

INT - 404

Artificial Intelligence

Assignment

TASK:

To make a program in Python to predict the price of Used Cars.

Submitted to: Shabnam Ma'am
Submitted by: <ul style="list-style-type: none">● Sumit Kumar Rajput 11802025● Rhituraj Sarkar 11802139● Siddharth Sahoo 11802128
Section: K18GA
GitHub link : https://github.com/sumit-rk/AI-Project

Techniques Employed in Solving the problem :

- Linear Regression
- Decision Tree
- Random Forest Regression

The problem is being solved using the above techniques, and then the results are compared based on the accuracy of each method.

The training technique used is test and split, so that a particular part of the data is used to train the model and the other part of the data is used for the predictions.

And then based on the Root Mean Squared Error, the best model is selected.

Solving Methodology:

- First, the data is taken from the CSV file.
- Then, the data is inspected in order to determine the attributes that have a significant role in determining the price, like
- How many kilometers is the vehicle used for?
- When was the vehicle purchased?
- If there is any damage? etcetera...

Then based on these attributes the model is trained. And the predictions are made using the models.

CODE:

```
In [1]: 1 import matplotlib.pyplot as plt
2 import pandas as pd
3 #from pandas.tools.plotting import scatter_matrix
4 import pylab as pl
5 import numpy as np
6 from sklearn.pipeline import Pipeline
7 from sklearn.pipeline import FeatureUnion
8 from sklearn.preprocessing import StandardScaler
9 from sklearn.base import BaseEstimator, TransformerMixin
10 from sklearn.linear_model import LinearRegression
11 from sklearn.model_selection import train_test_split
12 #from sklearn.preprocessing import Imputer
13 from sklearn.preprocessing import OneHotEncoder
14 import os
15 %matplotlib inline
```

```
In [2]: 1 data_frame = pd.read_csv("autos.csv",encoding="latin-1")
```

```
In [3]: 1 data_frame.head()
```

```
In [3]: 1 data_frame.head()
```

Out[3]:

	dateCrawled	name	seller	offerType	price	abtest	vehicleType	yearOfRegistration	gearbox	powerPS	model	kilometer	mont
0	2016-03-24 11:52:17	Golf_3_1.6	privat	Angebot	480	test	NaN	1993	manuell	0	golf	150000	
1	2016-03-24 10:58:45	A5_Sportback_2.7_Tdi	privat	Angebot	18300	test	coupe	2011	manuell	190	NaN	125000	
2	2016-03-14 12:52:21	Jeep_Grand_Cherokee_"Overland"	privat	Angebot	9800	test	suv	2004	automatik	163	grand	125000	
3	2016-03-17 16:54:04	GOLF_4_1_4__3TÜRER	privat	Angebot	1500	test	kleinwagen	2001	manuell	75	golf	150000	
4	2016-03-31 17:25:20	Skoda_Fabia_1.4_TDI_PD_Classic	privat	Angebot	3600	test	kleinwagen	2008	manuell	69	fabia	90000	

```
In [4]: 1 data_frame.describe()
```

Out[4]:

	price	yearOfRegistration	powerPS	kilometer	monthOfRegistration	nrOfPictures	postalCode
count	3.715280e+05	371528.000000	371528.000000	371528.000000	371528.000000	371528.0	371528.000000
mean	1.729514e+04	2004.577997	115.549477	125618.688228	5.734445	0.0	50820.66764
std	3.587954e+06	92.866598	192.139578	40112.337051	3.712412	0.0	25799.08247
min	0.000000e+00	1000.000000	0.000000	5000.000000	0.000000	0.0	1067.000000
25%	1.150000e+03	1999.000000	70.000000	125000.000000	3.000000	0.0	30459.000000
50%	2.950000e+03	2003.000000	105.000000	150000.000000	6.000000	0.0	49610.000000
75%	7.200000e+03	2008.000000	150.000000	150000.000000	9.000000	0.0	71546.000000
max	2.147484e+09	9999.000000	20000.000000	150000.000000	12.000000	0.0	99998.000000

```
In [5]: 1 missing_values = data_frame.isnull().sum()
2 missing_values
```

Out[5]:

dateCrawled	0
name	0
seller	0
offerType	0
price	0
abtest	0
vehicleType	37869
yearOfRegistration	0
gearbox	20209
powerPS	0
model	20484
kilometer	0
monthOfRegistration	0
fuelType	33386
brand	0
notRepairedDamage	72060
dateCreated	0
nrOfPictures	0
postalCode	0
lastSeen	0

dtype: int64

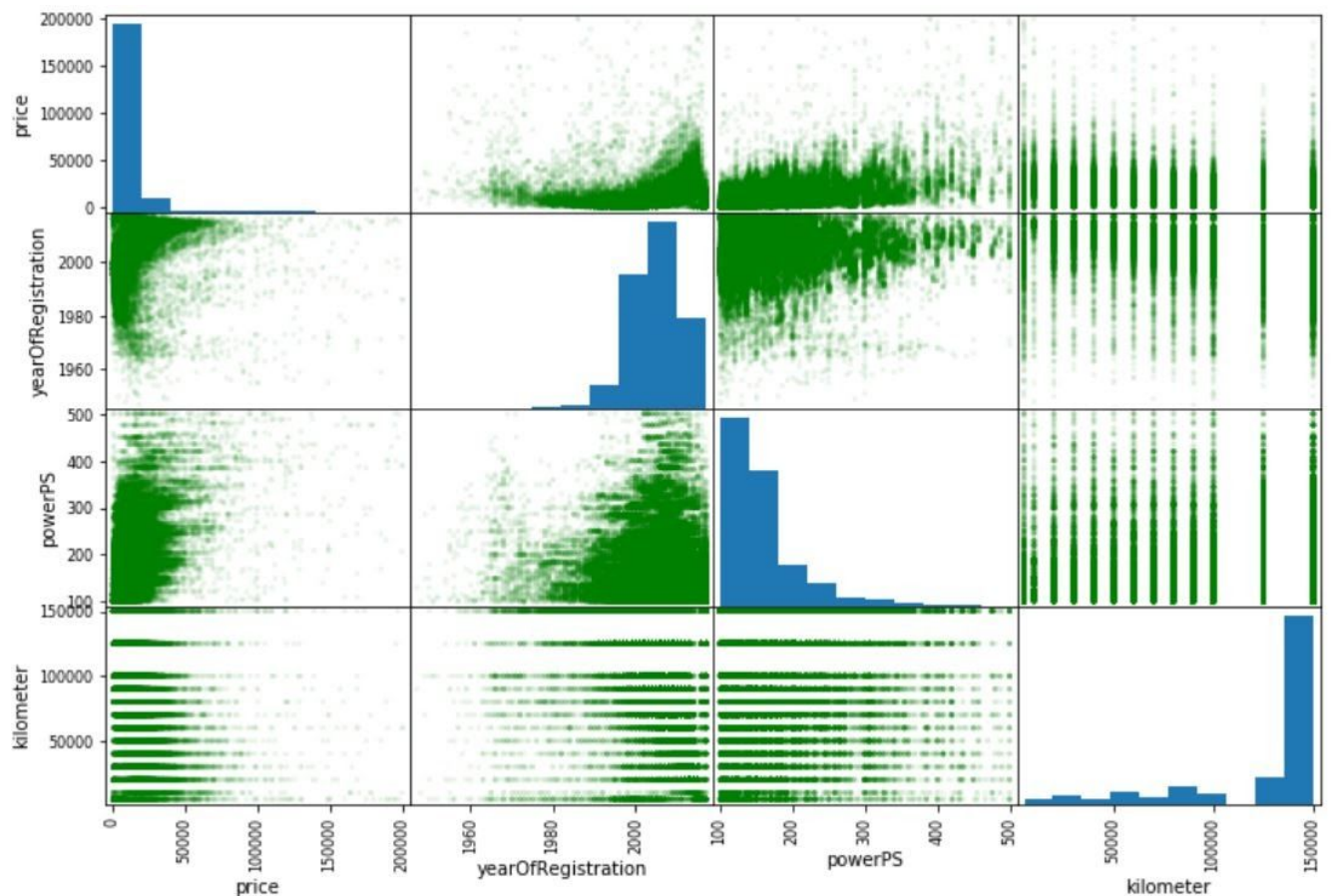
```
In [6]: 1 cat_val = ["seller", "offerType", "abtest", "gearbox", "fuelType", "notRepairedDamage", "nrOfPictures"]
2
3 for col in cat_val:
4     print ([col], ":", data_frame[col].unique())
```

```
['seller'] : ['privat' 'gewerblich']
['offerType'] : ['Angebot' 'Gesuch']
['abtest'] : ['test' 'control']
['gearbox'] : ['manuell' 'automatik' nan]
['fuelType'] : ['benzin' 'diesel' nan 'lpg' 'andere' 'hybrid' 'cng' 'elektro']
['notRepairedDamage'] : [nan 'ja' 'nein']
['nrOfPictures'] : [0]
```

```
In [7]: 1 cars_cl = data_frame.copy()
2 cars_cl = cars_cl [
3
4     (cars_cl["yearOfRegistration"].between(1945,2017,inclusive=True)) &
5     (cars_cl["powerPS"].between(100,500,inclusive=True)) &
6     (cars_cl["price"].between(100,200000,inclusive=True))
7 ]
```

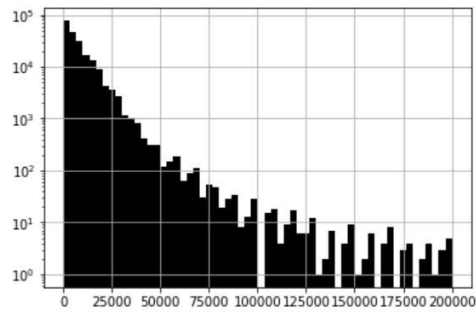
```
In [8]: 1 num_attr = ["price", "yearOfRegistration", "powerPS", "kilometer"]
2 %matplotlib inline
3 pd.plotting.scatter_matrix(cars_cl[num_attr],figsize=(12,8),alpha = 0.08,color='green')
```

```
Out[8]: array([[<matplotlib.axes._subplots.AxesSubplot object at 0x000001CEFC7C1608>,
<matplotlib.axes._subplots.AxesSubplot object at 0x000001CEF0BE8F88>,
<matplotlib.axes._subplots.AxesSubplot object at 0x000001CEF0C29048>,
<matplotlib.axes._subplots.AxesSubplot object at 0x000001CEF0C5F148>],
[<matplotlib.axes._subplots.AxesSubplot object at 0x000001CEF0C96248>,
<matplotlib.axes._subplots.AxesSubplot object at 0x000001CEF0CCF348>,
<matplotlib.axes._subplots.AxesSubplot object at 0x000001CEF0D07448>,
<matplotlib.axes._subplots.AxesSubplot object at 0x000001CEF0D40548>],
[<matplotlib.axes._subplots.AxesSubplot object at 0x000001CEF0D4C148>,
<matplotlib.axes._subplots.AxesSubplot object at 0x000001CEF0D84348>,
<matplotlib.axes._subplots.AxesSubplot object at 0x000001CEF0DE9848>,
<matplotlib.axes._subplots.AxesSubplot object at 0x000001CEF0E20988>],
[<matplotlib.axes._subplots.AxesSubplot object at 0x000001CEF0E59A88>,
<matplotlib.axes._subplots.AxesSubplot object at 0x000001CEF0E90B88>,
<matplotlib.axes._subplots.AxesSubplot object at 0x000001CEF0EC8C88>,
<matplotlib.axes._subplots.AxesSubplot object at 0x000001CEF0F02D48>]],
dtype=object)
```



```
In [9]: 1 cars_cl["price"].hist(bins = 60, log = True,color='black')
```

```
Out[9]: <matplotlib.axes_subplots.AxesSubplot at 0x1cef0646548>
```



```
In [10]: 1 cars_clean = data_frame.copy()
2
3 # Filtering the irrelevant data.
4
5 cars_clean = cars_clean[
6     (cars_clean["yearOfRegistration"].between(1945, 2017, inclusive=True)) &
7     (cars_clean["powerPS"].between(100, 500, inclusive=True)) &
8     (cars_clean["price"].between(100, 200000, inclusive=True)) &
9     (cars_clean["offerType"] == "Angebot")
10 ]
11
12 # Replacing the empty/ NaN values in the data.
13
14 cars_clean['vehicleType'].fillna(value='blank', inplace=True)
15 cars_clean['gearbox'].fillna(value='blank', inplace=True)
16 cars_clean['model'].fillna(value='blank', inplace=True)
17 cars_clean['fuelType'].fillna(value='blank', inplace=True)
18 cars_clean['notRepairedDamage'].fillna(value='blank', inplace=True)
19
20 # Change categorical attributes dtype to category
21
22 for col in cars_clean:
23     if cars_clean[col].dtype == "object":
24         cars_clean[col] = cars_clean[col].astype('category')
25
26 # Assign codes to categorical attributes instead of strings
27
28 cat_columns = cars_clean.select_dtypes(['category']).columns
29
30 cars_clean[cat_columns] = cars_clean[cat_columns].apply(lambda x: x.cat.codes)
31
32
33 # Dropping useless columns
34
35 drop_cols = ["dateCrawled", "abtest", "dateCreated", "nrOfPictures", "lastSeen"]
36
37 cars_clean = cars_clean.drop(drop_cols, axis=1)
```



```
In [11]: 1 cars_clean.head()
```

Out[11]:

	name	seller	offerType	price	vehicleType	yearOfRegistration	gearbox	powerPS	model	kilometer	monthOfRegistration	fuelType	brand	notRepairedDa
1	2136	1	0	18300	4	2011	2	190	50	125000	5	4	1	
2	56312	1	0	9800	8	2004	0	163	115	125000	8	4	14	
5	20652	1	0	650	7	1995	2	102	11	150000	10	1	2	
6	87933	1	0	2200	3	2004	2	109	8	150000	8	1	25	
8	44386	1	0	14500	2	2014	2	125	59	30000	8	1	10	

```
In [12]: 1 # Getting the train and test sets
2 train_set, test_set = train_test_split(cars_clean, test_size = 0.2, random_state = 42)
3
4 # Separation of Features and Labels
5
6 cars_price = train_set["price"].copy()
7 cars = train_set.drop("price", axis=1)
8 #cars_price
```

```
In [13]: 1
2 # Create a class to select numerical or categorical columns
3 # since Scikit-Learn doesn't handle DataFrames yet
4 class DFSelector(BaseEstimator, TransformerMixin):
5     def __init__(self, attribute_names):
6         self.attribute_names = attribute_names
7     def fit(self, X, y=None):
8         return self
9     def transform(self, X):
10        return X[self.attribute_names].values
```

```
In [14]: 1 # Setting categorical and numerical attributes
2
3 cat_attribs = ["name", "seller", "offerType", "vehicleType", "fuelType", "brand", "notRepairedDamage"]
4 num_attribs = list(cars.drop(cat_attribs, axis=1))
5
6 # Building the Pipelines
7
8 num_pipeline = Pipeline([
9     ("selector", DFSelector(num_attribs)),
10    ("std_scaler", StandardScaler())
11 ])
12
13 cat_pipeline = Pipeline([
14     ("selector", DFSelector(cat_attribs)),
15     ("encoder", OneHotEncoder(sparse=True))
16 ])
17
18 full_pipeline = FeatureUnion(transformer_list=[
19     ("num_pipeline", num_pipeline),
20     ("cat_pipeline", cat_pipeline)
21 ])
22
```

```
In [15]: 1 cars_prepared = full_pipeline.fit_transform(cars)
```

```
In [16]: 1 lin_reg = LinearRegression()
2 lin_reg.fit(cars_prepared, cars_price)
```


Out[16]: LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)

```
In [17]: 1 from sklearn.metrics import mean_squared_error
```

```
In [19]: 1 cars_predictions = lin_reg.predict(cars_prepared)
2 lin_mse = mean_squared_error(cars_price, cars_predictions)
3 lin_rmse = np.sqrt(lin_mse)
4 lin_rmse
```

Out[19]: 1827.3421981985368

```
In [20]: 1 cars_predictions[0:4]
```

Out[20]: array([25749.93597856, 4000.11431413, 1118.64918339, 2599.98254437])

```
In [21]: 1 list(cars_price[0:4])
```

Out[21]: [25750, 4000, 800, 2600]

```
In [23]: 1 # Decision Tree Approach
2 from sklearn.tree import DecisionTreeRegressor
3
4 tree_reg = DecisionTreeRegressor(random_state=42)
5 tree_reg.fit(cars_prepared, cars_price)
```

Out[23]: 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, presort='deprecated', random_state=42, splitter='best')

```
In [:]: 1 cars_predictions = tree_reg.predict(cars_prepared)
2 tree_mse = mean_squared_error(cars_price, cars_predictions)
3 tree_rmse = np.sqrt(tree_mse)
4 tree_rmse
```

```
In [*]: 1 # Random Forest Approach
2 from sklearn.ensemble import RandomForestRegressor
3
4 forest_reg = RandomForestRegressor(random_state=42, n_jobs=-1, max_depth=30)
5 forest_reg.fit(cars_prepared, cars_price)
```

```
In [*]: 1 cars_predictions = forest_reg.predict(cars_prepared)
2 forest_mse = mean_squared_error(cars_price, cars_predictions)
3 forest_rmse = np.sqrt(forest_mse)
4 forest_rmse
```

```
In [*]: 1 # Cross Validation
2
3 from sklearn.model_selection import cross_val_score
4
5 def display_scores(scores):
6     print("Scores:", scores)
7     print("Mean:", scores.mean())
8     print("Standard deviation:", scores.std())
```

```
In [*]: 1 # LinReg - CrossValidation
2 scores = cross_val_score(lin_reg, cars_prepared, cars_price,
3                           scoring="neg_mean_squared_error", cv=4)
4 lin_rmse_scores = np.sqrt(-scores)
5
6 display_scores(lin_rmse_scores)
```

```
In [*]: 1 # Decision Tree - CrossValidation
2
3 scores = cross_val_score(tree_reg, cars_prepared, cars_price,
4                           scoring="neg_mean_squared_error", cv=4)
5 tree_rmse_scores = np.sqrt(-scores)
6
7 display_scores(tree_rmse_scores)
```

```

In [*]: 1 #Random Forest- CrossValidation
2
3 from sklearn.model_selection import cross_val_score
4
5 scores = cross_val_score(forest_reg, cars_prepared, cars_price,
6                          scoring="neg_mean_squared_error", cv=2)
7 forest_rmse_scores = np.sqrt(-scores)
8
9 display_scores(forest_rmse_scores)

In [*]: 1 # To see the importance of the features
2 feature_importances = forest_reg.feature_importances_
3 feature_importances

In [*]: 1 cat_encoder = cat_pipeline.named_steps["encoder"]
2 #cat_one_hot_attribs = list(cat_encoder.categories_[0])
3 attributes = num_attribs #+ cat_encoder
4 sorted(zip(feature_importances, attributes), reverse=True)

In [*]: 1 # Final Model
2 final_model = forest_reg
3
4 cars_test = test_set.drop("price", axis = 1)
5 cars_price_test = test_set["price"].copy()
6
7 cars_test_prepared = full_pipeline.transform(cars_test) ## call transform NOT fit_transform
8
9
10 from sklearn.metrics import mean_squared_error
11 final_predictions = final_model.predict(cars_test_prepared)
12
13 final_mse = mean_squared_error(cars_price_test, final_predictions)
14
15 final_rmse = np.sqrt(final_mse)
16 #to see the results.
17 #final_predictions

```

Using test and split technique, we can say that the **Forest Regression Approach** is most accurate, but time consuming.