

Linear Regression- Vehicle selection

In [1]: *#Step1: Importing All the Required Libraries*

```
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn import preprocessing, svm
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
```

In [2]: *#Step2: Reading the Dataset*

```
df=pd.read_csv(r"C:\Users\sneha\Downloads\fiat500_VehicleSelection_Dataset (2)
df
```

Out[2]:

	ID	model	engine_power	age_in_days	km	previous_owners	lat	lon
0	1	lounge	51	882	25000	1	44.907242	8.611560
1	2	pop	51	1186	32500	1	45.666359	12.241890
2	3	sport	74	4658	142228	1	45.503300	11.417840
3	4	lounge	51	2739	160000	1	40.633171	17.634609
4	5	pop	73	3074	106880	1	41.903221	12.495650
...
1533	1534	sport	51	3712	115280	1	45.069679	7.704920
1534	1535	lounge	74	3835	112000	1	45.845692	8.666870
1535	1536	pop	51	2223	60457	1	45.481541	9.413480
1536	1537	lounge	51	2557	80750	1	45.000702	7.682270
1537	1538	pop	51	1766	54276	1	40.323410	17.568270

1538 rows × 9 columns



In [3]:

```
df=df[['age_in_days','price']]
df.columns=['age_in_days','price']
```

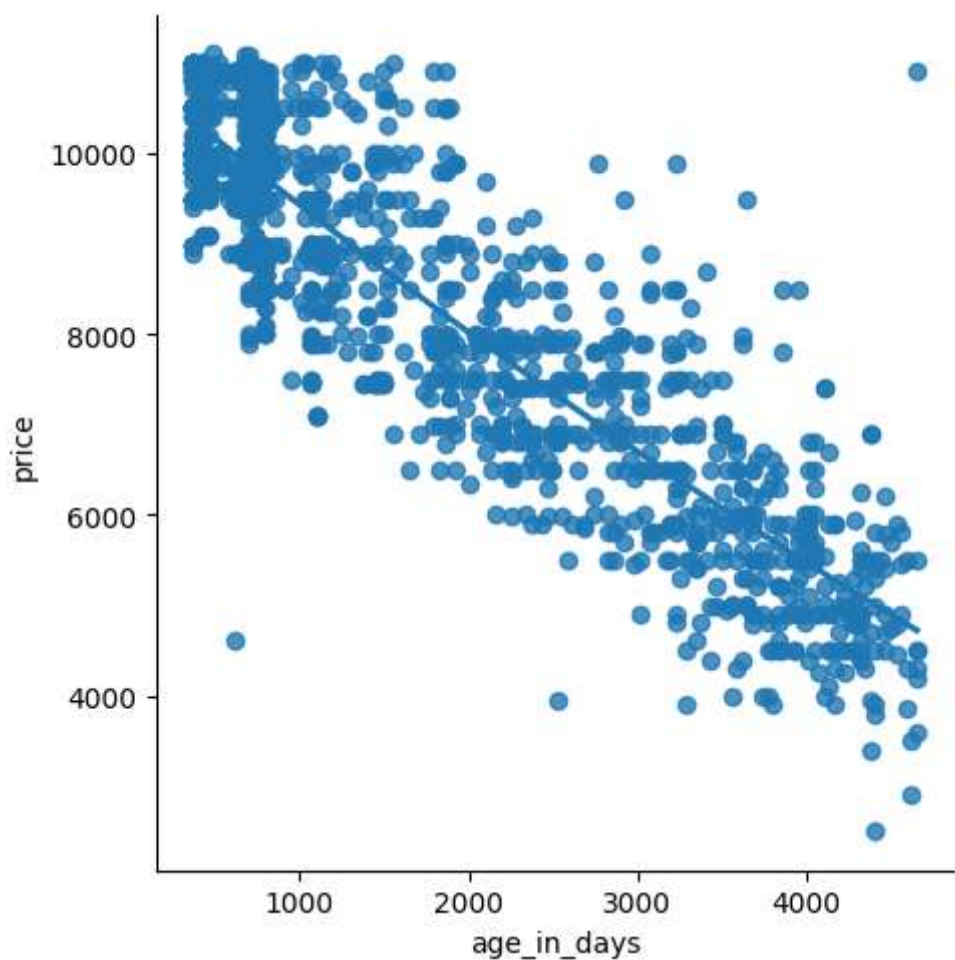
```
In [4]: df.head(10)
```

Out[4]:

	age_in_days	price
0	882	8900
1	1186	8800
2	4658	4200
3	2739	6000
4	3074	5700
5	3623	7900
6	731	10750
7	1521	9190
8	4049	5600
9	3653	6000

```
In [5]: #step3: Exploring the Data Scatter . Plotting the data scatter  
sns.lmplot(x="age_in_days", y="price", data=df, order=2, ci=None)
```

Out[5]: <seaborn.axisgrid.FacetGrid at 0x1d50858dba0>



In [6]: `df.describe()`

Out[6]:

	age_in_days	price
count	1538.000000	1538.000000
mean	1650.980494	8576.003901
std	1289.522278	1939.958641
min	366.000000	2500.000000
25%	670.000000	7122.500000
50%	1035.000000	9000.000000
75%	2616.000000	10000.000000
max	4658.000000	11100.000000

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

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1538 entries, 0 to 1537
Data columns (total 2 columns):
#   Column      Non-Null Count  Dtype
---  -
0   age_in_days  1538 non-null   int64
1   price       1538 non-null   int64
dtypes: int64(2)
memory usage: 24.2 KB
```

In [8]: *#Step4: Data Cleaning- Eliminating NaN or missing input numbers*

```
df.fillna(method= 'ffill', inplace = True)
```

C:\Users\sneha\AppData\Local\Temp\ipykernel_10652\1577053232.py:3: SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
df.fillna(method= 'ffill', inplace = True)
```

In [9]: *# Step-5: Training Our Model*

```
X = np.array(df['age_in_days']).reshape(-1, 1)
```

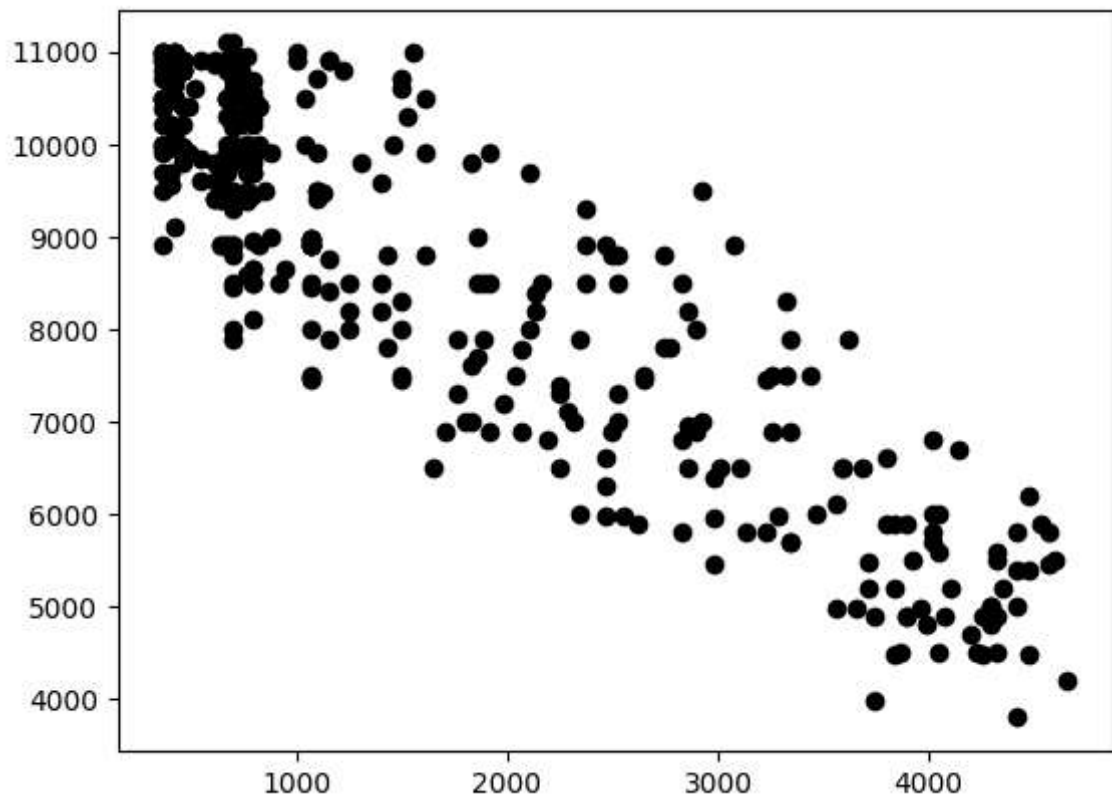
```
y = np.array(df['price']).reshape(-1, 1)
```

```
#Seperating the data into independent and dependent variables and convert
#Now each dataset contains only one column
```

```
In [10]: X_train,X_test,y_train,y_test = train_test_split(X, y, test_size = 0.25)
# Splitting the data into training data and test data
regr = LinearRegression()
regr.fit(X_train, y_train)
print(regr.score(X_test, y_test))
```

0.8017595752654839

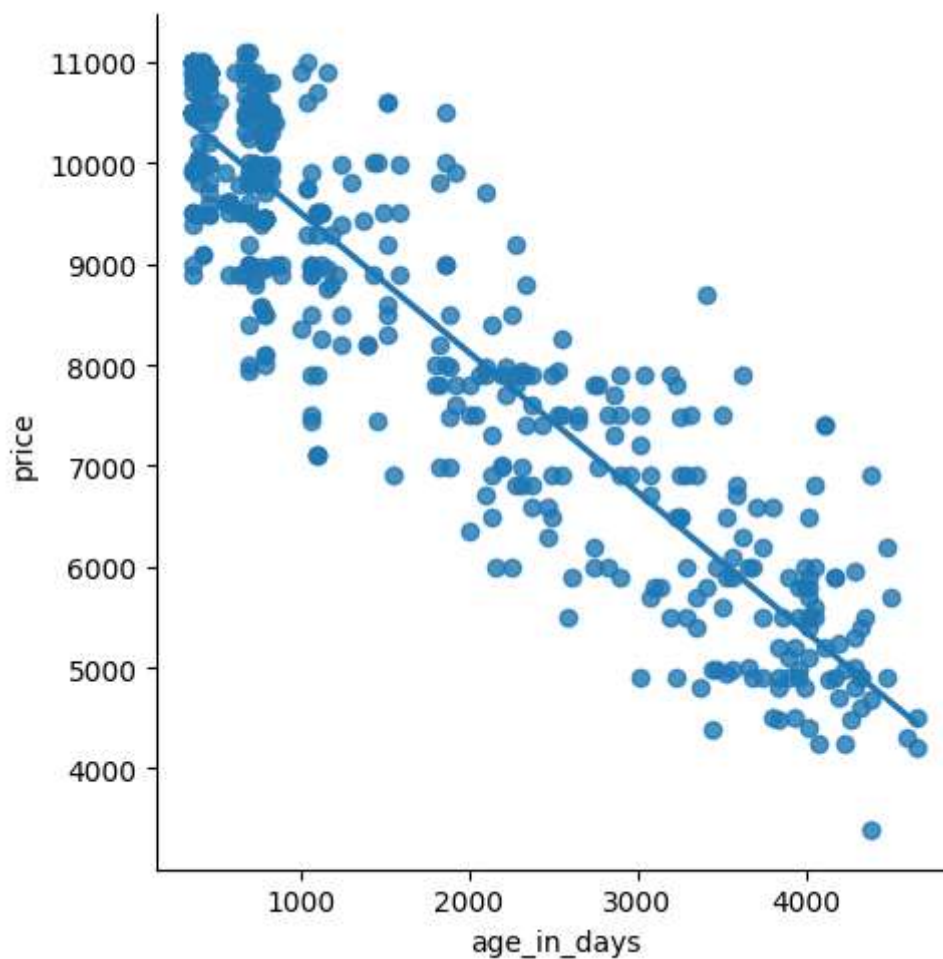
```
In [11]: #step-6: Exploring Our Results
y_pred = regr.predict(X_test)
plt.scatter(X_test, y_test, color = 'k')
plt.show()
# Data scatter of predicted values
```



```
In [12]: #Step7: Working with a Smaller dataset
df500 = df[:][:500]

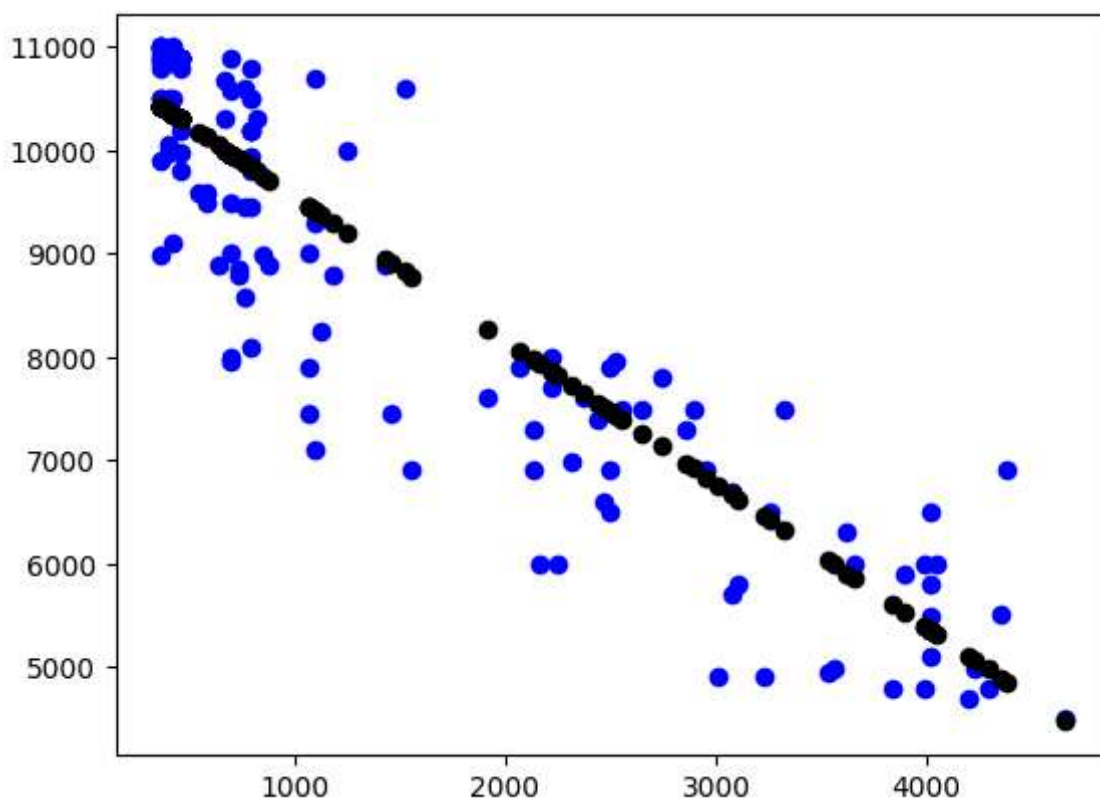
#Selecting the 1st 500 rowss of the data
sns.lmplot(x= "age_in_days", y="price", data= df500, order = 1, ci = None)
```

Out[12]: <seaborn.axisgrid.FacetGrid at 0x1d508c9b520>



```
In [13]: df500.fillna(method = 'ffill', inplace = True)
X = np.array(df500['age_in_days']).reshape(-1, 1)
y = np.array(df500['price']).reshape(-1, 1)
df500.dropna(inplace = True)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25)
regr = LinearRegression()
regr.fit(X_train, y_train)
print("Regression:", regr.score(X_test, y_test))
y_pred = regr.predict(X_test)
plt.scatter(X_test, y_test, color = 'b')
plt.scatter(X_test, y_pred, color = 'k')
plt.show()
```

Regression: 0.8154849953038915



```
In [14]: #Step-8: Evaluation of model

from sklearn.linear_model import LinearRegression

from sklearn.metrics import r2_score

#Train the model

model = LinearRegression()

model.fit(X_train, y_train)

#Evaluating the model on the test set

y_pred = model.predict(X_test)

r2 = r2_score(y_test, y_pred)

print("R2 score:",r2)
```

R2 score: 0.8154849953038915

#Step-9:Conclusion:

Dataset we have taken is poor for Linear Model, but with the smaller data works well with Linear Model.

In []: