

# PROBLEM STATEMENT:- TO PREDICT THE RAIN FALL BASED ON VARIOUS FEATURES OF THE DATASET

```
In [2]: import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.linear_model import LinearRegression
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
```

## Data collection

```
In [3]: rain_df=pd.read_csv(r"C:\Users\LENOVO\Desktop\rainfall in india 1901-2015.csv")
rain_df
```

Out[3]:

	SUBDIVISION	YEAR	JAN	FEB	MAR	APR	MAY	JUN	JUL	AUG	SEP	OCT
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
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
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
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
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
...	...	...	...	...	...	...	...	...	...	...	...	...
4111	LAKSHADWEEP	2011	5.1	2.8	3.1	85.9	107.2	153.6	350.2	254.0	255.2	117.4
4112	LAKSHADWEEP	2012	19.2	0.1	1.6	76.8	21.2	327.0	231.5	381.2	179.8	145.9
4113	LAKSHADWEEP	2013	26.2	34.4	37.5	5.3	88.3	426.2	296.4	154.4	180.0	72.8
4114	LAKSHADWEEP	2014	53.2	16.1	4.4	14.9	57.4	244.1	116.1	466.1	132.2	169.2
4115	LAKSHADWEEP	2015	2.2	0.5	3.7	87.1	133.1	296.6	257.5	146.4	160.4	165.4

4116 rows × 19 columns



In [4]:

rain\_df.shape

Out[4]: (4116, 19)

In [5]:

rain\_df.head()

Out[5]:

	SUBDIVISION	YEAR	JAN	FEB	MAR	APR	MAY	JUN	JUL	AUG	SEP	OCT	NOV
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
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
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
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
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

In [6]:

rain\_df.tail()

Out[6]:

	SUBDIVISION	YEAR	JAN	FEB	MAR	APR	MAY	JUN	JUL	AUG	SEP	OCT	N
4111	LAKSHADWEEP	2011	5.1	2.8	3.1	85.9	107.2	153.6	350.2	254.0	255.2	117.4	18
4112	LAKSHADWEEP	2012	19.2	0.1	1.6	76.8	21.2	327.0	231.5	381.2	179.8	145.9	1
4113	LAKSHADWEEP	2013	26.2	34.4	37.5	5.3	88.3	426.2	296.4	154.4	180.0	72.8	1
4114	LAKSHADWEEP	2014	53.2	16.1	4.4	14.9	57.4	244.1	116.1	466.1	132.2	169.2	5
4115	LAKSHADWEEP	2015	2.2	0.5	3.7	87.1	133.1	296.6	257.5	146.4	160.4	165.4	25

In [7]: `rain_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
```

## To find null values

```
In [8]: rain_df.fillna(method="ffill",inplace=True)
rain_df
```

Out[8]:

	SUBDIVISION	YEAR	JAN	FEB	MAR	APR	MAY	JUN	JUL	AUG	SEP	OCT
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
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
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
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
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
...	...	...	...	...	...	...	...	...	...	...	...	...
4111	LAKSHADWEEP	2011	5.1	2.8	3.1	85.9	107.2	153.6	350.2	254.0	255.2	117.4
4112	LAKSHADWEEP	2012	19.2	0.1	1.6	76.8	21.2	327.0	231.5	381.2	179.8	145.9
4113	LAKSHADWEEP	2013	26.2	34.4	37.5	5.3	88.3	426.2	296.4	154.4	180.0	72.8
4114	LAKSHADWEEP	2014	53.2	16.1	4.4	14.9	57.4	244.1	116.1	466.1	132.2	169.2
4115	LAKSHADWEEP	2015	2.2	0.5	3.7	87.1	133.1	296.6	257.5	146.4	160.4	165.4

4116 rows × 19 columns



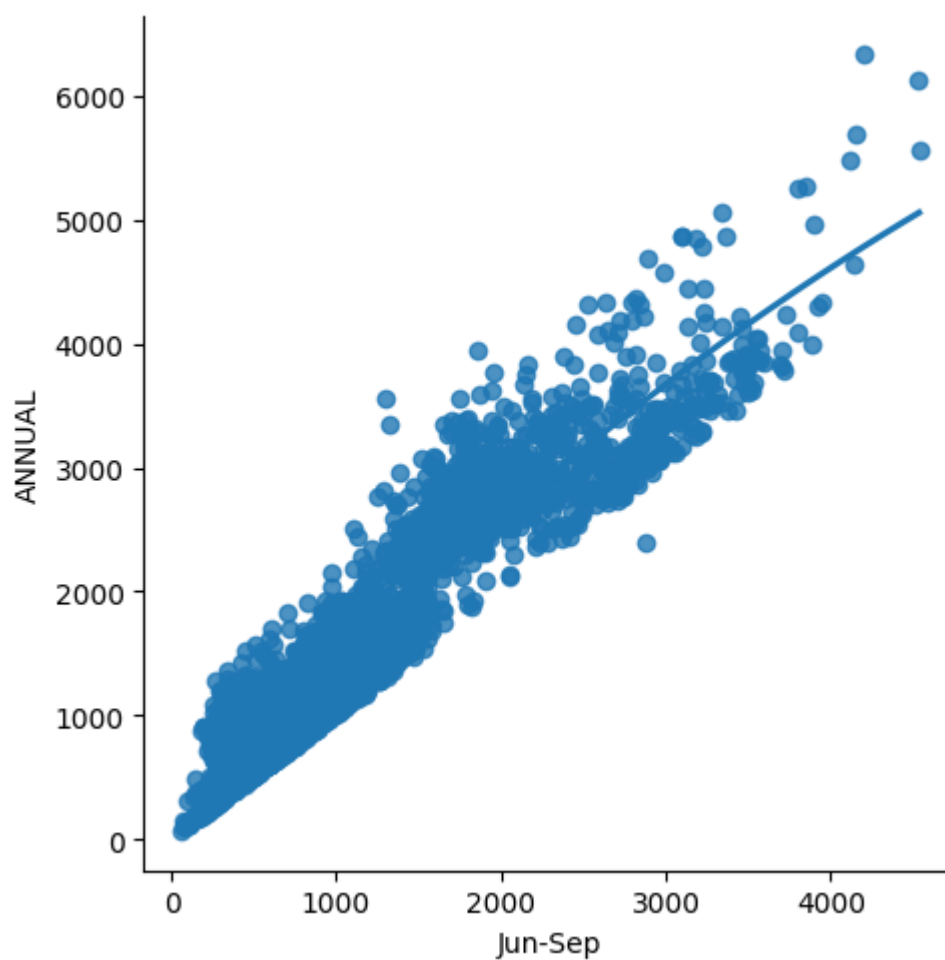
```
In [9]: rain_df.info()
```

