# PROBLEM STATEMENT:To predict the rainfall based on various features of the dataset

```
In [1]:  #importing libraries
         import numpy as np
         import pandas as pd
         from sklearn.linear_model import LinearRegression
         from sklearn import preprocessing,svm
         from sklearn.model_selection import train_test_split
         import matplotlib.pyplot as plt
         import seaborn as sns
```

```
In [2]:  df=pd.read_csv(r"C:\Users\user\Downloads\rainfall in india 1901-2015.csv")
         df
```

Out[2]:

| | SUBDIVISION | YEAR | JAN | FEB | MAR | APR | MAY | JUN | JUL | AUG | SEP | OCT | NOV | DEC | ANNUAL | Jan-Feb | Mar-May |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 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 |
| 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 |
| 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 |
| 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 |
| 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 |

# Data cleaning and preprocessing

```
In [3]:  df.head()
```

Out[3]:

| | SUBDIVISION | YEAR | JAN | FEB | MAR | APR | MAY | JUN | JUL | AUG | SEP | OCT | NOV | DEC | ANNUAL | Jan-Feb | Mar-May | Jun-Sep |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 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. |
| 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. |
| 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. |
| 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. |
| 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. |

In [4]: `df.tail()`

Out[4]:

| | SUBDIVISION | YEAR | JAN | FEB | MAR | APR | MAY | JUN | JUL | AUG | SEP | OCT | NOV | DEC | ANNUAL | Jan-Feb | Mar-May | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 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 | 10 |
| 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 | 11 |
| 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 | 10 |
| 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 | 9 |
| 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 | 8 |

In [5]: `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 [6]: `df.shape`

Out[6]: `(4116, 19)`

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

Out[7]:

| | YEAR | JAN | FEB | MAR | APR | MAY | JUN | JUL | AUG | |
|---|---|---|---|---|---|---|---|---|---|---|
| count | 4116.000000 | 4112.000000 | 4113.000000 | 4110.000000 | 4112.000000 | 4113.000000 | 4111.000000 | 4109.000000 | 4112.000000 | 4110.000 |
| mean | 1958.218659 | 18.957320 | 21.805325 | 27.359197 | 43.127432 | 85.745417 | 230.234444 | 347.214334 | 290.263497 | 197.361 |
| std | 33.140898 | 33.585371 | 35.909488 | 46.959424 | 67.831168 | 123.234904 | 234.710758 | 269.539667 | 188.770477 | 135.408 |
| min | 1901.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.400000 | 0.000000 | 0.000000 | 0.100 |
| 25% | 1930.000000 | 0.600000 | 0.600000 | 1.000000 | 3.000000 | 8.600000 | 70.350000 | 175.600000 | 155.975000 | 100.525 |
| 50% | 1958.000000 | 6.000000 | 6.700000 | 7.800000 | 15.700000 | 36.600000 | 138.700000 | 284.800000 | 259.400000 | 173.900 |
| 75% | 1987.000000 | 22.200000 | 26.800000 | 31.300000 | 49.950000 | 97.200000 | 305.150000 | 418.400000 | 377.800000 | 265.800 |
| max | 2015.000000 | 583.700000 | 403.500000 | 605.600000 | 595.100000 | 1168.600000 | 1609.900000 | 2362.800000 | 1664.600000 | 1222.000 |

```python
In [8]: df.isnull().sum()
```

```
Out[8]: 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
```
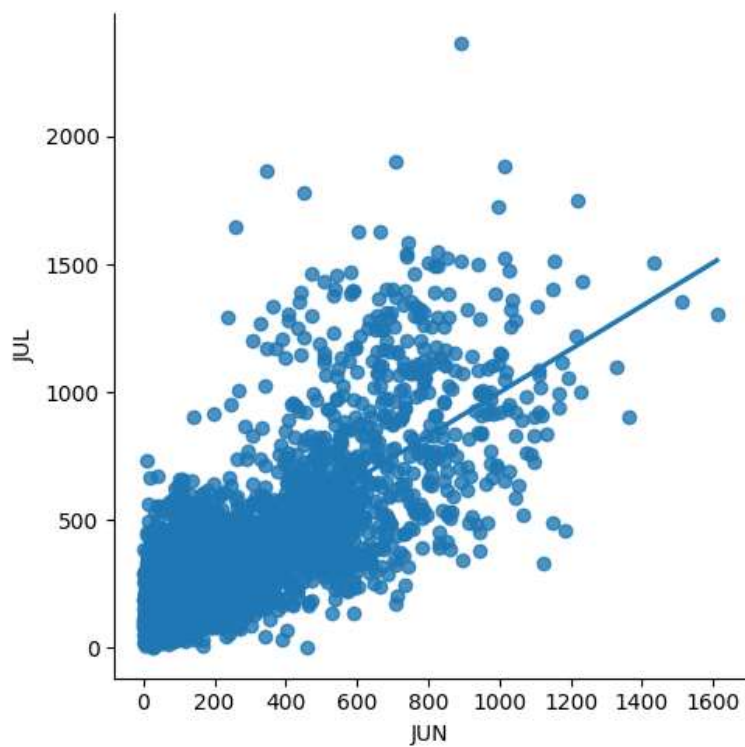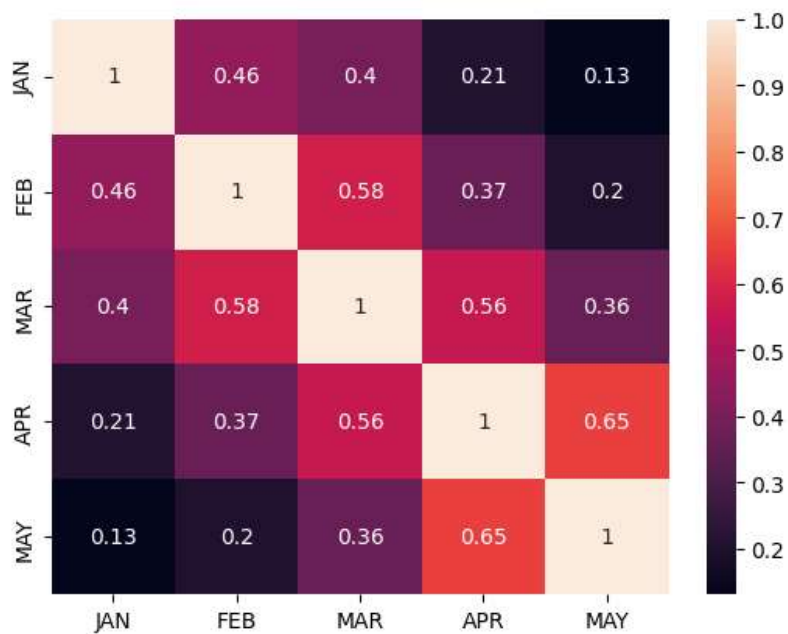
```python
In [9]: df.fillna(method="ffill",inplace=True)
```

```python
In [10]: df['YEAR'].value_counts()
```

```
Out[10]: 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
```
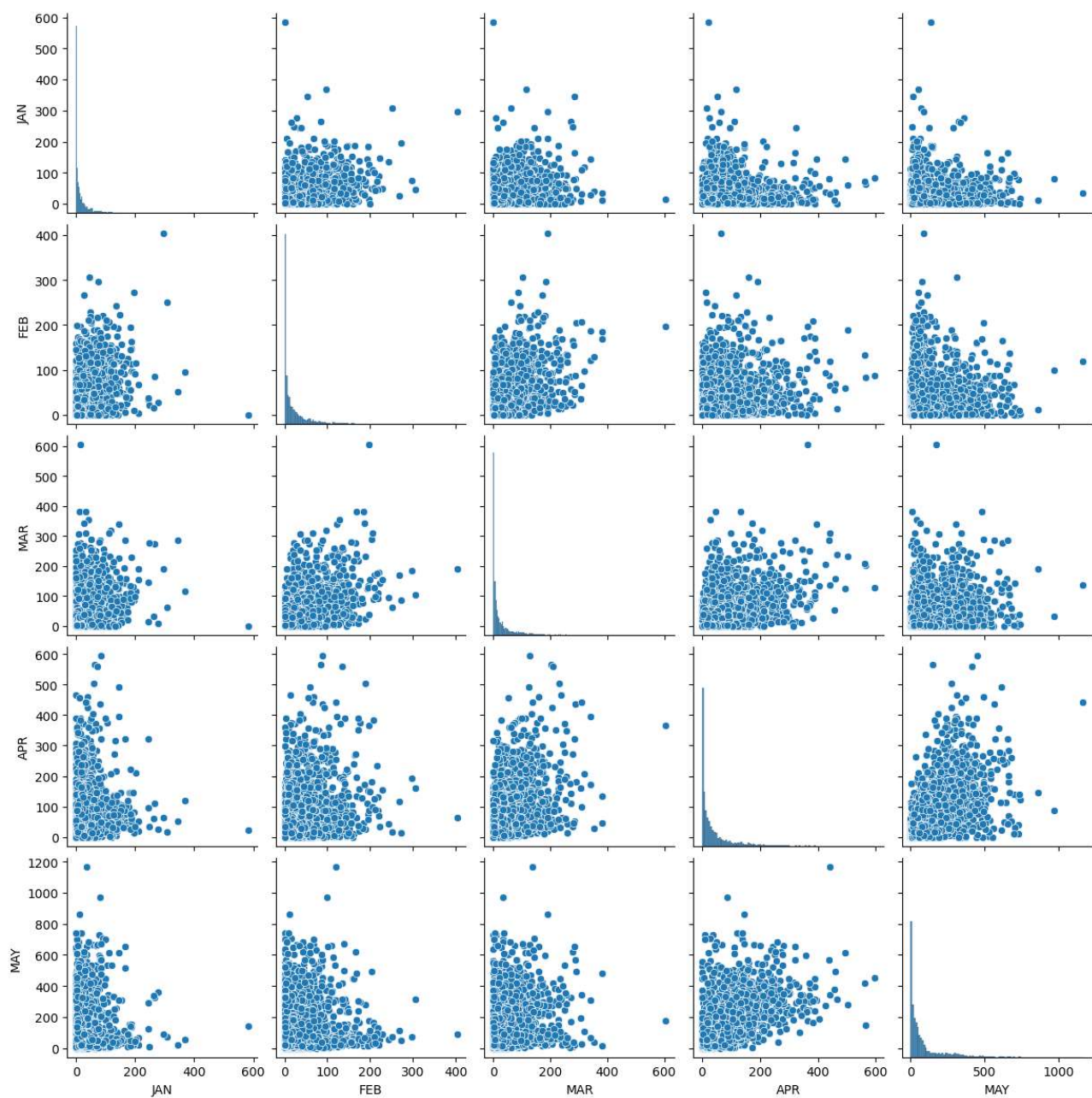
In [11]:
```python
sns.lmplot(x='JUN',y='JUL',order=2,data=df,ci=None)
plt.show()
```



In [12]:
```python
df=df[['JAN','FEB','MAR','APR','MAY']]
sns.heatmap(df.corr(),annot=True)
plt.show()
```

```
In [13]: sns.pairplot(df)
         plt.show()
```
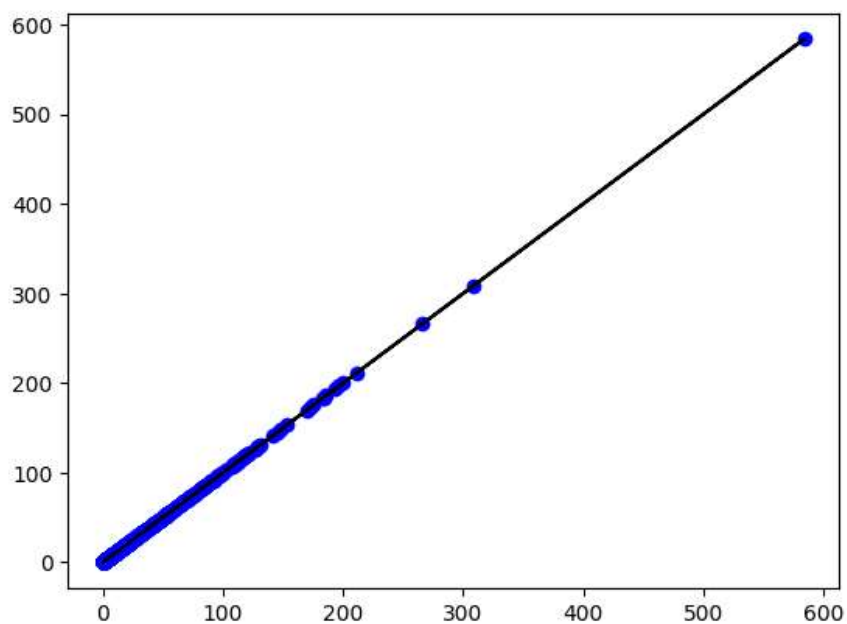


```
In [14]: x=np.array(df['FEB']).reshape(-1,1)
         y=x=np.array(df['JAN']).reshape(-1,1)
```

```
In [15]: x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.30)
```
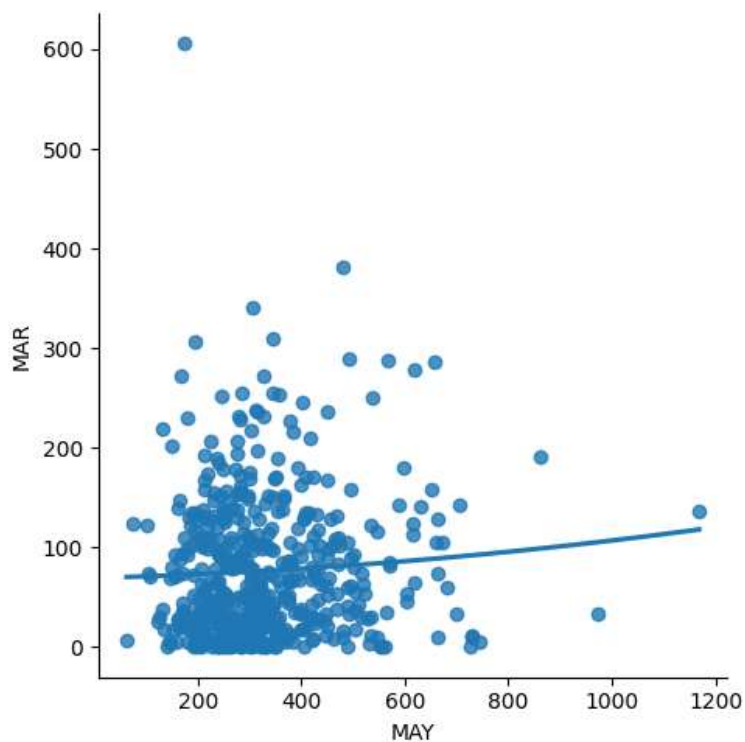
```
In [16]: regr=LinearRegression()
         regr.fit(x_train,y_train)
         print(regr.score(x_train,y_train))
```

```
1.0
```

In [17]:
```python
y_pred=regr.predict(x_test)
plt.scatter(x_test,y_test,color='blue')
plt.plot(x_test,y_pred,color='black')
plt.show()
```



In [18]:
```python
df500=df[:][:500]
sns.lmplot(x='MAY',y='MAR',order=2,ci=None,data=df500)
plt.show()
```
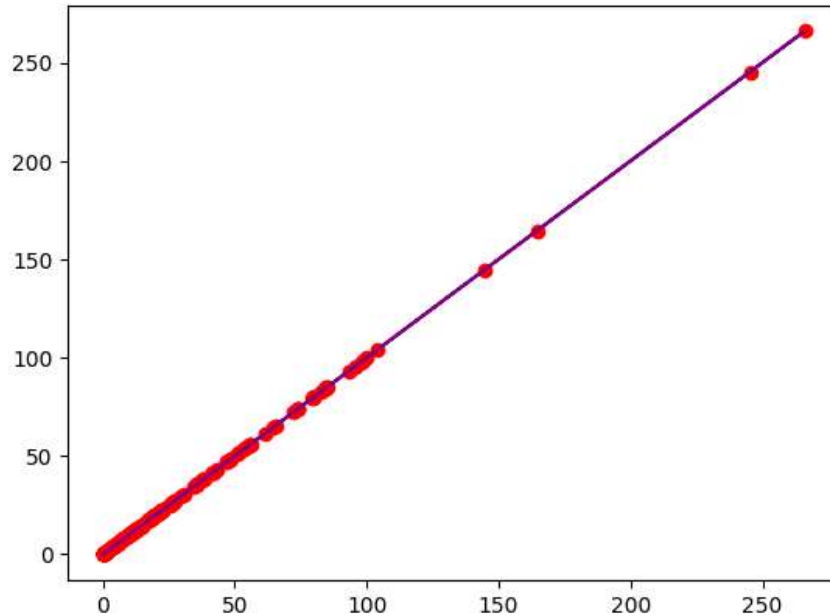


In [19]:
```python
df500.fillna(method='ffill',inplace=True)
```

In [20]:
```python
x=np.array(df500['FEB']).reshape(-1,1)
y=x=np.array(df500['JAN']).reshape(-1,1)
```

```
In [21]: X_train,X_test,y_train,y_test=train_test_split(x,y,test_size=0.33)
         regr.fit(X_train,y_train)
         regr.fit(X_test,y_test)
```

Out[21]:    ▼ LinearRegression

           LinearRegression()

```
In [22]: y_pred=regr.predict(X_test)
         plt.scatter(X_test,y_test,color='red')
         plt.plot(X_test,y_pred,color='purple')
         plt.show()
```



```
In [23]: from sklearn.metrics import r2_score
         model=LinearRegression()
         model.fit(X_train,y_train)
         y_pred=model.predict(X_test)
         r2=r2_score(y_test,y_pred)
         print("R2 Score:",r2)
```

R2 Score: 1.0

# RIGDE MODEL

```
In [24]: from sklearn.linear_model import Lasso,Ridge
         from sklearn.preprocessing import StandardScaler
```

```
In [25]: features= df.columns[0:4]
         target= df.columns[-4]
```

```
In [26]: x=np.array(df['FEB']).reshape(-1,1)
         y=np.array(df['JAN']).reshape(-1,1)
```

```
In [27]: x= df[features].values
         y= df[target].values
         x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.4,random_state=100)
```

```
In [28]:   ridgeReg=Ridge(alpha=10)
           ridgeReg.fit(x_train,y_train)
           train_score_ridge=ridgeReg.score(x_train,y_train)
           test_score_ridge=ridgeReg.score(x_test,y_test)
```

```
In [29]:   print("\n Ridge Model:\n")
           print("the train score for ridge model is{}".format(train_score_ridge))
           print("the test score for ridge model is{}".format(test_score_ridge))
```

```
 Ridge Model:

the train score for ridge model is0.9999999999828922
the test score for ridge model is0.9999999999836248
```

```
In [30]:   lr=LinearRegression()
```

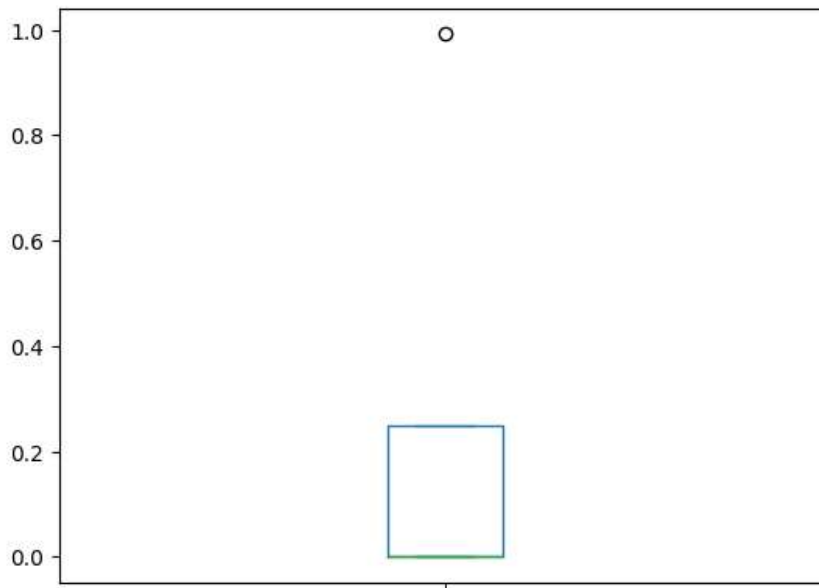# LASSO MODEL

```
In [32]:   print("\n Lasso Model:\n")
           lasso=Lasso(alpha=10)
           lasso.fit(x_train,y_train)
           train_score_ls=lasso.score(x_train,y_train)
           test_score_ls=lasso.score(x_test,y_test)
           print("The train score for ls model is {}".format(train_score_ls))
           print("The test score for ls model is{}".format(test_score_ls))
```

```
 Lasso Model:

The train score for ls model is 0.999936226881545
The test score for ls model is0.9999362250950667
```

```
In [34]:   pd.Series(lasso.coef_,features).sort_values(ascending=True).plot(kind="box")
```
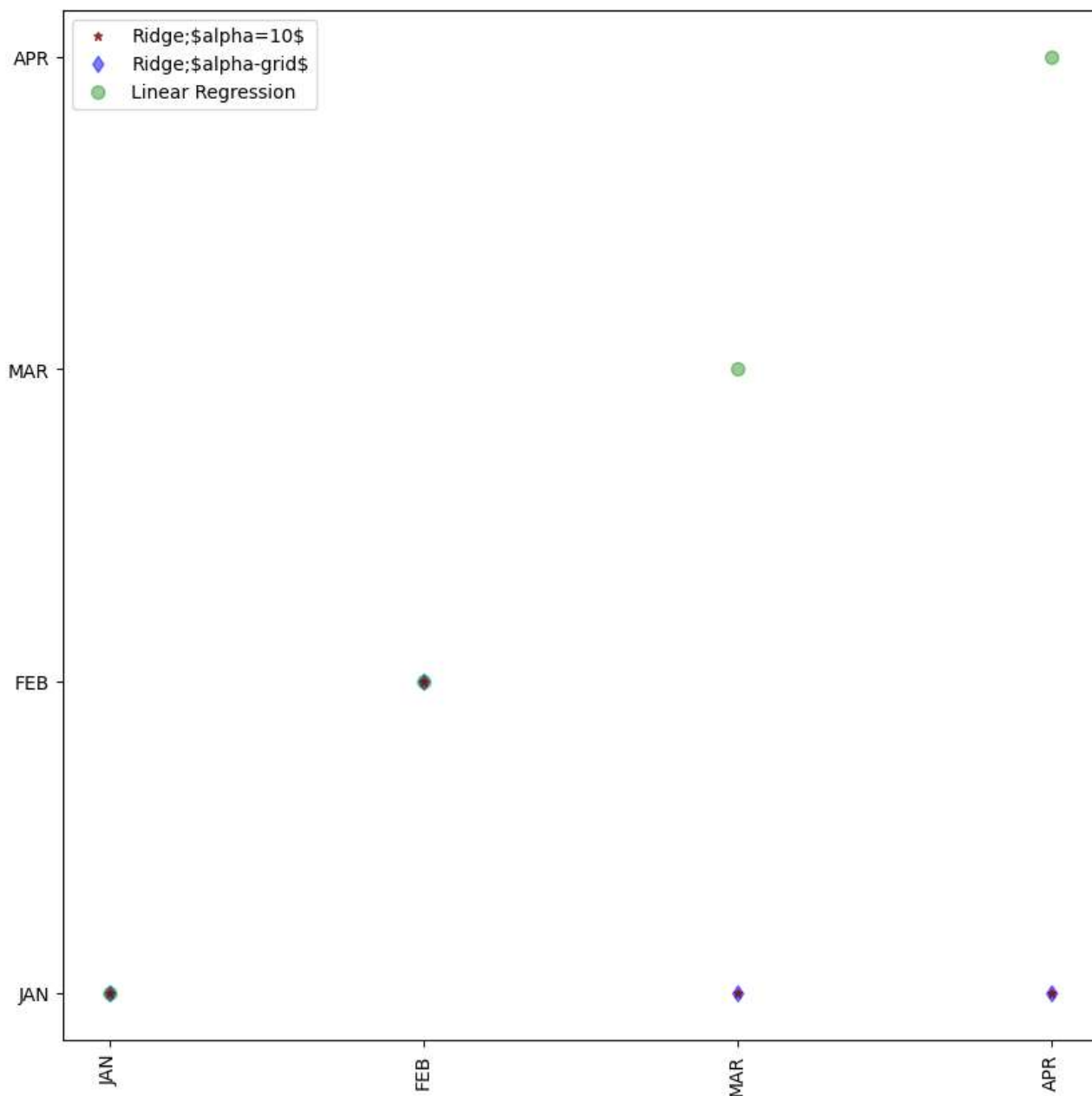
Out[34]:   <Axes: >



```
In [35]:   from sklearn.linear_model import LassoCV
           lasso_cv=LassoCV(alphas=[0.0001,0.001,0.01,1,10],random_state=0).fit(x_train,y_train)
           print(lasso_cv.score(x_train,y_train))
           print(lasso_cv.score(x_test,y_test))
```

```
0.9999999994453043
0.9999999993985869
```

```
In [37]: gure(figsize= (10,10))
         ot(features,ridgeReg.coef_,alpha=0.7,linestyle='none',marker="*",markersize=5,color='maroon',label=r'Ridge;\$
         ot(lasso_cv.coef_,alpha=0.5,linestyle='none',marker='d',markersize=6,color='blue',label=r'Ridge;\$alpha-grid$
         ot(features,alpha=0.4,linestyle='none',marker='o',markersize=7,color="green",label='Linear Regression')
         icks(rotation = 90)
         gend()
         ow()
```



## ELASTIC NET

```
In [38]: from sklearn.linear_model import ElasticNet
         eln=ElasticNet()
         eln.fit(x,y)
         print(eln.coef_)
         print(eln.intercept_)
         print(eln.score(x,y))
```

```
[0.00000000e+00 9.99136156e-01 1.56645396e-04 0.00000000e+00]
0.01455738923123917
0.9999994172977724
```

```
In [39]: y_pred_elastic = eln.predict(x_train)
         mean_squared_error=np.mean((y_pred_elastic - y_train)**2)
         print(mean_squared_error)
```

```
0.000725960274965855
```

## Conclusion:The above implemented models"Lasso and Ridge" regression is high accuracy compared to them

```
In [ ]:
```