

```
In [1]: 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 [4]: df=pd.read_csv(r"C:\Users\ubinl\Downloads\fiat500_VehicleSelection_Dataset.csv")
df
```

Out[4]:

| | 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



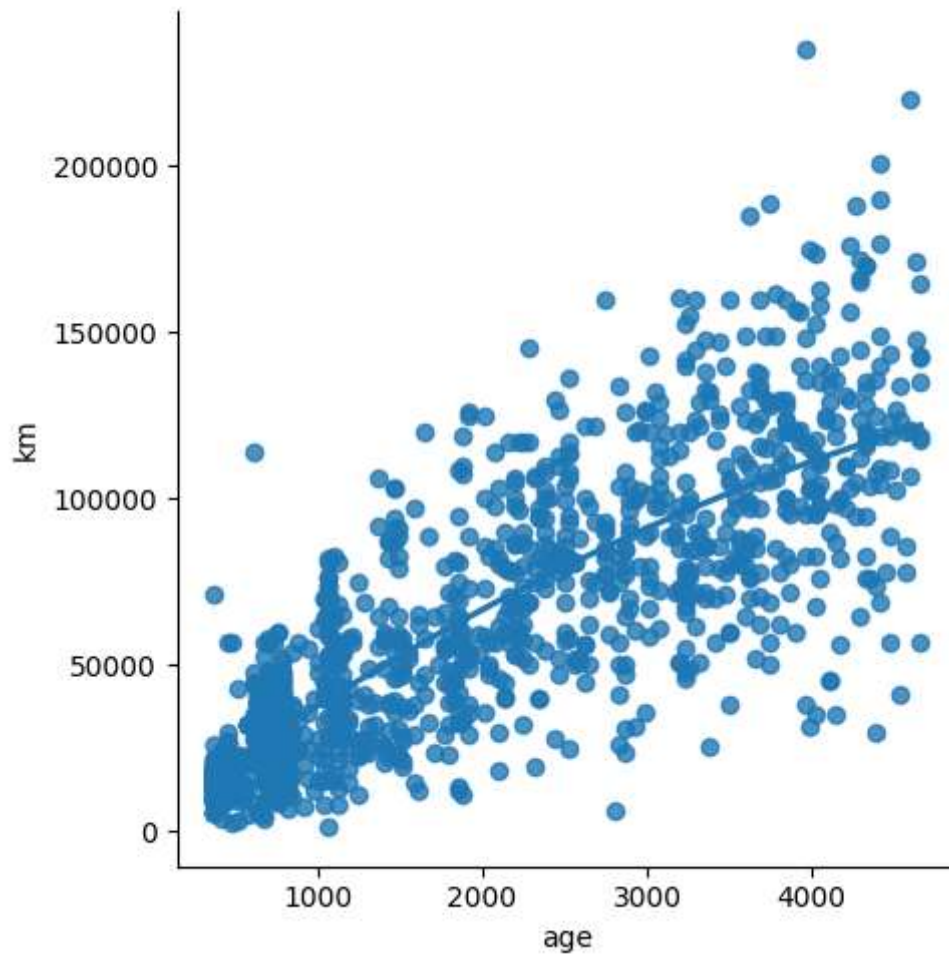
```
In [5]: df=df[['age_in_days','km']]  
df.columns=['age','km']  
df.head(10)
```

Out[5]:

| | age | km |
|---|------|--------|
| 0 | 882 | 25000 |
| 1 | 1186 | 32500 |
| 2 | 4658 | 142228 |
| 3 | 2739 | 160000 |
| 4 | 3074 | 106880 |
| 5 | 3623 | 70225 |
| 6 | 731 | 11600 |
| 7 | 1521 | 49076 |
| 8 | 4049 | 76000 |
| 9 | 3653 | 89000 |

```
In [6]: #step 3:exploring
sns.lmplot(x="age",y="km",data=df,order=2,ci=None)
```

```
Out[6]: <seaborn.axisgrid.FacetGrid at 0x22a07ca1f90>
```



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

```
Out[7]:
```

| | age | km |
|-------|-------------|---------------|
| count | 1538.000000 | 1538.000000 |
| mean | 1650.980494 | 53396.011704 |
| std | 1289.522278 | 40046.830723 |
| min | 366.000000 | 1232.000000 |
| 25% | 670.000000 | 20006.250000 |
| 50% | 1035.000000 | 39031.000000 |
| 75% | 2616.000000 | 79667.750000 |
| max | 4658.000000 | 235000.000000 |

In [8]: 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      1538 non-null    int64
 1   km       1538 non-null    int64
dtypes: int64(2)
memory usage: 24.2 KB
```

In [10]: *#step 4:*
df.fillna(method='ffill',inplace=True)

C:\Users\ubini\AppData\Local\Temp\ipykernel_35812\3632936489.py:2: 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)

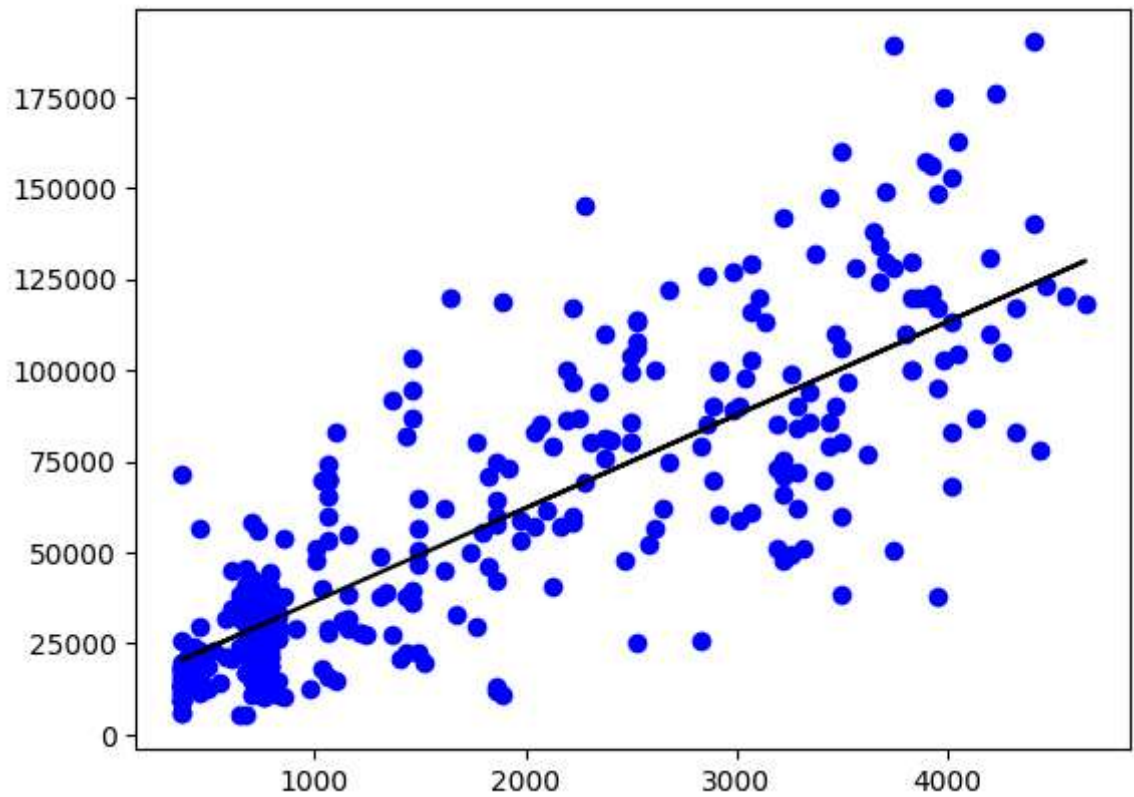
In [11]: *#step 5:training model*
x=np.array(df['age']).reshape(-1,1)
y=np.array(df['km']).reshape(-1,1)
#seperating
#column
df.dropna(inplace=True)
#dropping values
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.25)
#splitting data
regr=LinearRegression()
regr.fit(x_train,y_train)
print(regr.score(x_test,y_test))

0.7037357369026189

C:\Users\ubini\AppData\Local\Temp\ipykernel_35812\3196656996.py:6: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

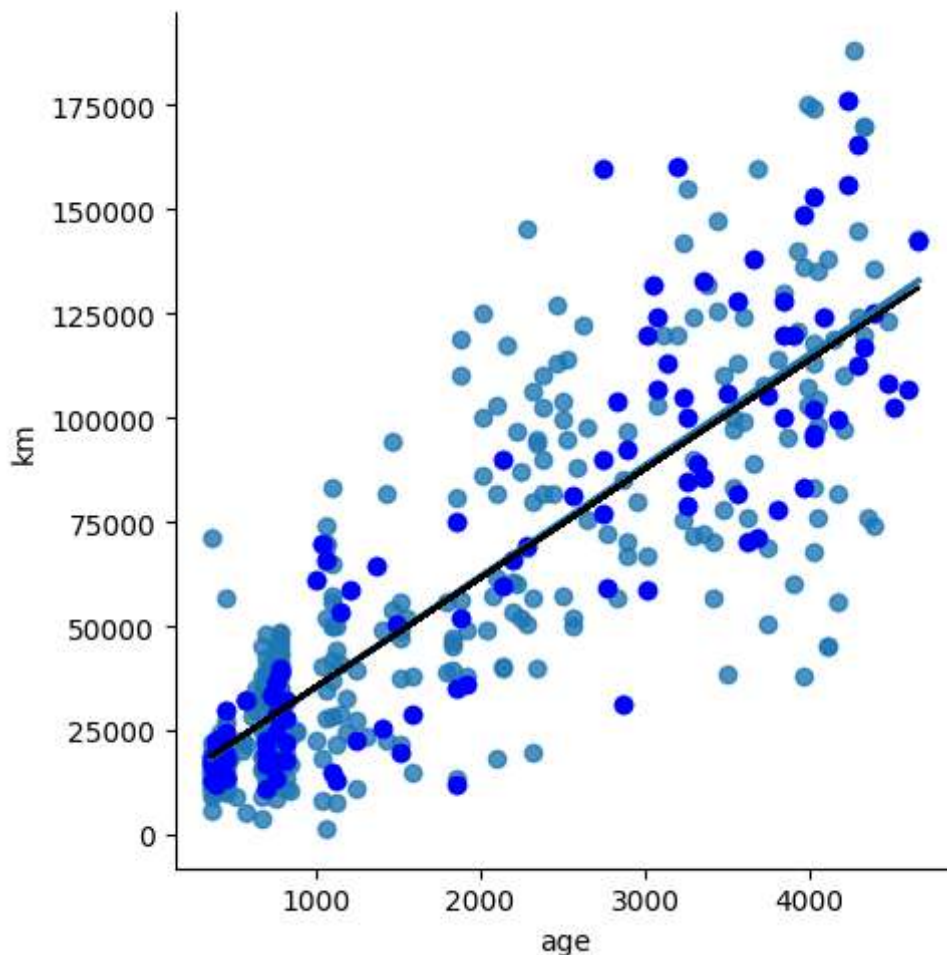
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)
df.dropna(inplace=True)

```
In [12]: #step 6:exploring results
y_pred=regr.predict(x_test)
plt.scatter(x_test,y_test,color='b')
plt.plot(x_test,y_pred,color='k')
plt.show()
#scatter
```



```
In [13]: #step7:working with a smaller data set
df500=df[:][:500]
#selecting
sns.lmplot(x="age",y="km",data=df500,order=1,ci=None)
df500.fillna(method='ffill',inplace=True)
x=np.array(df500['age']).reshape(-1,1)
y=np.array(df500['km']).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.plot(x_test,y_pred,color='k')
plt.show()
```

Regression: 0.7872639373865098



```
In [14]: #step 8:
        from sklearn.linear_model import LinearRegression
        from sklearn.metrics import r2_score
        #train
        model=LinearRegression()
        model.fit(x_train,y_train)
        #evaluate
        y_pred=model.predict(x_test)
        r2=r2_score(y_test,y_pred)
        print("r2_score:",r2)
```

r2_score: 0.7872639373865098

In []: