# problem statement:

TO find the relation between the price and the seats available e in the car.

```python
In [18]: import numpy as np
         import pandas as pd
         import matplotlib.pyplot as plt
         import seaborn as sns
         from sklearn import preprocessing,svm
         from sklearn.model_selection import train_test_split
         from sklearn.linear_model import LinearRegression
```

In [19]:
```python
df=pd.read_csv(r"C:\Users\my pc\Downloads\useD_cars_data.csv")
df
```

Out[19]:

|  | S.No. | Name | Location | Year | Kilometers_Driven | Fuel_Type | Transmission | Owner_T |
|---|---|---|---|---|---|---|---|---|
| **0** | 0 | Maruti Wagon R LXI CNG | Mumbai | 2010 | 72000 | CNG | Manual | |
| **1** | 1 | Hyundai Creta 1.6 CRDi SX Option | Pune | 2015 | 41000 | Diesel | Manual | |
| **2** | 2 | Honda Jazz V | Chennai | 2011 | 46000 | Petrol | Manual | |
| **3** | 3 | Maruti Ertiga VDI | Chennai | 2012 | 87000 | Diesel | Manual | |
| **4** | 4 | Audi A4 New 2.0 TDI Multitronic | Coimbatore | 2013 | 40670 | Diesel | Automatic | Sec |
| **...** | ... | ... | ... | ... | ... | ... | ... | |
| **7248** | 7248 | Volkswagen Vento Diesel Trendline | Hyderabad | 2011 | 89411 | Diesel | Manual | |
| **7249** | 7249 | Volkswagen Polo GT TSI | Mumbai | 2015 | 59000 | Petrol | Automatic | |
| **7250** | 7250 | Nissan Micra Diesel XV | Kolkata | 2012 | 28000 | Diesel | Manual | |
| **7251** | 7251 | Volkswagen Polo GT TSI | Pune | 2013 | 52262 | Petrol | Automatic | T |
| **7252** | 7252 | Mercedes-Benz E-Class 2009-2013 E 220 CDI Avan... | Kochi | 2014 | 72443 | Diesel | Automatic | |

7253 rows × 14 columns

In [20]: `df.head()`

Out[20]:

| | S.No. | Name | Location | Year | Kilometers_Driven | Fuel_Type | Transmission | Owner_Type |
|---|---|---|---|---|---|---|---|---|
| **0** | 0 | Maruti Wagon R LXI CNG | Mumbai | 2010 | 72000 | CNG | Manual | First |
| **1** | 1 | Hyundai Creta 1.6 CRDi SX Option | Pune | 2015 | 41000 | Diesel | Manual | First |
| **2** | 2 | Honda Jazz V | Chennai | 2011 | 46000 | Petrol | Manual | First |
| **3** | 3 | Maruti Ertiga VDI | Chennai | 2012 | 87000 | Diesel | Manual | First |
| **4** | 4 | Audi A4 New 2.0 TDI Multitronic | Coimbatore | 2013 | 40670 | Diesel | Automatic | Second |

In [21]: `df.tail()`

Out[21]:

| | S.No. | Name | Location | Year | Kilometers_Driven | Fuel_Type | Transmission | Owner_Ty |
|---|---|---|---|---|---|---|---|---|
| **7248** | 7248 | Volkswagen Vento Diesel Trendline | Hyderabad | 2011 | 89411 | Diesel | Manual | F |
| **7249** | 7249 | Volkswagen Polo GT TSI | Mumbai | 2015 | 59000 | Petrol | Automatic | F |
| **7250** | 7250 | Nissan Micra Diesel XV | Kolkata | 2012 | 28000 | Diesel | Manual | F |
| **7251** | 7251 | Volkswagen Polo GT TSI | Pune | 2013 | 52262 | Petrol | Automatic | T |
| **7252** | 7252 | Mercedes-Benz E-Class 2009-2013 E 220 CDI Avan... | Kochi | 2014 | 72443 | Diesel | Automatic | F |

In [22]: `df.size`

Out[22]: 101542

In [23]: `df.shape`

Out[23]: (7253, 14)

In [25]: `df.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 7253 entries, 0 to 7252
Data columns (total 14 columns):
 #   Column             Non-Null Count  Dtype
---  ------             --------------  -----
 0   S.No.              7253 non-null   int64
 1   Name               7253 non-null   object
 2   Location           7253 non-null   object
 3   Year               7253 non-null   int64
 4   Kilometers_Driven  7253 non-null   int64
 5   Fuel_Type          7253 non-null   object
 6   Transmission       7253 non-null   object
 7   Owner_Type         7253 non-null   object
 8   Mileage            7251 non-null   object
 9   Engine             7207 non-null   object
 10  Power              7207 non-null   object
 11  Seats              7200 non-null   float64
 12  New_Price          1006 non-null   object
 13  Price              6019 non-null   float64
dtypes: float64(2), int64(3), object(9)
memory usage: 793.4+ KB
```

In [26]: `df.describe()`

Out[26]:

|  | S.No. | Year | Kilometers_Driven | Seats | Price |
|---|---|---|---|---|---|
| count | 7253.000000 | 7253.000000 | 7.253000e+03 | 7200.000000 | 6019.000000 |
| mean | 3626.000000 | 2013.365366 | 5.869906e+04 | 5.279722 | 9.479468 |
| std | 2093.905084 | 3.254421 | 8.442772e+04 | 0.811660 | 11.187917 |
| min | 0.000000 | 1996.000000 | 1.710000e+02 | 0.000000 | 0.440000 |
| 25% | 1813.000000 | 2011.000000 | 3.400000e+04 | 5.000000 | 3.500000 |
| 50% | 3626.000000 | 2014.000000 | 5.341600e+04 | 5.000000 | 5.640000 |
| 75% | 5439.000000 | 2016.000000 | 7.300000e+04 | 5.000000 | 9.950000 |
| max | 7252.000000 | 2019.000000 | 6.500000e+06 | 10.000000 | 160.000000 |

In [27]: `df.isna().any()`

Out[27]:
```
S.No.                 False
Name                  False
Location              False
Year                  False
Kilometers_Driven     False
Fuel_Type             False
Transmission          False
Owner_Type            False
Mileage                True
Engine                 True
Power                  True
Seats                  True
New_Price              True
Price                  True
dtype: bool
```

In [28]: `df.isnull().sum()`

Out[28]:
```
S.No.                    0
Name                     0
Location                 0
Year                     0
Kilometers_Driven        0
Fuel_Type                0
Transmission             0
Owner_Type               0
Mileage                  2
Engine                  46
Power                   46
Seats                   53
New_Price             6247
Price                 1234
dtype: int64
```

In [29]: `df=df[['Seats','Price']]`
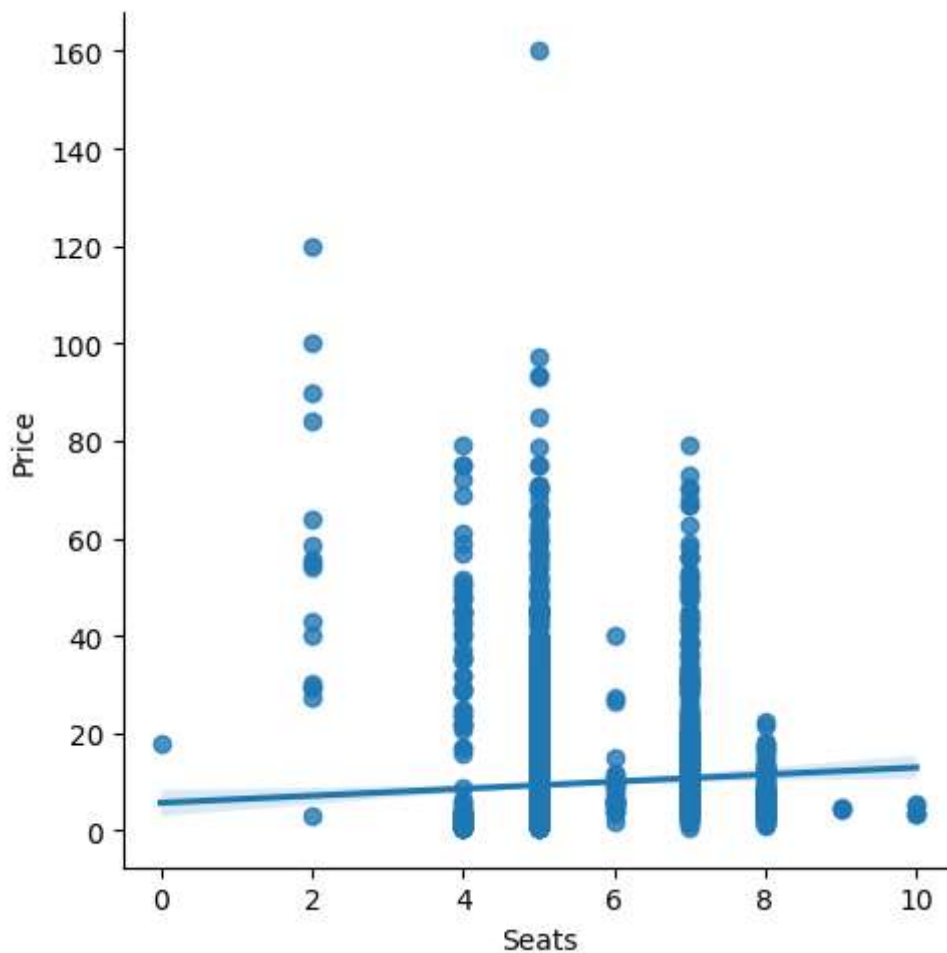
In [30]:
```python
df.columns=['Seats','Price']
df
```

Out[30]:

|  | Seats | Price |
|---|---|---|
| **0** | 5.0 | 1.75 |
| **1** | 5.0 | 12.50 |
| **2** | 5.0 | 4.50 |
| **3** | 7.0 | 6.00 |
| **4** | 5.0 | 17.74 |
| **...** | ... | ... |
| **7248** | 5.0 | NaN |
| **7249** | 5.0 | NaN |
| **7250** | 5.0 | NaN |
| **7251** | 5.0 | NaN |
| **7252** | 5.0 | NaN |

7253 rows × 2 columns

In [31]: `sns.lmplot(x='Seats',y='Price',data=df,order=1)`

Out[31]: `<seaborn.axisgrid.FacetGrid at 0x1da270d3820>`



In [32]: `df.fillna(method='ffill',inplace=True)`

```
C:\Users\my pc\AppData\Local\Temp\ipykernel_11652\4116506308.py:1: SettingWit
hCopyWarning:
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/s
table/user_guide/indexing.html#returning-a-view-versus-a-copy (https://panda
s.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-ver
sus-a-copy)
  df.fillna(method='ffill',inplace=True)
```
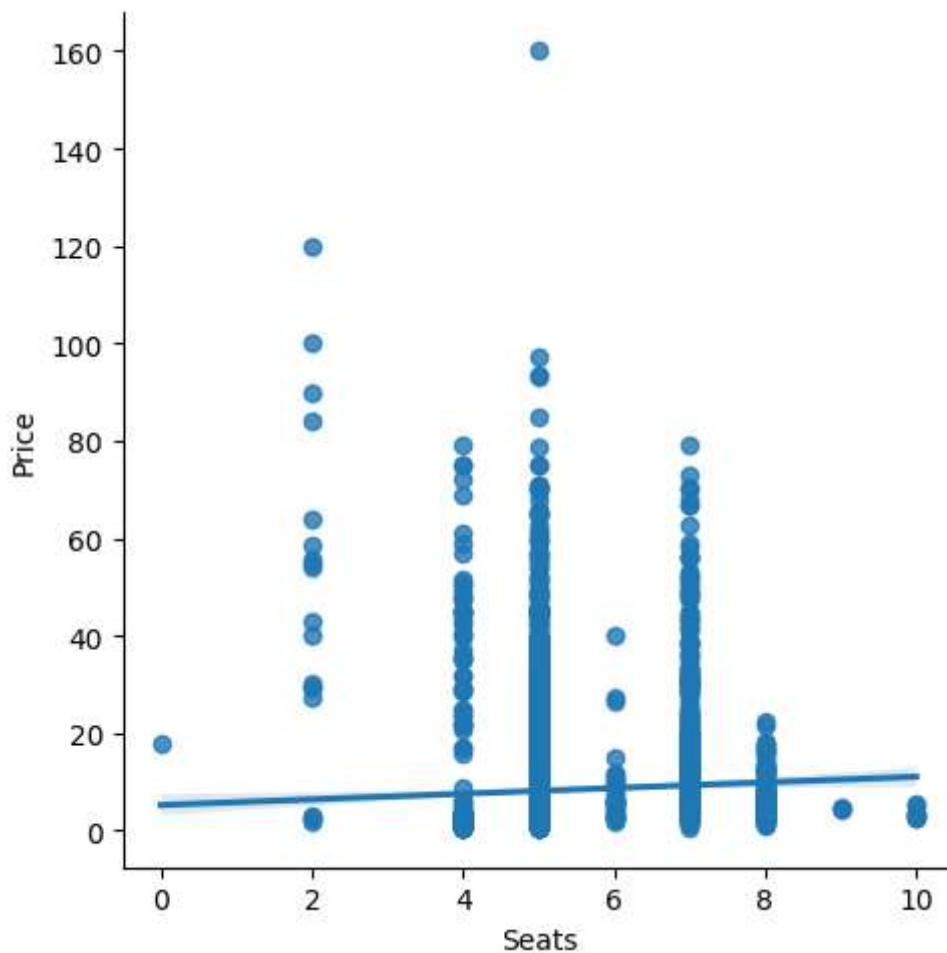
In [33]: `df.isna().any()`

Out[33]:
```
Seats    False
Price    False
dtype: bool
```

In [34]: `sns.lmplot(x='Seats',y='Price',data=df,order=1)`

Out[34]: `<seaborn.axisgrid.FacetGrid at 0x1da234bf340>`



In [35]:
```python
#step 5: Training the model
x=np.array(df['Seats']).reshape(-1,1)
y=np.array(df['Price']).reshape(-1,1)
df.dropna(inplace=True) #Dropping any rows with null values
```

```
C:\Users\my pc\AppData\Local\Temp\ipykernel_11652\1929999587.py:4: SettingWit
hCopyWarning:
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/s
table/user_guide/indexing.html#returning-a-view-versus-a-copy (https://panda
s.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-ver
sus-a-copy)
  df.dropna(inplace=True) #Dropping any rows with null values
```
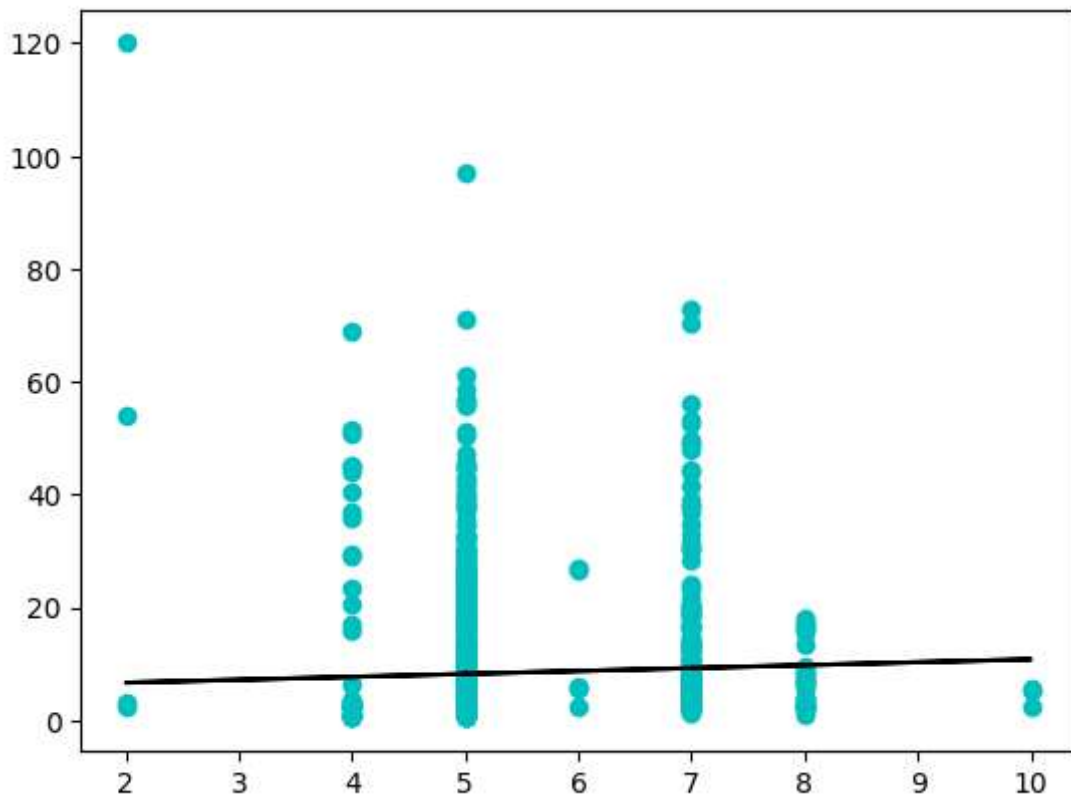
In [36]:
```python
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.25) #Spliting t
regr=LinearRegression()
regr.fit(x_train,y_train)
print(regr.score(x_test,y_test))
```

```
0.002745468949376615
```

In [37]: `#step6: Exploring our results`
```
y_pred=regr.predict(x_test)
plt.scatter(x_test,y_test,color='c')
plt.plot(x_test,y_pred,color='k')
plt.show()
```



# CONCLUSION:

The above dataset is not reliable or not suitable for LinearReg
ression. Because the data is not fitted properly
and also the r2 score is low. Therefore, it is a poor model