


```
[24]: for u in data_1.col_counts:
      if (Type[u]=="object"):
          print(data_1[u].value_counts())

gas      185
diesel   20
Name: Fuel Type, dtype: int64
std      168
turbo     37
Name: Aspiration, dtype: int64
four     115
two       90
Name: Door Number, dtype: int64
hatchback 70
wagon     25
hardtop    8
convertible 6
Name: Car Body, dtype: int64
fwd      129
rwd       76
Name: Drive Wheel, dtype: int64
front    202
rear       3
Name: Engine Location, dtype: int64
cnc      148
ohcf      15
ohcv      13
dohc     12
l         12
c         4
rotor     4
dohcv     1
Name: enginetype, dtype: int64
four     159
six       24
five      11
eight      5
two        4
three      1
twelve     1
Name: Cylinder Number, dtype: int64
mpgfi     94
dbsl      66
idb1      11
spdi       9
ebbl       3
mfi        1
mpfi       1
Name: Fuel System, dtype: int64

In [25]: from sklearn.preprocessing import LabelEncoder
         label=LabelEncoder()

In [26]: data_1["Fuel Type"]=label.fit_transform(data_1["Fuel Type"])
         data_1["Aspiration"]=label.fit_transform(data_1["Aspiration"])
         data_1["Door Number"]=label.fit_transform(data_1["Door Number"])
         data_1["Car Body"]=label.fit_transform(data_1["Car Body"])
         data_1["Drive Wheel"]=label.fit_transform(data_1["Drive Wheel"])
         data_1["Engine Location"]=label.fit_transform(data_1["Engine Location"])
         data_1["enginetype"]=label.fit_transform(data_1["enginetype"])
         data_1["Cylinder Number"]=label.fit_transform(data_1["Cylinder Number"])
         data_1["Fuel System"]=label.fit_transform(data_1["Fuel System"])

In [27]: data_1.dtypes

Out[27]: car_ID              int64
         symboling       int64
         Fuel Type       int32
         Aspiration       int32
         Door Number     int32
         Car Body        int32
         Drive Wheel     int32
         Engine Location float64
         wheelbase       float64
         carlength       float64
         carwidth        float64
         carheight       float64
         curbweight      int64
         enginetype      int32
         Cylinder Number int32
         engineize       int64
         Fuel System     int32
         Bore/ratio      float64
         stroke          float64
         compressionratio float64
         horsepower      int64
         peakrpm         int64
         citympg         int64
         highway(mpg)    int64
         price           float64
         dtype: object

Result: Almost every feature has been converted to integer format for easy access of the model.
```

7.2 Separating train and test data:

The train-test split is used to estimate the performance of machine learning algorithms that are applicable for prediction-based Algorithms/Applications. This method is a fast and easy procedure to perform such that we can compare our own machine learning model results to machine results.

```
In [28]: X=data_1.drop(["price"],axis=1)
         y=data_1["price"]

In [29]: 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=45)

In [30]: X_train.shape,X_test.shape

Out[30]: (143, 24), (62, 24)


```

Result: The data is separated 70% train data and 30% test data for model performance.

7.3 Feature Scaling:

Feature Scaling is a technique to standardize the independent features present in the data in a fixed range. So the X_train and X_test data should be standardized for easy performance. In standardizing, the mean and range should be 0 to 1.

```
In [31]: from sklearn.preprocessing import StandardScaler
         scaler=StandardScaler()
         X_train_scaler=scaler.fit_transform(X_train)
         X_test_scaler=scaler.transform(X_test)

In [32]: X_test_scaler

Out[32]: array([[ 0.7751997, -0.61608312,  0.31622777, ..., -0.73763795,
        -0.3217526, -0.24523878],
       [-0.63737597, -0.61608312,  0.31622777, ..., -0.73763795,
        0.11230987, -0.17223105],
       [ 1.14516, -1.03063973,  0.31622777, ..., -0.73763795,
        -0.17706511, -0.1060755 ],
       ...,
       [-0.33468118,  1.85400715,  0.31622777, ..., -0.31905768,
        -0.90050257, -0.94105516],
       [ 1.48148754, -1.85400715,  0.31622777, ...,  0.72739298,
        -0.17706511, -0.24523878],
       [ 0.85928159, -0.61608312,  0.31622777, ..., -0.73763795,
        0.83574733,  0.86806743]])

8.Modeling

8.1 Defining baseline model:

A baseline model is essentially a simple model that acts as a reference in a machine learning project. Its main function is to contextualize the results of trained models. Baseline models usually lack complexity and may have little predictive power.
```

```
In [33]: from sklearn.neighbors import KNeighborsRegressor # Importing KNN
         from sklearn.ensemble import RandomForestRegressor # Importing Random Forest Regressor
         from sklearn.ensemble import BaggingRegressor # Importing Bagging Regressor
         from sklearn.tree import DecisionTreeRegressor # Importing DecisionTreeRegressor
         from sklearn.linear_model import LinearRegression # Importing LinearRegression

# Defining the scores list
model_scores = []

# Defining a list of useful regressors
regressors = [RandomForestRegressor(),
              KNeighborsRegressor(),
              BaggingRegressor(),
              LinearRegression()]

In [34]: for i in regressors:

         #Importing model name
         model_name=type(i).__name__

         # Fitting the train data in model
         i.fit(X_train,Y_train)

         #Predict the train data using model
         y_pred_train=i.predict(X_train)

         #Predict the test data using model
         y_pred_test=i.predict(X_test)

         #Calculating train R2 score
         R2_score_train=r2_score(Y_train,y_pred_train)

         #calculating test R2 score
         R2_score_test=r2_score(Y_test,y_pred_test)

         print("Model name:",model_name)
         print("#R2_score_train:", R2_score_train)
         print("# R2_score_test:", R2_score_test)

         model_scores.append((model_name,R2_score_train,R2_score_test))

Model name: RandomForestRegressor
R2_score_train: 0.9893602862253864
R2_score_test: 0.9158547953548067
Model name: KNeighborsRegressor
R2_score_train: 0.9078959465171196
R2_score_test: 0.8158167202989743
Model name: BaggingRegressor
R2_score_train: 0.98744884177938
R2_score_test: 0.8898246057466188
Model name: LinearRegression
R2_score_train: 0.9038905143120971
R2_score_test: 0.8456306246862354

R2 score is the proportion of the variance in the dependent variable that is predictable from the independent variable(s).
```

```
In [35]: model_scores

Out[35]: [('RandomForestRegressor', 0.9893602862253864, 0.9158547953548067),
         ('KNeighborsRegressor', 0.9078959465171196, 0.8158167202989743),
         ('BaggingRegressor', 0.98744884177938, 0.8898246057466188),
         ('LinearRegression', 0.9038905143120971, 0.8456306246862354)]

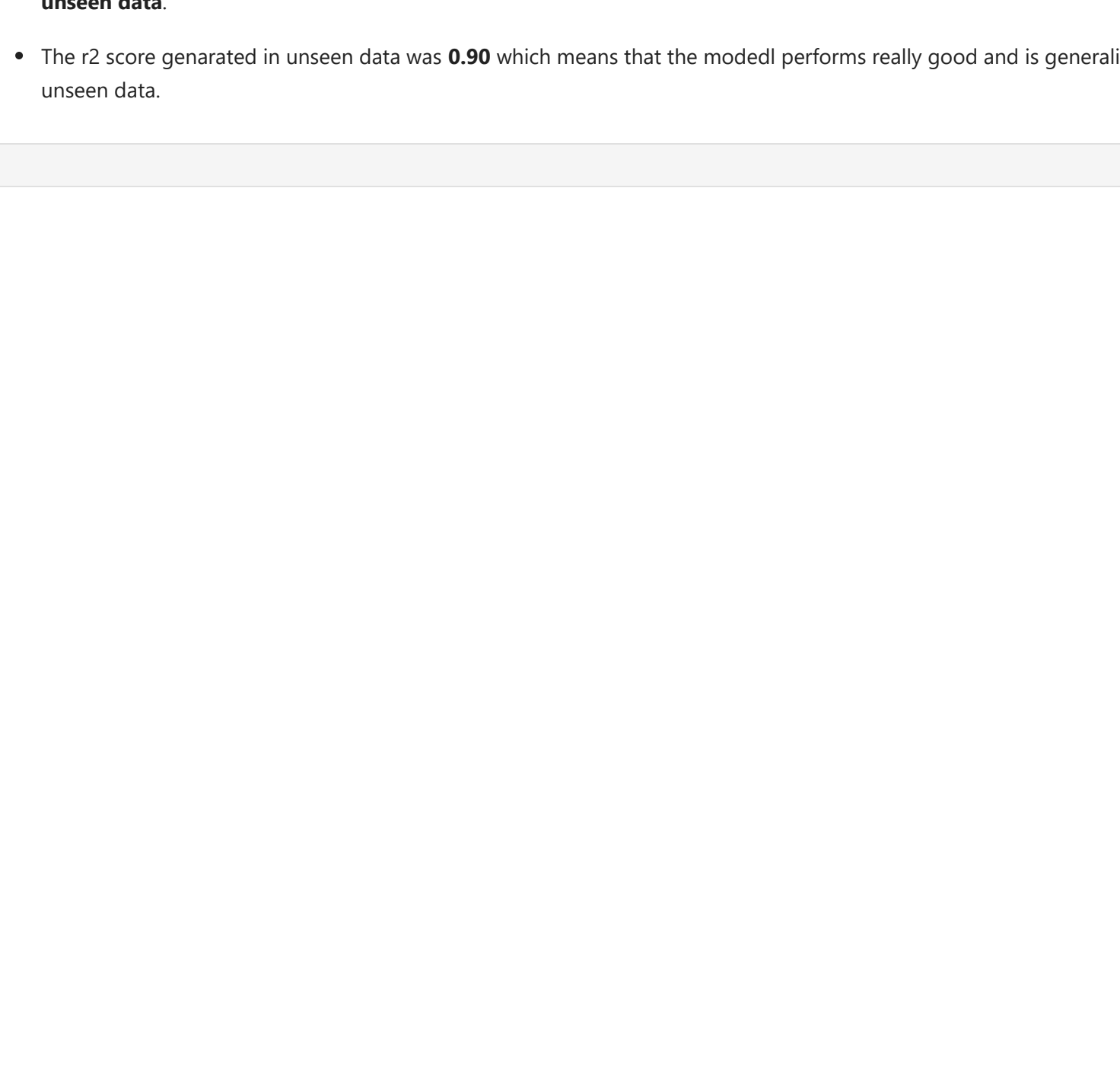
In [36]: Model=pd.DataFrame(data=model_scores,columns=["model_name","R2_score_train","R2_score_test"])

In [37]: Model

Out[37]:   model_name  R2_score_train  R2_score_test
0  RandomForestRegressor    0.989360    0.915855
1      KNeighborsRegressor    0.907996    0.815817
2      BaggingRegressor    0.987449    0.889825
3      LinearRegression    0.903891    0.845631

# To find the best test R2 score(plot bar chart)

In [39]: plt.figure(figsize=(12,6))
         sns.barplot(Model.model_name,Model.R2_score_test)


```

Observation:-

- RandomForest Regression gives the best r2 score and Low Basis and low variance (low overfitting) among the models, so the RandomForest Regression will be used for Hyperparameter tuning.

8.2 Hyperparameters Tuning:

Hyperparameters are parameters whose values control the learning process and determine the values of model parameters that a learning algorithm ends up learning and also used to find out the best model.

- We will be using Random Search in order to find the best values.
- We will consider RandomForest Regression as they have given best results

```
In [53]: para_grid={'n_estimators': [10, 17, 25, 33, 41, 48, 56, 64, 72, 80],
                 'max_features': ['auto', 'sqrt'], 'max_depth': [2, 4],
                 'min_samples_split': [2, 5],
                 'min_samples_leaf': [1, 2], 'bootstrap': [True, False]}

RFR =RandomForestRegressor(random_state=42)
random_tuning_model = RandomizedSearchCV(estimator = RFR,
                                         param_distributions = para_grid,
                                         scoring = 'r2',
                                         cv = 5)

random_tuning_model.fit(X_train, Y_train)

Y_pred=random_tuning_model.predict(X_train)

# Printing metrics
print("Hyperparameters: ", random_tuning_model.best_params_)
print("Train Score: ", random_tuning_model.best_score_)
print("Validation Score: ", r2_score(Y_test, y_pred_tuned))

(Hyperparameters): {'n_estimators': 33, 'min_samples_split': 5, 'min_samples_leaf': 1, 'max_features': 'auto',
'Train Score): 0.9106638786861891
'Validation Score): 0.966948182463496
```

Observation:-

we found the best [hyperparameters]: {'n_estimators': 33, 'min_samples_split': 5, 'min_samples_leaf': 1, 'max_features': 'auto', 'max_depth': 4, 'bootstrap': True} [Train Score]: 0.9146482211854652 [Validation Score]: 0.969134496856467

9. Test set

```
In [54]: # Predicting with the best fit parameters
         best_fit = random_tuning_model.best_estimator_

In [55]: # Fitting the train data in best fit model
         best_fit.fit(X_train,Y_train)

Out[55]: RandomForestRegressor(max_depth=4, min_samples_split=5, n_estimators=33,
                               random_state=42)

In [56]: #Predicting the unseen data
         y_pred_tuned=best_fit.predict(X_test)

In [60]: #R2_score of unseen data
         best_fit_R2_score=r2_score(Y_test,y_pred_tuned)
         best_fit_R2_score

Out[60]: 0.8937774595844317

In [58]: # Making submission file
         submission = pd.DataFrame({"car_ID": X_test['car_ID'], 'Salary': y_pred_tuned})
         submission.head(10)

Out[58]:   car_ID  Salary
148   149   9468.919434
64    65   9095.618948
170   171  14362.109066
72    73  34135.200505
25    26   6354.463336
13    14  17790.377239
79    80   8713.626674
83    84  15939.809462
3     4  11048.954316
84    85  15939.809462
```

10. Conclusion

- In this case study the given data was analysed and on top of that a regression model was built.
- It can be found that enginysize,stroke and curbweight are making the price high in the American Market.
- The model chosen for this case study was a RandomForest Regression as it was returning the least overfitting and best r2 score on unseen data.
- The r2 score generated in unseen data was 0.90 which means that the model performs really good and is generalizing well on unseen data.

```
In [ ]:
```