

Problem statement

To test which model is best fit for this dataset

In [3]:

```
1 #Importing Libraries
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 from sklearn.preprocessing import StandardScaler
7 from sklearn.model_selection import train_test_split
8 from sklearn.linear_model import LinearRegression
```

In [5]:

```
1 #Reading data
2 df=pd.read_csv(r"C:\Users\teppa\Desktop\AHISAI\rainfall in india 1901-2015.csv")
3 df
```

Out[5]:

	SUBDIVISION	YEAR	JAN	FEB	MAR	APR	MAY	JUN	JUL	AUG	SEP	OCT	NOV	DEC	ANNUAL	Jan-Feb	Mar-May	Jun-Sep	Oct-Dec
0	ANDAMAN & NICOBAR ISLANDS	1901	49.2	87.1	29.2	2.3	528.8	517.5	365.1	481.1	332.6	388.5	558.2	33.6	3373.2	136.3	560.3	1696.3	980.3
1	ANDAMAN & NICOBAR ISLANDS	1902	0.0	159.8	12.2	0.0	446.1	537.1	228.9	753.7	666.2	197.2	359.0	160.5	3520.7	159.8	458.3	2185.9	716.7
2	ANDAMAN & NICOBAR ISLANDS	1903	12.7	144.0	0.0	1.0	235.1	479.9	728.4	326.7	339.0	181.2	284.4	225.0	2957.4	156.7	236.1	1874.0	690.6
3	ANDAMAN & NICOBAR ISLANDS	1904	9.4	14.7	0.0	202.4	304.5	495.1	502.0	160.1	820.4	222.2	308.7	40.1	3079.6	24.1	506.9	1977.6	571.0
4	ANDAMAN & NICOBAR ISLANDS	1905	1.3	0.0	3.3	26.9	279.5	628.7	368.7	330.5	297.0	260.7	25.4	344.7	2566.7	1.3	309.7	1624.9	630.8
...
4111	LAKSHADWEEP	2011	5.1	2.8	3.1	85.9	107.2	153.6	350.2	254.0	255.2	117.4	184.3	14.9	1533.7	7.9	196.2	1013.0	316.6
4112	LAKSHADWEEP	2012	19.2	0.1	1.6	76.8	21.2	327.0	231.5	381.2	179.8	145.9	12.4	8.8	1405.5	19.3	99.6	1119.5	167.1
4113	LAKSHADWEEP	2013	26.2	34.4	37.5	5.3	88.3	426.2	296.4	154.4	180.0	72.8	78.1	26.7	1426.3	60.6	131.1	1057.0	177.6
4114	LAKSHADWEEP	2014	53.2	16.1	4.4	14.9	57.4	244.1	116.1	466.1	132.2	169.2	59.0	62.3	1395.0	69.3	76.7	958.5	290.5
4115	LAKSHADWEEP	2015	2.2	0.5	3.7	87.1	133.1	296.6	257.5	146.4	160.4	165.4	231.0	159.0	1642.9	2.7	223.9	860.9	555.4

4116 rows × 19 columns

2) Data cleaning and preprocessing

In [6]:

```
1 df.head()
```

Out[6]:

	SUBDIVISION	YEAR	JAN	FEB	MAR	APR	MAY	JUN	JUL	AUG	SEP	OCT	NOV	DEC	ANNUAL	Jan-Feb	Mar-May	Jun-Sep	Oct-Dec
0	ANDAMAN & NICOBAR ISLANDS	1901	49.2	87.1	29.2	2.3	528.8	517.5	365.1	481.1	332.6	388.5	558.2	33.6	3373.2	136.3	560.3	1696.3	980.3
1	ANDAMAN & NICOBAR ISLANDS	1902	0.0	159.8	12.2	0.0	446.1	537.1	228.9	753.7	666.2	197.2	359.0	160.5	3520.7	159.8	458.3	2185.9	716.7
2	ANDAMAN & NICOBAR ISLANDS	1903	12.7	144.0	0.0	1.0	235.1	479.9	728.4	326.7	339.0	181.2	284.4	225.0	2957.4	156.7	236.1	1874.0	690.6
3	ANDAMAN & NICOBAR ISLANDS	1904	9.4	14.7	0.0	202.4	304.5	495.1	502.0	160.1	820.4	222.2	308.7	40.1	3079.6	24.1	506.9	1977.6	571.0
4	ANDAMAN & NICOBAR ISLANDS	1905	1.3	0.0	3.3	26.9	279.5	628.7	368.7	330.5	297.0	260.7	25.4	344.7	2566.7	1.3	309.7	1624.9	630.8

In [7]:

```
1 df.tail()
```

Out[7]:

	SUBDIVISION	YEAR	JAN	FEB	MAR	APR	MAY	JUN	JUL	AUG	SEP	OCT	NOV	DEC	ANNUAL	Jan-Feb	Mar-May	Jun-Sep	Oct-Dec
4111	LAKSHADWEEP	2011	5.1	2.8	3.1	85.9	107.2	153.6	350.2	254.0	255.2	117.4	184.3	14.9	1533.7	7.9	196.2	1013.0	316.6
4112	LAKSHADWEEP	2012	19.2	0.1	1.6	76.8	21.2	327.0	231.5	381.2	179.8	145.9	12.4	8.8	1405.5	19.3	99.6	1119.5	167.1
4113	LAKSHADWEEP	2013	26.2	34.4	37.5	5.3	88.3	426.2	296.4	154.4	180.0	72.8	78.1	26.7	1426.3	60.6	131.1	1057.0	177.6
4114	LAKSHADWEEP	2014	53.2	16.1	4.4	14.9	57.4	244.1	116.1	466.1	132.2	169.2	59.0	62.3	1395.0	69.3	76.7	958.5	290.5
4115	LAKSHADWEEP	2015	2.2	0.5	3.7	87.1	133.1	296.6	257.5	146.4	160.4	165.4	231.0	159.0	1642.9	2.7	223.9	860.9	555.4

```
In [8]: 1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4116 entries, 0 to 4115
Data columns (total 19 columns):
#   Column      Non-Null Count  Dtype
---  -
0   SUBDIVISION  4116 non-null   object
1   YEAR         4116 non-null   int64
2   JAN          4112 non-null   float64
3   FEB          4113 non-null   float64
4   MAR          4110 non-null   float64
5   APR          4112 non-null   float64
6   MAY          4113 non-null   float64
7   JUN          4111 non-null   float64
8   JUL          4109 non-null   float64
9   AUG          4112 non-null   float64
10  SEP          4110 non-null   float64
11  OCT          4109 non-null   float64
12  NOV          4105 non-null   float64
13  DEC          4106 non-null   float64
14  ANNUAL       4090 non-null   float64
15  Jan-Feb      4110 non-null   float64
16  Mar-May      4107 non-null   float64
17  Jun-Sep      4106 non-null   float64
18  Oct-Dec      4103 non-null   float64
dtypes: float64(17), int64(1), object(1)
memory usage: 611.1+ KB
```

```
In [9]: 1 df.describe()
```

```
Out[9]:
```

	YEAR	JAN	FEB	MAR	APR	MAY	JUN	JUL	AUG	SEP	OCT	
count	4116.000000	4112.000000	4113.000000	4110.000000	4112.000000	4113.000000	4111.000000	4109.000000	4112.000000	4110.000000	4109.000000	4105.000000
mean	1958.218659	18.957320	21.805325	27.359197	43.127432	85.745417	230.234444	347.214334	290.263497	197.361922	95.507009	39.861905
std	33.140898	33.585371	35.909488	46.959424	67.831168	123.234904	234.710758	269.539667	188.770477	135.408345	99.519134	68.661905
min	1901.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.400000	0.000000	0.000000	0.100000	0.000000	0.000000
25%	1930.000000	0.600000	0.600000	1.000000	3.000000	8.600000	70.350000	175.600000	155.975000	100.525000	14.600000	0.700000
50%	1958.000000	6.000000	6.700000	7.800000	15.700000	36.600000	138.700000	284.800000	259.400000	173.900000	65.200000	9.500000
75%	1987.000000	22.200000	26.800000	31.300000	49.950000	97.200000	305.150000	418.400000	377.800000	265.800000	148.400000	46.100000
max	2015.000000	583.700000	403.500000	605.600000	595.100000	1168.600000	1609.900000	2362.800000	1664.600000	1222.000000	948.300000	648.900000

```
In [10]: 1 df.isnull().sum()
```

```
Out[10]: SUBDIVISION    0
YEAR                0
JAN                 4
FEB                 3
MAR                 6
APR                 4
MAY                 3
JUN                 5
JUL                 7
AUG                 4
SEP                 6
OCT                 7
NOV                11
DEC                10
ANNUAL             26
Jan-Feb            6
Mar-May            9
Jun-Sep           10
Oct-Dec           13
dtype: int64
```

```
In [11]: 1 df.fillna(method="ffill", inplace=True)
```

```
In [12]: 1 df.isnull().sum()
```

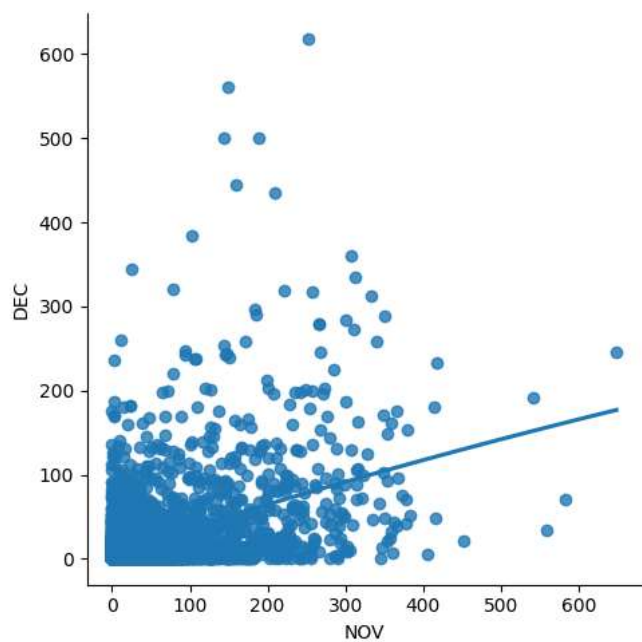
```
Out[12]: SUBDIVISION    0
YEAR                0
JAN                 0
FEB                 0
MAR                 0
APR                 0
MAY                 0
JUN                 0
JUL                 0
AUG                 0
SEP                 0
OCT                 0
NOV                 0
DEC                 0
ANNUAL              0
Jan-Feb             0
Mar-May             0
Jun-Sep             0
Oct-Dec             0
dtype: int64
```

```
In [13]: 1 df['YEAR'].value_counts()
```

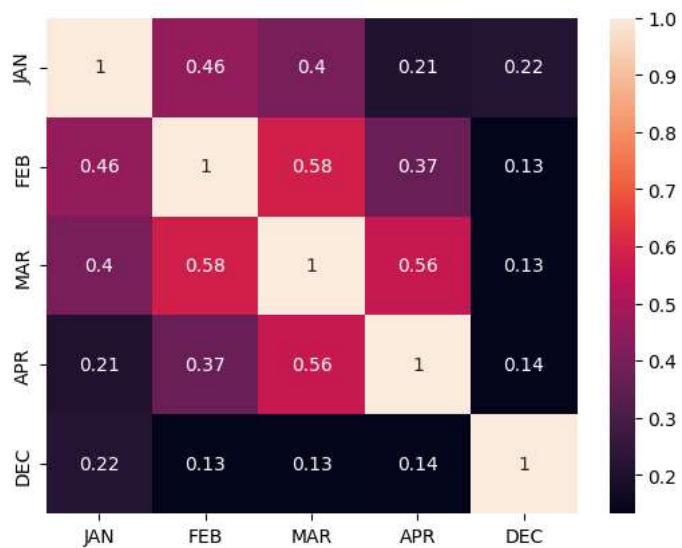
```
Out[13]: YEAR
1963     36
2002     36
1976     36
1975     36
1974     36
..
1915     35
1918     35
1954     35
1955     35
1909     34
Name: count, Length: 115, dtype: int64
```

3) Exploratory Data Analysis

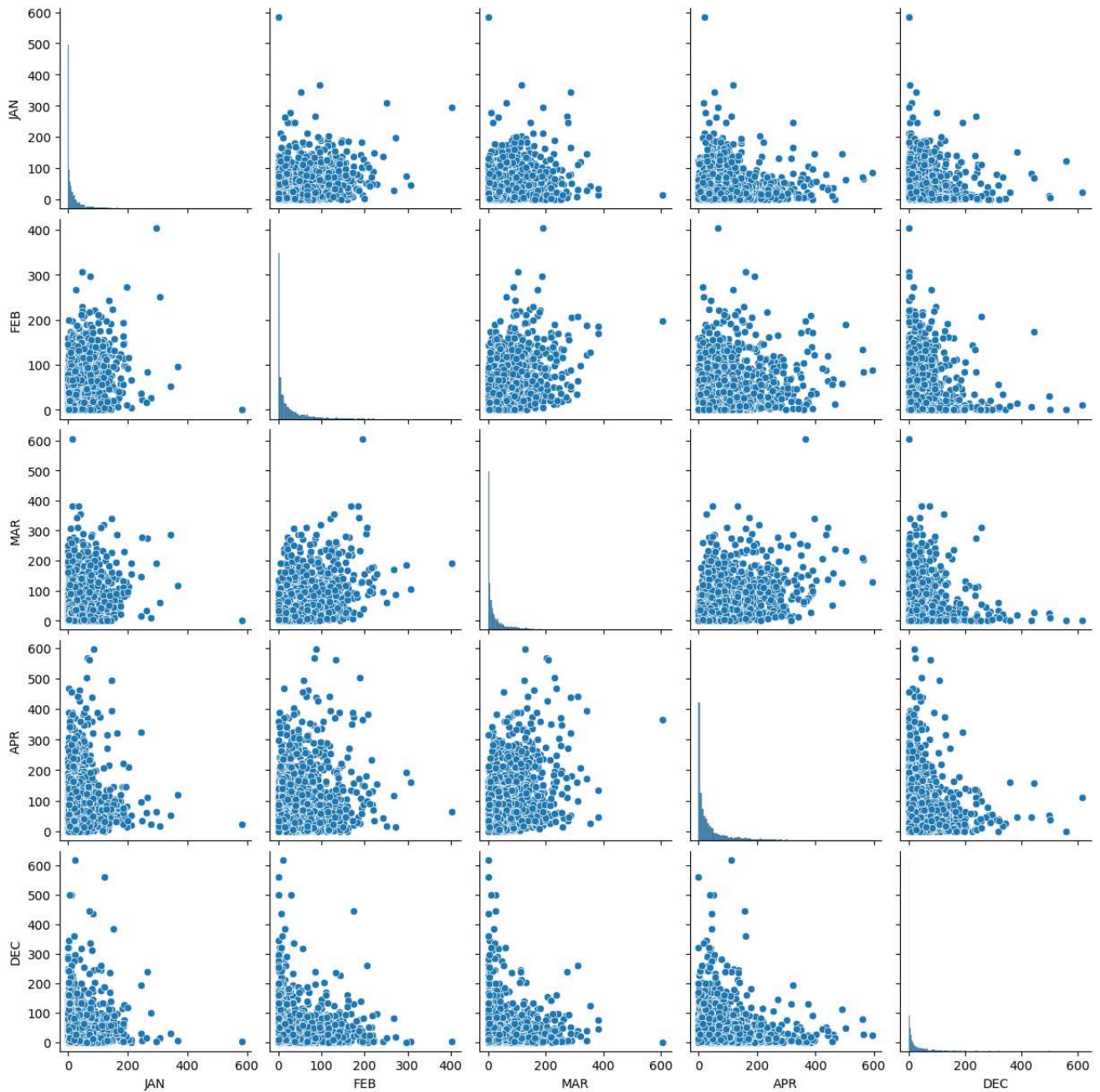
```
In [14]: 1 sns.lmplot(x='NOV',y='DEC',order=2,data=df,ci=None)
2 plt.show()
```



```
In [15]: 1 df=df[['JAN','FEB','MAR','APR','DEC']]
2 sns.heatmap(df.corr(),annot=True)
3 plt.show()
4
```



```
In [16]: 1 sns.pairplot(df)
2 plt.show()
```



4) Training our Model

```
In [17]: 1 x=np.array(df['FEB']).reshape(-1,1)
2 y=x=np.array(df['JAN']).reshape(-1,1)
```

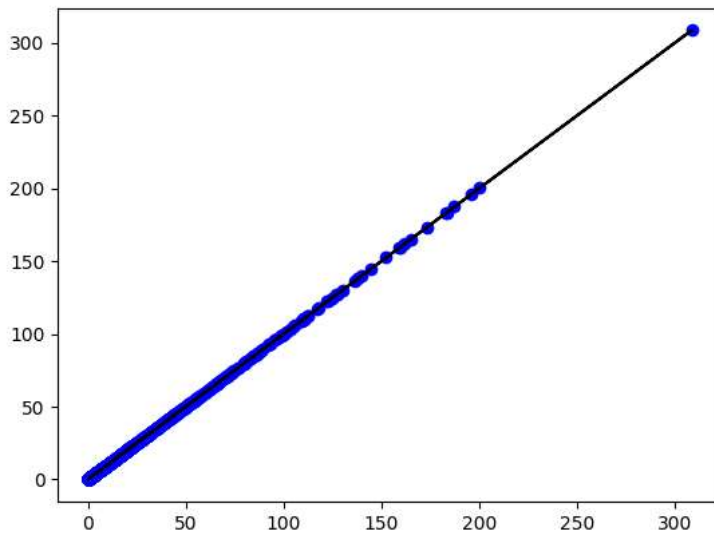
```
In [18]: 1 x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.30)
```

```
In [19]: 1 lin=LinearRegression()
2 lin.fit(x_train,y_train)
3 print(lin.score(x_train,y_train))
```

1.0

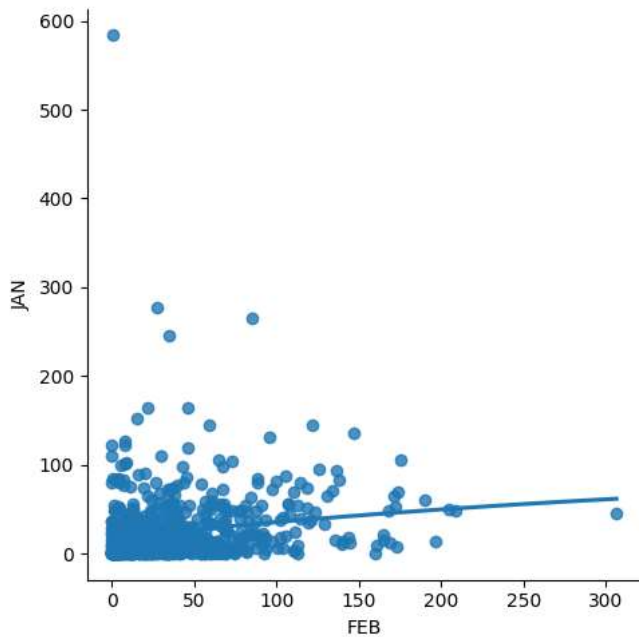
5) Exploring our Results

```
In [20]: 1 y_pred=lin.predict(x_test)
2 plt.scatter(x_test,y_test,color='blue')
3 plt.plot(x_test,y_pred,color='black')
4 plt.show()
5
```



7)Working with subset of data

```
In [21]: 1 df700=df[:700]
2 sns.lmplot(x='FEB',y='JAN',order=2,ci=None,data=df700)
3 plt.show()
```



```
In [22]: 1 df700.fillna(method='ffill',inplace=True)
2
```

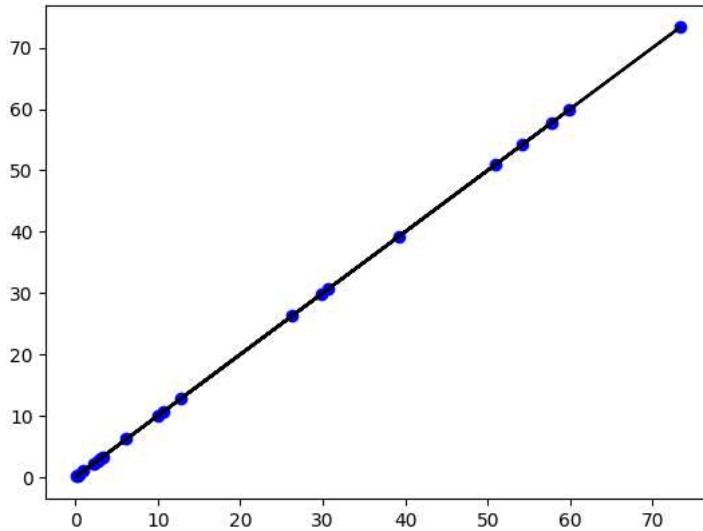
```
In [23]: 1 x=np.array(df700['FEB']).reshape(-1,1)
2 y=x*np.array(df700['JAN']).reshape(-1,1)
```

```
In [24]: 1 df700.dropna(inplace=True)
```

```
In [25]: 1 x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.03)
2 lr=LinearRegression()
3 lr.fit(x_train,y_train)
4 print(lr.score(x_test,y_test))
```

1.0

```
In [26]: 1 y_pred=lr.predict(x_test)
2 plt.scatter(x_test,y_test,color='b')
3 plt.plot(x_test,y_pred,color='k')
4 plt.show()
```



```
In [27]: 1 from sklearn.linear_model import LinearRegression
2 from sklearn.metrics import r2_score
```

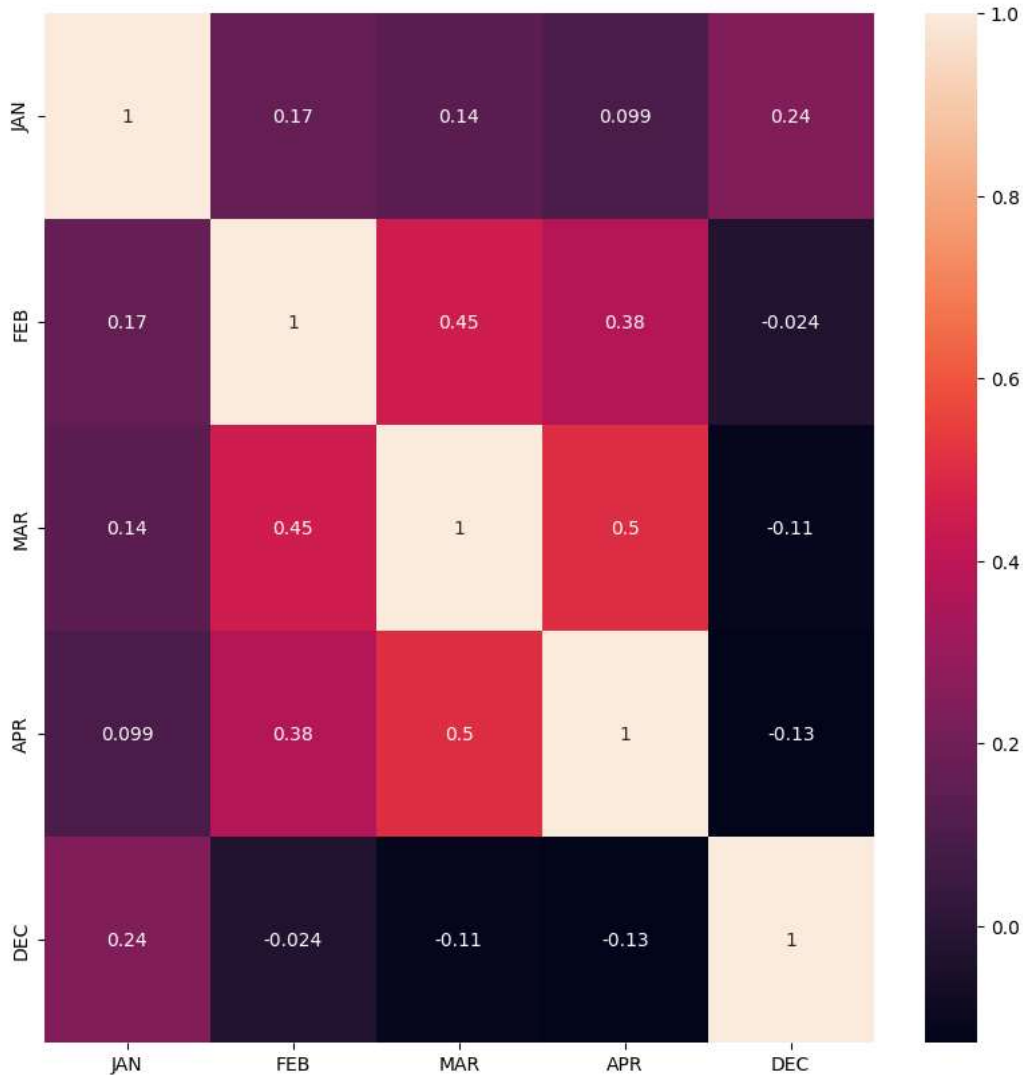
```
In [28]: 1 lr=LinearRegression()
2 lr.fit(x_train,y_train)
3 y_pred=lr.predict(x_test)
4 r2=r2_score(y_test,y_pred)
5 print("R2 score:",r2)
```

R2 score: 1.0

Ridge Regression

```
In [29]: 1 #Importing Libraries
2 from sklearn.linear_model import Lasso,Ridge
3 from sklearn.preprocessing import StandardScaler
```

```
In [30]: 1 plt.figure(figsize=(10,10))
2         sns.heatmap(df700.corr(),annot=True)
3         plt.show()
```



```
In [31]: 1 features=df.columns[0:5]
2         target=df.columns[-5]
```

```
In [32]: 1 x=df[features].values
2         y=df[target].values
3         x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.30,random_state=1)
4         print("The dimension of X_train is {}".format(x_train.shape))
5         print("The dimension of X_test is {}".format(x_test.shape))
```

The dimension of X_train is (2881, 5)
The dimension of X_test is (1235, 5)

```
In [33]: 1 lr = LinearRegression()
2         #Fit model
3         lr.fit(x_train, y_train)
4         #predict
5         actual = y_test
6         train_score_lr = lr.score(x_train, y_train)
7         test_score_lr = lr.score(x_test, y_test)
8         print("\nLinear Regression Model:\n")
9         print("The train score for lr model is {}".format(train_score_lr))
10        print("The test score for lr model is {}".format(test_score_lr))
```

Linear Regression Model:

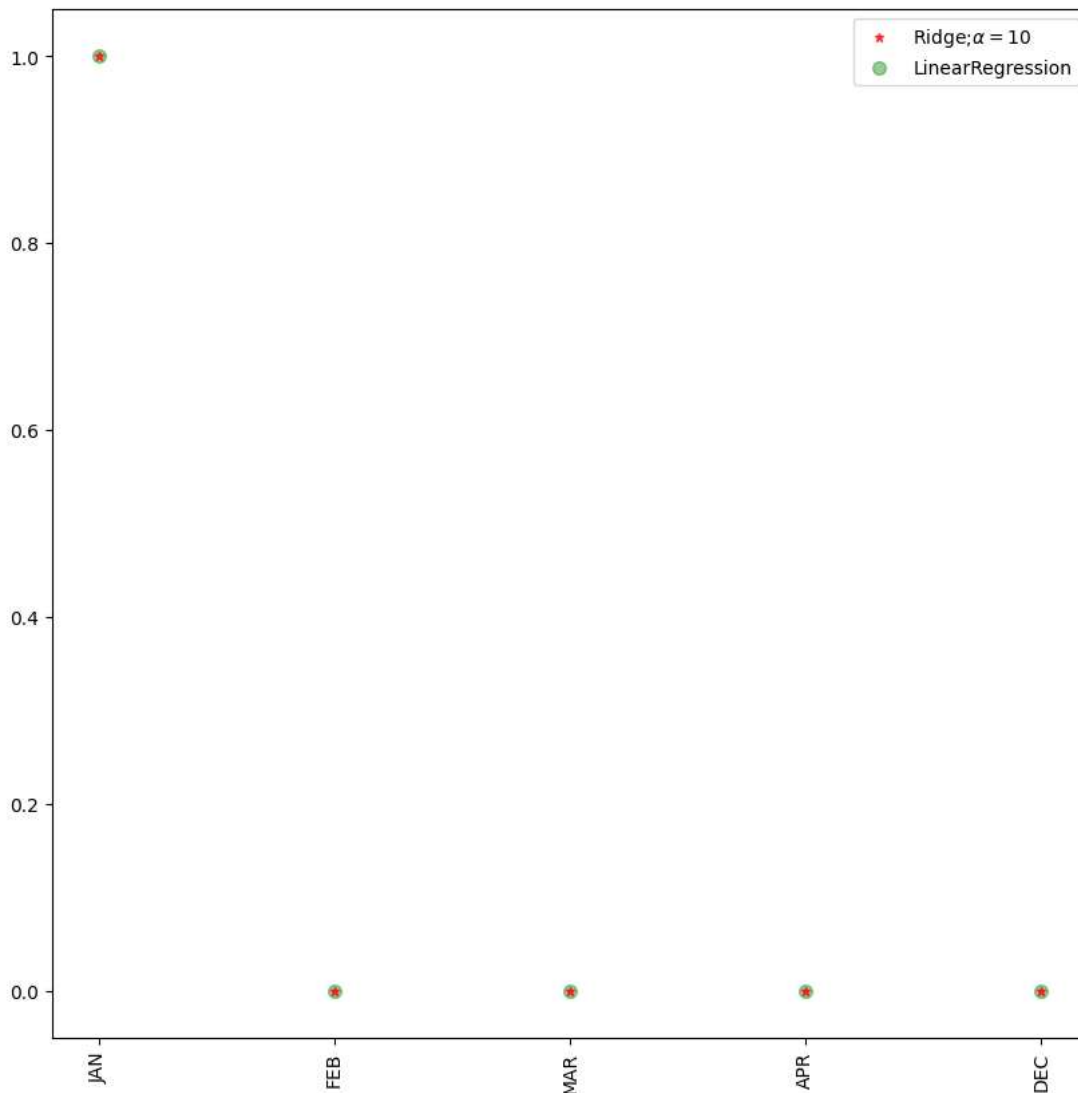
The train score for lr model is 1.0
The test score for lr model is 1.0


```
In [34]: 1 ridgeReg = Ridge(alpha=10)
2 ridgeReg.fit(x_train,y_train)
3 #train and test score for ridge regression
4 train_score_ridge = ridgeReg.score(x_train, y_train)
5 test_score_ridge = ridgeReg.score(x_test, y_test)
6 print("\nRidge Model:\n")
7 print("The train score for ridge model is {}".format(train_score_ridge))
8 print("The test score for ridge model is {}".format(test_score_ridge))
```

Ridge Model:

The train score for ridge model is 0.9999999999856335
The test score for ridge model is 0.9999999999840021

```
In [38]: .figure(figsize=(10,10))
.plot(features,ridgeReg.coef_,alpha=0.7,linestyle='none',marker="*",markersize=5,color='red',label=r'Ridge;$\alpha=10$',zorder=7)
.plot(features,lr.coef_,alpha=0.4,linestyle='none',marker="o",markersize=7,color='green',label='LinearRegression')
.xticks(rotation=90)
.legend()
.show()
```



Lasso Regression

```
In [40]: 1 #Importing Libraries
2 lasso= Lasso(alpha=10)
3 lasso.fit(x_train,y_train)
4 #train and test scorefor ridge regression
5 train_score_ls = lasso.score(x_train, y_train)
6 test_score_ls= lasso.score(x_test, y_test)
7 print("\nLasso Model:\n")
8 print("The train score for lasso model is {}".format(train_score_ls))
9 print("The test score for lasso model is {}".format(test_score_ls))
```

Lasso Model:

The train score for lasso model is 0.9999147271297208
The test score for lasso model is 0.9999147248375002

```
In [41]: 1 plt.figure(figsize=(10,10))
```

Out[41]: <Figure size 1000x1000 with 0 Axes>
<Figure size 1000x1000 with 0 Axes>

```
In [42]: 1 from sklearn.linear_model import LassoCV
```

```
In [43]: 1 #using the linear cv model
2 from sklearn.linear_model import RidgeCV
3 #cross validation
4 ridge_cv=RidgeCV(alphas =[0.0001,0.001,0.01,0.1,1,10]).fit(x_train,y_train)
5 #score
6 print(ridge_cv.score(x_train,y_train))
7 print(ridge_cv.score(x_test,y_test))
```

0.999999999261034
0.9999999993719254

```
In [44]: 1 #using the linear cv model
2 from sklearn.linear_model import LassoCV
3 #cross validation
4 lasso_cv=LassoCV(alphas =[0.0001,0.001,0.01,0.1,1,10]).fit(x_train,y_train)
5 #score
6 print(lasso_cv.score(x_train,y_train))
7 print(lasso_cv.score(x_test,y_test))
```

0.9999999999999915
0.9999999999999915

Elastic Regression

```
In [45]: 1 from sklearn.linear_model import ElasticNet
2
```

```
In [46]: 1 el=ElasticNet()
2 el.fit(x_train,y_train)
3 print(el.coef_)
4 print(el.intercept_)
5 el.score(x,y)
```

[9.99044548e-01 1.38835344e-05 4.58897515e-05 0.00000000e+00
0.00000000e+00]
0.01656567968369771

Out[46]: 0.9999991435191248

```
In [47]: 1 y_pred_elastic=el.predict(x_train)
```

```
In [48]: 1 mean_squared_error=np.mean((y_pred_elastic-y_train)**2)
2 print(mean_squared_error)
```

0.0009226812593703956

conclusion

from the above dataset we have got that rigidRegression is bestfit

In []:

1