7 High Cardinality Features False
8 High Cardinality Method None
9 Transformed Train Set (12898, 14)
10 Transformed Test Set (3225, 14)
11 Shuffle Train-Test True
12 Stratify Train-Test False
13 Fold Generator KFold
14 Fold Number 10
15 CPU Jobs -1
16 Use GPU False
17 Log Experiment False
18 Experiment Name reg-default-name
19 USI 72bb
20 Imputation Type simple
21 Iterative Imputation Iteration None
22 Numeric Imputer mean
23 Iterative Imputation Numeric Model None
24 Categorical Imputer constant
25 Iterative Imputation Categorical Model None
26 Unknown Categoricals Handling least_frequent
27 Normalize True
28 Normalization Method zscore
29 Transformation False
30 Transformation Method None
31 PCA False
32 PCA Method None
33 PCA Components None
34 Ignore Low Variance False
35 Combine Rare Levels False
36 Rare Level Threshold None
37 Numeric Binning False
38 Remove Outliers False
39 Outliers Threshold None
40 Remove Multicollinearity False
41 Multicollinearity Threshold None
42 Remove Perfect Collinearity True
43 Clustering False
44 Clustering Iteration None
45 Polynomial Features False
46 Polynomial Degree None
47 Trigonometry Features False
48 Polynomial Threshold None
49 Group Features False
50 Feature Selection False
51 Feature Selection Method classic
52 Feature Selection Threshold None
53 Feature Interaction False
54 Feature Ratio False
55 Interaction Threshold None
56 Transform Target False
57 Transform Target Method box-cox

In (71):
compare_models(exclude=['omp','lbc','lact','lpat','lansac','lct','huber','lct','ada','mlp','xgboost','lightgbm'],to

Model MAE MSE RMSE R2 RMSLE MAPE TT (Sec)
dt Decision Tree Regressor 0.0000 0.0000 0.0000 1.0000 0.0000 0.0000 0.0300
rf Random Forest Regressor 0.0000 0.0000 0.0000 1.0000 0.0000 0.0000 0.6700
et Extra Trees Regressor 0.0012 0.0002 0.0143 1.0000 0.0006 0.0000 0.6300
gbr Gradient Boosting Regressor 0.1999 0.0641 0.2526 0.9993 0.0119 0.0099 0.6300
knn K Neighbors Regressor 1.2292 7.4005 2.7193 0.9150 0.1186 0.0613 0.1100
lr Linear Regression 4.7568 35.7548 5.9787 0.5896 0.2802 0.2526 0.1400
ridge Ridge Regression 4.7568 35.7548 5.9787 0.5896 0.2802 0.2526 0.1400
lar Least Angle Regression 4.7568 35.7548 5.9787 0.5896 0.2802 0.2525 0.0160
lasso Lasso Regression 7.2695 73.5511 8.5757 0.1557 0.4035 0.4177 0.0160
en Elastic Net 7.4659 75.9660 8.7154 0.1280 0.4095 0.4280 0.0140
llar Lasso Least Angle Regression 8.0209 87.1459 9.3348 -0.0004 0.4341 0.4594 0.0140
dummy Dummy Regressor 8.0209 87.1459 9.3348 -0.0004 0.4341 0.4594 0.0140

Out(71):
DecisionTreeRegressor(ccp_alpha=0.0, criterion='mse', max_depth=None,
max_features=None, max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, n_estimators=100,
n_iter_no_change=None, presort='deprecated',
random_state=None, subsample=1.0, tol=0.0001,
validation_fraction=0.1, verbose=0, warm_start=False)

In (72):
model_selected = create_model('gbr')

MAE MSE RMSE R2 RMSLE MAPE
0 0.2004 0.0635 0.2521 0.9993 0.0124 0.0103
1 0.2116 0.0699 0.2644 0.9992 0.0124 0.0103
2 0.1990 0.0635 0.2520 0.9993 0.0111 0.0095
3 0.2467 0.0929 0.3048 0.9990 0.0139 0.0119
4 0.4998 0.1055 0.3249 0.9988 0.0158 0.0127
5 0.2336 0.0787 0.2806 0.9991 0.0130 0.0115
6 0.1989 0.0628 0.2506 0.9993 0.0128 0.0104
7 0.2326 0.0851 0.2917 0.9990 0.0139 0.0117
8 0.2344 0.0899 0.2998 0.9990 0.0133 0.0112
9 0.2175 0.0690 0.2623 0.9992 0.0118 0.0106
Mean 0.2225 0.0781 0.2783 0.9991 0.0130 0.0110
SD 0.0185 0.0140 0.0247 0.0002 0.0012 0.0009

In (73):
print(model_selected)

GradientBoostingRegressor(alpha=0.9, ccp_alpha=0.0, criterion='friedman_mse',
init=None, learning_rate=0.1, loss='ls', max_depth=3,
max_features=None, max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, n_estimators=100,
n_iter_no_change=None, presort='deprecated',
random_state=None, subsample=1.0, tol=0.0001,
validation_fraction=0.1, verbose=0, warm_start=False)

In (74):
plot_model(model_selected)

Residuals for GradientBoostingRegressor
Train R^2 = 0.999
Test R^2 = 0.999
Residuals
Predictor Value
Distribution

plot_model(model_selected, plot = 'error')

Prediction Error for GradientBoostingRegressor
best fit
identity
R^2 = 0.999
y-hat
y
y-hat

In (76):
plot_model(model_selected, plot='feature')

Feature Importance Plot
day
dayofweek_0
dayofweek_1
dayofweek_2
dayofweek_4
Features
dayofweek_6
dayofweek_5
dayofweek_3
hour
AveSpeed
Variable Importance

Using Regression or Classification Models

In (77):
reg_model = GradientBoostingRegressor()

In (78):
reg_model.fit(X_train_scaled,y_train)

Out(78):
GradientBoostingRegressor(alpha=0.9, ccp_alpha=0.0, criterion='friedman_mse',
init=None, learning_rate=0.1, loss='ls', max_depth=3,
max_features=None, max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, n_estimators=100,
n_iter_no_change=None, presort='deprecated',
random_state=None, subsample=1.0, tol=0.0001,
validation_fraction=0.1, verbose=0, warm_start=False)

In (79):
y_pred = reg_model.predict(X_test_scaled)

In (80):
y_pred

Out(80):
array([37.08322843, 20.02907741, 13.81554489, ..., 24.40743378,
24.67452898, 36.79374546])

Model Evaluation

In (81):
mse = mean_squared_error(y_test,y_pred)
mse

Out(81):
0.1080166158851585

In (82):
rmse = np.sqrt(mse)
rmse

Out(82):
0.3286588137950335

In (83):
r2score = r2_score(y_test,y_pred)
r2score

Out(83):
0.9987575637140792

In (84):
fig, ax = plt.subplots(figsize=(10,8))
sns.regplot(x=y_test, y=y_pred, ax=ax)
plt.title("Plot to compare actual vs predicted", fontsize=20)
plt.ylabel("Predicted")
plt.xlabel("Actual")
plt.show()

Plot to compare actual vs predicted
Predicted
Actual

Plot Feature Importances

reg_model.feature_importances_

array([0., 0., 0., 0., 0., 0., 0., 0., 0.])

feat_importances = pd.Series(reg_model.feature_importances_, index=X.columns)

feat_importances

Passengers 0.00
Distance 0.00
TotalCharge 0.00
AveSpeed 0.00
trip_duration 0.00
day 0.62
dayofweek 0.23
hour 0.00
dtype: float64

In (88):
feat_importances.nlargest(10).plot(kind='barh', figsize=(10,10))
plt.title("Feature Importances")
plt.show()

Feature Importances
Passengers
TotalCharge
trip_duration
Distance
hour
AveSpeed
day
dayofweek

Cross-Validation

In (89):
cv = cross_val_score(reg_model,X,y,cv=5,verbose=1,scoring='neg_root_mean_squared_error')

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
[Parallel(n_jobs=1)]: Done 5 out of 5 | elapsed: 7.5s finished

In (90):
cv.mean()

Out(90):
-6.323263197229591

Model Application, Results, and Analysis

In (91):
test_data = pd.DataFrame([10,3,30,0,20,0,5,6,12]).T

In (92):
test_data

Out(92):
0 10 3 30 0 20 0 5 6 12
1 40 10 30 10 20 10 5 6 10 7

In (93):
trip_count_prediction = reg_model.predict(test_data)

In (94):
trip_count_prediction

Out(94):
array([32.79617035])

Using Gradient Boosting Regressor, the model was able to achieve less than 4.9 RMSE score. We found out day and day of week is most important for taxi trip counts. Therefore deployment of taxi fleets must be targeted on customers who need them on special days or day of the week.
```