```
<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                   4116 non-null  float64
3   FEB                   4116 non-null  float64
4   MAR                   4116 non-null  float64
5   APR                   4116 non-null  float64
6   MAY                   4116 non-null  float64
7   JUN                   4116 non-null  float64
8   JUL                   4116 non-null  float64
9   AUG                   4116 non-null  float64
10  SEP                   4116 non-null  float64
11  OCT                   4116 non-null  float64
12  NOV                   4116 non-null  float64
13  DEC                   4116 non-null  float64
14  ANNUAL                4116 non-null  float64
15  Jan-Feb               4116 non-null  float64
16  Mar-May               4116 non-null  float64
17  Jun-Sep               4116 non-null  float64
18  Oct-Dec               4116 non-null  float64
dtypes: float64(17), int64(1), object(1)
memory usage: 611.1+ KB
```

## Visualization

```
In [10]: sns.lmplot(x='Jun-Sep',y='ANNUAL',data=rain_df,order=2,ci=None)
```

```
Out[10]: <seaborn.axisgrid.FacetGrid at 0x2974684b150>
```



```
In [11]: rain_df['SUBDIVISION'].value_counts()
```

```
Out[11]: SUBDIVISION
WEST MADHYA PRADESH          115
EAST RAJASTHAN               115
COASTAL KARNATAKA            115
TAMIL NADU                   115
RAYALSEEMA                   115
TELANGANA                    115
COASTAL ANDHRA PRADESH       115
CHHATTISGARH                 115
VIDARBHA                     115
MATATHWADA                   115
MADHYA MAHARASHTRA           115
KONKAN & GOA                  115
SAURASHTRA & KUTCH           115
GUJARAT REGION               115
EAST MADHYA PRADESH          115
KERALA                        115
WEST RAJASTHAN                115
SOUTH INTERIOR KARNATAKA     115
JAMMU & KASHMIR               115
HIMACHAL PRADESH             115
PUNJAB                       115
HARYANA DELHI & CHANDIGARH   115
UTTARAKHAND                  115
WEST UTTAR PRADESH           115
EAST UTTAR PRADESH           115
BIHAR                        115
JHARKHAND                     115
ORISSA                        115
GANGETIC WEST BENGAL         115
SUB HIMALAYAN WEST BENGAL & SIKKIM 115
NAGA MANI MIZO TRIPURA      115
ASSAM & MEGHALAYA             115
NORTH INTERIOR KARNATAKA     115
LAKSHADWEEP                  114
ANDAMAN & NICOBAR ISLANDS     110
ARUNACHAL PRADESH            97
Name: count, dtype: int64
```

```
In [12]: rain_df.columns
```

```
Out[12]: Index(['SUBDIVISION', 'YEAR', 'JAN', 'FEB', 'MAR', 'APR', 'MAY', 'JUN', 'JUL',
               'AUG', 'SEP', 'OCT', 'NOV', 'DEC', 'ANNUAL', 'Jan-Feb', 'Mar-May',
               'Jun-Sep', 'Oct-Dec'],
              dtype='object')
```

## Feature selection

```
In [13]: x=np.array(rain_df['Jun-Sep']).reshape(-1,1)
y=np.array(rain_df['ANNUAL']).reshape(-1,1)
```

```
In [14]: X_train,X_test,y_train,y_test=train_test_split(x,y,test_size=0.3,random_state=1)
```

## Linear Regression

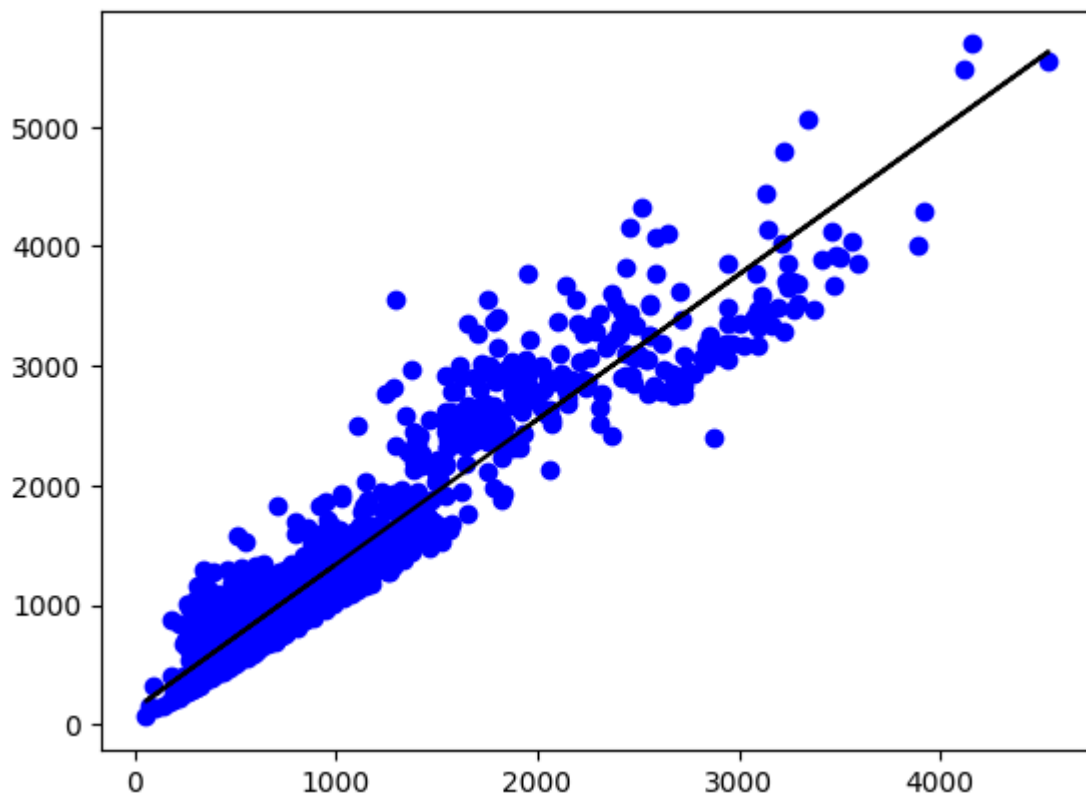
```
In [15]: regr=LinearRegression()
X_train,X_test,y_train,y_test=train_test_split(x,y,test_size=0.3)
regr.fit(X_train,y_train)
regr.fit(X_train,y_train)
```

```
Out[15]: ▾ LinearRegression
LinearRegression()
```

```
In [16]: print(regr.score(X_test,y_test))
```

0.889280453409494

```
In [17]: 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]: 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: 0.889280453409494

## Ridge model

```
In [19]: from sklearn.linear_model import Lasso,Ridge
from sklearn.preprocessing import StandardScaler
```

```
In [20]: x=np.array(rain_df['Jun-Sep']).reshape(-1,1)
y=np.array(rain_df['ANNUAL']).reshape(-1,1)
```

```
In [21]: x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3,random_state=1)
```

```
In [22]: 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 [23]: print("the train score for ridge model is{}".format(train_score_ridge))
print("the test score for ridge model is{}".format(test_score_ridge))
```

the train score for ridge model is0.88425882286469  
the test score for ridge model is0.897227115927482

## Lasso model

```
In [24]: 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 0.8842588226165713  
The test score for ls model is0.8972268430634212

## Elastic Net

```
In [25]: from sklearn.linear_model import ElasticNet
elnet=ElasticNet()
elnet.fit(x,y)
print(elnet.coef_)
print(elnet.intercept_)
print(elnet.score(x,y))
```

```
[1.20838561]
[129.62392906]
0.8882690674973049
```

## Conclusion

**score for linearregression is 0.889280453409494 ,score for Ridge model is 0.897227115927482,score for Lasso model is 0.8972268430634212,score for elastic net is 0.8882690674973049.From all the above models we can conclude that ridge model is the best model to predict the rainfall based on various features of the dataset.**

In [ ]:

In [ ]:

In [ ]: