

## PROBLEM STATEMENT:

TO CHECK HOW BEST FIT IS IT?

## importing required libraries

```
In [64]: ▶ import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
```

## data collection

```
In [65]: df=pd.read_csv(r"C:\Users\MY HOME\Downloads\rainfall in india 1901-2015.csv")
df
```

Out[65]:

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

4116 rows × 19 columns



```
In [66]: df1=pd.read_csv(r"C:\Users\MY HOME\Downloads\district wise rainfall normal.csv")
df1
```

Out[66]:

	STATE_UT_NAME	DISTRICT	JAN	FEB	MAR	APR	MAY	JUN	JUL	AUG	SEP	OCT	NOV	DEC	ANNUAL	Jan Feb
0	ANDAMAN And NICOBAR ISLANDS	NICOBAR	107.3	57.9	65.2	117.0	358.5	295.5	285.0	271.9	354.8	326.0	315.2	250.9	2805.2	165.2
1	ANDAMAN And NICOBAR ISLANDS	SOUTH ANDAMAN	43.7	26.0	18.6	90.5	374.4	457.2	421.3	423.1	455.6	301.2	275.8	128.3	3015.7	69.7
2	ANDAMAN And NICOBAR ISLANDS	N & M ANDAMAN	32.7	15.9	8.6	53.4	343.6	503.3	465.4	460.9	454.8	276.1	198.6	100.0	2913.3	48.6
3	ARUNACHAL PRADESH	LOHIT	42.2	80.8	176.4	358.5	306.4	447.0	660.1	427.8	313.6	167.1	34.1	29.8	3043.8	123.0
4	ARUNACHAL PRADESH	EAST SIANG	33.3	79.5	105.9	216.5	323.0	738.3	990.9	711.2	568.0	206.9	29.5	31.7	4034.7	112.8
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
636	KERALA	IDUKKI	13.4	22.1	43.6	150.4	232.6	651.6	788.9	527.3	308.4	343.2	172.9	48.1	3302.5	35.2
637	KERALA	KASARGOD	2.3	1.0	8.4	46.9	217.6	999.6	1108.5	636.3	263.1	234.9	84.6	18.4	3621.6	3.2
638	KERALA	PATHANAMTHITTA	19.8	45.2	73.9	184.9	294.7	556.9	539.9	352.7	266.2	359.4	213.5	51.3	2958.4	65.0
639	KERALA	WAYANAD	4.8	8.3	17.5	83.3	174.6	698.1	1110.4	592.9	230.7	213.1	93.6	25.8	3253.1	13.2
640	LAKSHADWEEP	LAKSHADWEEP	20.8	14.7	11.8	48.9	171.7	330.2	287.7	217.5	163.1	157.1	117.7	58.8	1600.0	35.2

641 rows × 19 columns



## data cleaning

In [67]: `df.head()`

Out[67]:

	SUBDIVISION	YEAR	JAN	FEB	MAR	APR	MAY	JUN	JUL	AUG	SEP	OCT	NOV	DEC	ANNUAL	Jan-Feb	Mar-May	Jun-Sep	C
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	98
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	71
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	69
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	57
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	63

```
In [68]: 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 [69]: df.describe()
```

```
Out[69]:
```

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

## finding missing values

```
In [70]: df.fillna(method="ffill",inplace=True)
df
```

Out[70]:

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

4116 rows × 19 columns



```
In [71]: 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
```



In [72]: `df1.head()`

Out[72]:

	STATE_UT_NAME	DISTRICT	JAN	FEB	MAR	APR	MAY	JUN	JUL	AUG	SEP	OCT	NOV	DEC	ANNUAL	Jan-Feb	Mar-May	
0	ANDAMAN And NICOBAR ISLANDS	NICOBAR	107.3	57.9	65.2	117.0	358.5	295.5	285.0	271.9	354.8	326.0	315.2	250.9	2805.2	165.2	540.7	1
1	ANDAMAN And NICOBAR ISLANDS	SOUTH ANDAMAN	43.7	26.0	18.6	90.5	374.4	457.2	421.3	423.1	455.6	301.2	275.8	128.3	3015.7	69.7	483.5	1
2	ANDAMAN And NICOBAR ISLANDS	N & M ANDAMAN	32.7	15.9	8.6	53.4	343.6	503.3	465.4	460.9	454.8	276.1	198.6	100.0	2913.3	48.6	405.6	1
3	ARUNACHAL PRADESH	LOHIT	42.2	80.8	176.4	358.5	306.4	447.0	660.1	427.8	313.6	167.1	34.1	29.8	3043.8	123.0	841.3	1
4	ARUNACHAL PRADESH	EAST SIANG	33.3	79.5	105.9	216.5	323.0	738.3	990.9	711.2	568.0	206.9	29.5	31.7	4034.7	112.8	645.4	3

In [73]: `df1.describe()`

Out[73]:

	JAN	FEB	MAR	APR	MAY	JUN	JUL	AUG	SEP	OCT	
<b>count</b>	641.000000	641.000000	641.000000	641.000000	641.000000	641.000000	641.000000	641.000000	641.000000	641.000000	641.000000
<b>mean</b>	18.355070	20.984399	30.034789	45.543214	81.535101	196.007332	326.033697	291.152262	194.609048	90.446334	34.111111
<b>std</b>	21.082806	27.729596	45.451082	71.556279	111.960390	196.556284	221.364643	152.647325	99.830540	74.990685	59.371111
<b>min</b>	0.000000	0.000000	0.000000	0.000000	0.900000	3.800000	11.600000	14.100000	8.600000	3.100000	1.200000
<b>25%</b>	6.900000	7.000000	7.000000	5.000000	12.100000	68.800000	206.400000	194.600000	128.800000	34.300000	6.600000
<b>50%</b>	13.300000	12.300000	12.700000	15.100000	33.900000	131.900000	293.700000	284.800000	181.300000	62.600000	12.900000
<b>75%</b>	19.200000	24.100000	33.200000	48.300000	91.900000	226.600000	374.800000	358.100000	234.100000	130.200000	32.300000
<b>max</b>	144.500000	229.600000	367.900000	554.400000	733.700000	1476.200000	1820.900000	1522.100000	826.300000	517.700000	475.100000

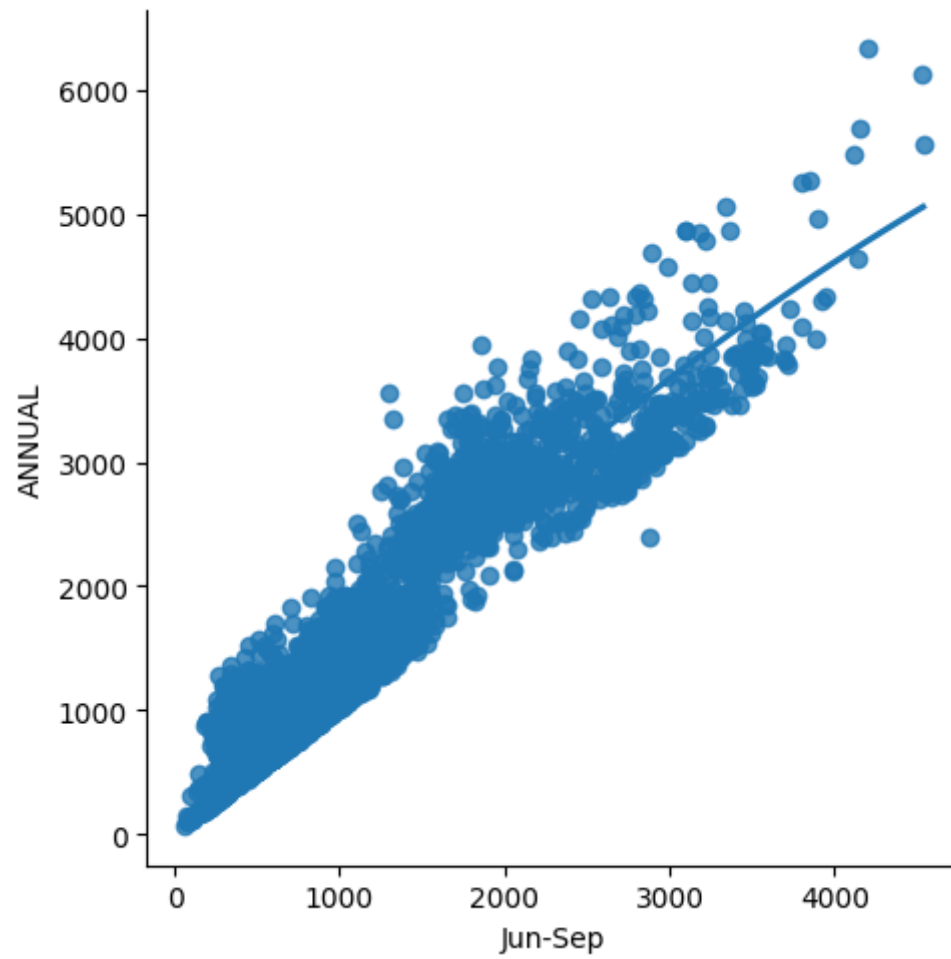


In [74]: `df1.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 641 entries, 0 to 640
Data columns (total 19 columns):
#   Column          Non-Null Count  Dtype
---  -
0   STATE_UT_NAME    641 non-null    object
1   DISTRICT         641 non-null    object
2   JAN              641 non-null    float64
3   FEB              641 non-null    float64
4   MAR              641 non-null    float64
5   APR              641 non-null    float64
6   MAY              641 non-null    float64
7   JUN              641 non-null    float64
8   JUL              641 non-null    float64
9   AUG              641 non-null    float64
10  SEP              641 non-null    float64
11  OCT              641 non-null    float64
12  NOV              641 non-null    float64
13  DEC              641 non-null    float64
14  ANNUAL           641 non-null    float64
15  Jan-Feb          641 non-null    float64
16  Mar-May          641 non-null    float64
17  Jun-Sep          641 non-null    float64
18  Oct-Dec          641 non-null    float64
dtypes: float64(17), object(2)
memory usage: 95.3+ KB
```

## data visualization

```
In [75]: ▶ sns.lmplot(x="Jun-Sep",y="ANNUAL",order=2,data=df,ci=None)  
plt.show()
```



```
In [76]: #taking selected columns from dataset  
df=df[['Jun-Sep', 'ANNUAL']]  
df
```

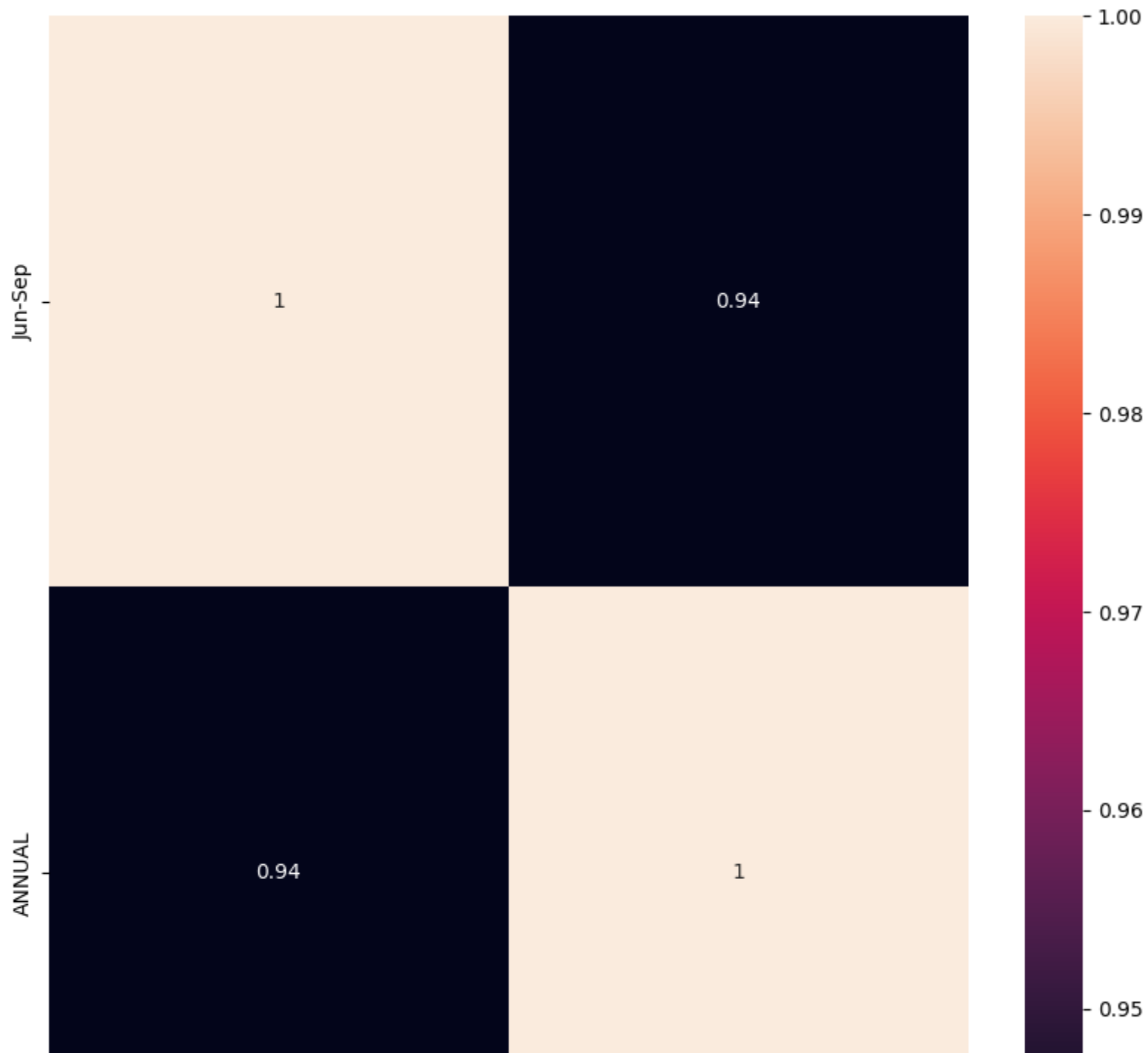
Out[76]:

	Jun-Sep	ANNUAL
0	1696.3	3373.2
1	2185.9	3520.7
2	1874.0	2957.4
3	1977.6	3079.6
4	1624.9	2566.7
...	...	...
4111	1013.0	1533.7
4112	1119.5	1405.5
4113	1057.0	1426.3
4114	958.5	1395.0
4115	860.9	1642.9

4116 rows × 2 columns

```
In [77]: ▶ from sklearn.preprocessing import StandardScaler  
plt.figure(figsize=(10,10))  
sns.heatmap(df.corr(),annot=True)  
plt.show()
```









```
In [78]: x=np.array(df['Jun-Sep']).reshape(-1,1)
y=np.array(df['ANNUAL']).reshape(-1,1)
#Dropping any rows with Nan values
df.dropna(inplace=True)
df
```

C:\Users\MY HOME\AppData\Local\Temp\ipykernel\_10092\2395007685.py:4: 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) ([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)
```

Out[78]:

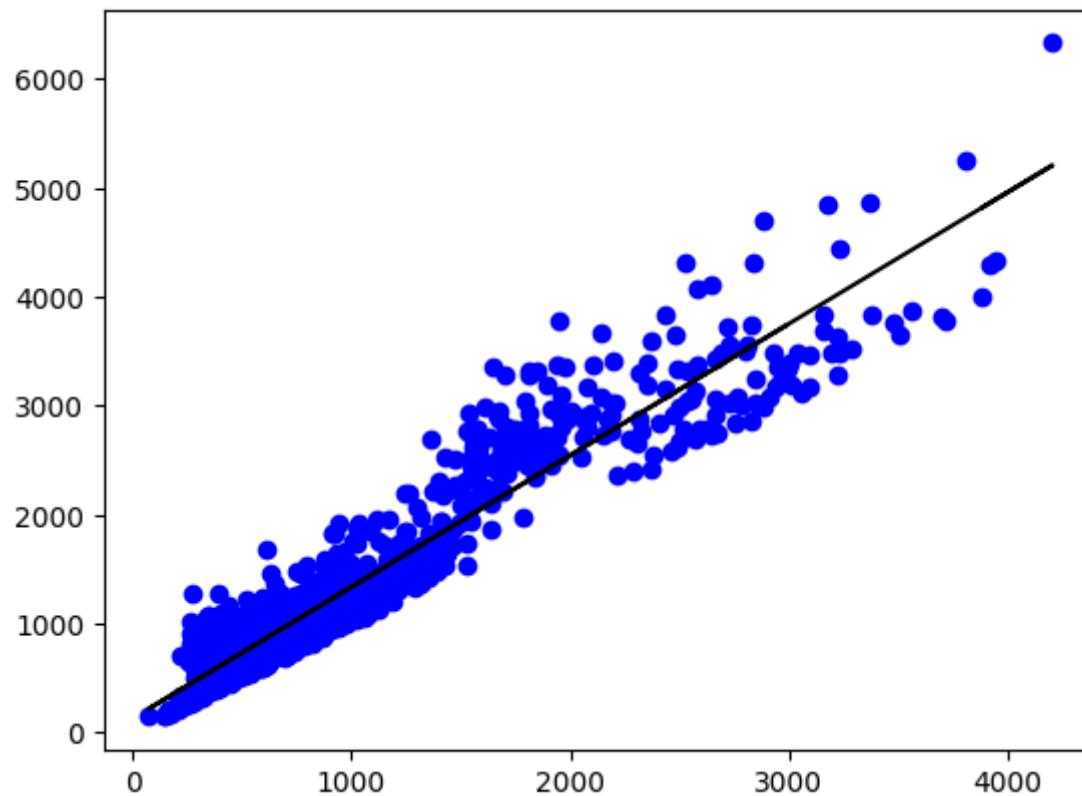
	Jun-Sep	ANNUAL
0	1696.3	3373.2
1	2185.9	3520.7
2	1874.0	2957.4
3	1977.6	3079.6
4	1624.9	2566.7
...	...	...
4111	1013.0	1533.7
4112	1119.5	1405.5
4113	1057.0	1426.3
4114	958.5	1395.0
4115	860.9	1642.9

4116 rows × 2 columns

```
In [79]: > 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(regr.score(x_test,y_test))
```

0.8923510009567511

```
In [80]: > y_pred=regr.predict(x_test)
> plt.scatter(x_test,y_test,color='b')
> plt.plot(x_test,y_pred,color='k')
> plt.show()
```



# RIDGE REGRESSION

In [ ]: ▶ ----->To check better accuracy/fit.

In [82]: ▶ `from sklearn.linear_model import Ridge, RidgeCV, Lasso`

In [83]: ▶ `ridge=Ridge(alpha=2)  
ridge.fit(x_train,y_train)  
train_score_ridge=ridge.score(x_train,y_train)  
test_score_ridge=ridge.score(x_test,y_test)  
print("\nLinearRegression\n", (train_score_ridge))  
print(test_score_ridge)`

```
LinearRegression  
0.8869000079295701  
0.8923510009454426
```

# LASSO REGRESSION

In [ ]: ▶ ----->To check **for** better model.

```
In [84]: ▶ #Lasso regression model
print("\nLasso 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 0.8869000076876236

The test score for ls model is 0.8923508560910814

## ELASTICNET

```
In [85]: ▶ from sklearn.linear_model import ElasticNet
```

```
In [86]: ▶ a=ElasticNet()
a.fit(x,y)
print(a.coef_)
print(a.intercept_)
```

[1.20838561]

[129.62392906]

```
In [87]: ► y_pred_elastic=a.predict(x_train)
mean_squared_error=np.mean((y_pred_elastic-y_train)**2)
print(mean_squared_error)
```

1556420.3818782303

In [ ]: ►

## # logistic Regression

```
In [31]: ► ridgeReg = Ridge(alpha=10)
ridgeReg.fit(x_train,y_train)
#train and test score for ridge regression
train_score_ridge = ridgeReg.score(x_train, y_train)
test_score_ridge = ridgeReg.score(x_test, y_test)
print("\nRidge 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 is 0.8890123548189879

The test score for ridge model is 0.8858366709919248

## CONCLUSION:

-----> compared to all models we did based on accuracies we conclude that the Linear Regression is the best model for this dataset with the accuracy of 88%.

# **DISTRICT WISE DATASET**

## **data collection**

```
In [50]: df1=pd.read_csv(r"C:\Users\MY HOME\Downloads\district wise rainfall normal.csv")
df1
```

Out[50]:

	STATE_UT_NAME	DISTRICT	JAN	FEB	MAR	APR	MAY	JUN	JUL	AUG	SEP	OCT	NOV	DEC	ANNUAL	Jan Feb
0	ANDAMAN And NICOBAR ISLANDS	NICOBAR	107.3	57.9	65.2	117.0	358.5	295.5	285.0	271.9	354.8	326.0	315.2	250.9	2805.2	165.2
1	ANDAMAN And NICOBAR ISLANDS	SOUTH ANDAMAN	43.7	26.0	18.6	90.5	374.4	457.2	421.3	423.1	455.6	301.2	275.8	128.3	3015.7	69.7
2	ANDAMAN And NICOBAR ISLANDS	N & M ANDAMAN	32.7	15.9	8.6	53.4	343.6	503.3	465.4	460.9	454.8	276.1	198.6	100.0	2913.3	48.6
3	ARUNACHAL PRADESH	LOHIT	42.2	80.8	176.4	358.5	306.4	447.0	660.1	427.8	313.6	167.1	34.1	29.8	3043.8	123.0
4	ARUNACHAL PRADESH	EAST SIANG	33.3	79.5	105.9	216.5	323.0	738.3	990.9	711.2	568.0	206.9	29.5	31.7	4034.7	112.8
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
636	KERALA	IDUKKI	13.4	22.1	43.6	150.4	232.6	651.6	788.9	527.3	308.4	343.2	172.9	48.1	3302.5	35.1
637	KERALA	KASARGOD	2.3	1.0	8.4	46.9	217.6	999.6	1108.5	636.3	263.1	234.9	84.6	18.4	3621.6	3.1
638	KERALA	PATHANAMTHITTA	19.8	45.2	73.9	184.9	294.7	556.9	539.9	352.7	266.2	359.4	213.5	51.3	2958.4	65.0
639	KERALA	WAYANAD	4.8	8.3	17.5	83.3	174.6	698.1	1110.4	592.9	230.7	213.1	93.6	25.8	3253.1	13.1
640	LAKSHADWEEP	LAKSHADWEEP	20.8	14.7	11.8	48.9	171.7	330.2	287.7	217.5	163.1	157.1	117.7	58.8	1600.0	35.1

641 rows × 19 columns



## data cleaning



In [51]: `df1.head()`

Out[51]:

	STATE_UT_NAME	DISTRICT	JAN	FEB	MAR	APR	MAY	JUN	JUL	AUG	SEP	OCT	NOV	DEC	ANNUAL	Jan-Feb	Mar-May	
0	ANDAMAN And NICOBAR ISLANDS	NICOBAR	107.3	57.9	65.2	117.0	358.5	295.5	285.0	271.9	354.8	326.0	315.2	250.9	2805.2	165.2	540.7	1
1	ANDAMAN And NICOBAR ISLANDS	SOUTH ANDAMAN	43.7	26.0	18.6	90.5	374.4	457.2	421.3	423.1	455.6	301.2	275.8	128.3	3015.7	69.7	483.5	1
2	ANDAMAN And NICOBAR ISLANDS	N & M ANDAMAN	32.7	15.9	8.6	53.4	343.6	503.3	465.4	460.9	454.8	276.1	198.6	100.0	2913.3	48.6	405.6	1
3	ARUNACHAL PRADESH	LOHIT	42.2	80.8	176.4	358.5	306.4	447.0	660.1	427.8	313.6	167.1	34.1	29.8	3043.8	123.0	841.3	1
4	ARUNACHAL PRADESH	EAST SIANG	33.3	79.5	105.9	216.5	323.0	738.3	990.9	711.2	568.0	206.9	29.5	31.7	4034.7	112.8	645.4	3

In [52]: `df1.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 641 entries, 0 to 640
Data columns (total 19 columns):
#   Column          Non-Null Count  Dtype
---  -
0   STATE_UT_NAME    641 non-null    object
1   DISTRICT         641 non-null    object
2   JAN              641 non-null    float64
3   FEB              641 non-null    float64
4   MAR              641 non-null    float64
5   APR              641 non-null    float64
6   MAY              641 non-null    float64
7   JUN              641 non-null    float64
8   JUL              641 non-null    float64
9   AUG              641 non-null    float64
10  SEP              641 non-null    float64
11  OCT              641 non-null    float64
12  NOV              641 non-null    float64
13  DEC              641 non-null    float64
14  ANNUAL           641 non-null    float64
15  Jan-Feb          641 non-null    float64
16  Mar-May          641 non-null    float64
17  Jun-Sep          641 non-null    float64
18  Oct-Dec          641 non-null    float64
dtypes: float64(17), object(2)
memory usage: 95.3+ KB
```

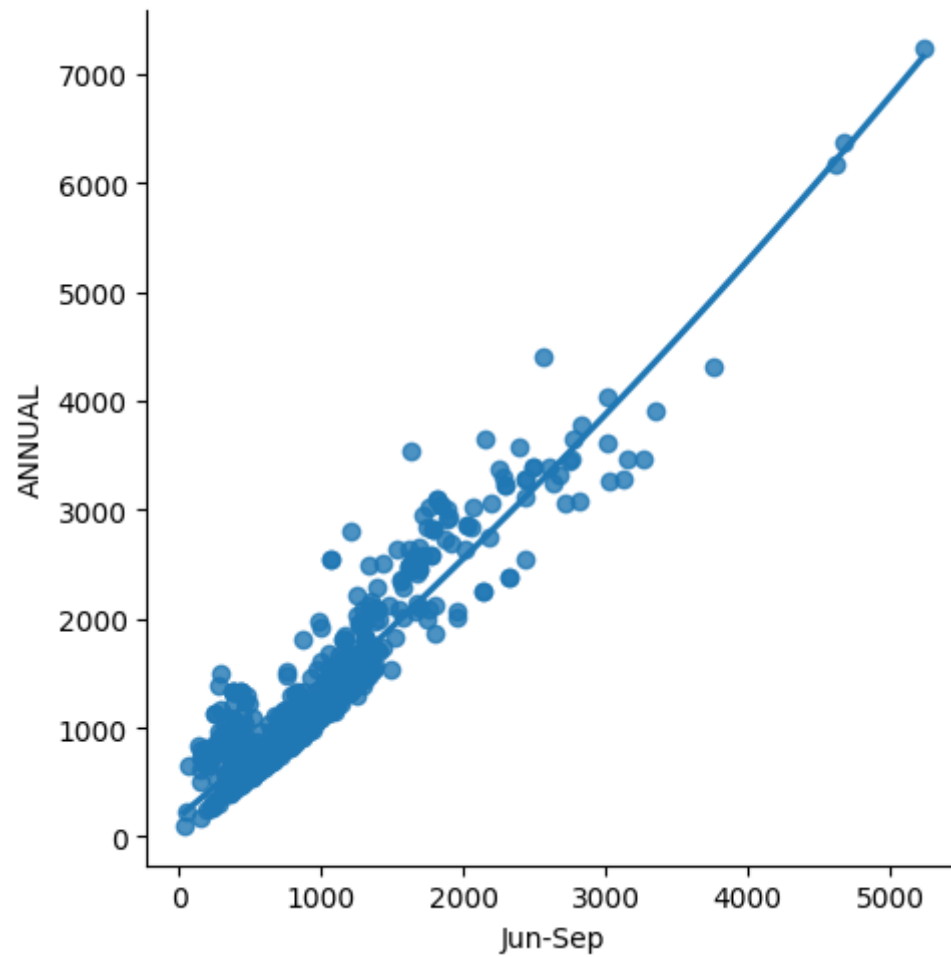
In [53]: `df1.describe()`

Out[53]:

	JAN	FEB	MAR	APR	MAY	JUN	JUL	AUG	SEP	OCT	
<b>count</b>	641.000000	641.000000	641.000000	641.000000	641.000000	641.000000	641.000000	641.000000	641.000000	641.000000	641.000000
<b>mean</b>	18.355070	20.984399	30.034789	45.543214	81.535101	196.007332	326.033697	291.152262	194.609048	90.446334	34.111111
<b>std</b>	21.082806	27.729596	45.451082	71.556279	111.960390	196.556284	221.364643	152.647325	99.830540	74.990685	59.371944
<b>min</b>	0.000000	0.000000	0.000000	0.000000	0.900000	3.800000	11.600000	14.100000	8.600000	3.100000	1.200000
<b>25%</b>	6.900000	7.000000	7.000000	5.000000	12.100000	68.800000	206.400000	194.600000	128.800000	34.300000	6.600000
<b>50%</b>	13.300000	12.300000	12.700000	15.100000	33.900000	131.900000	293.700000	284.800000	181.300000	62.600000	12.900000
<b>75%</b>	19.200000	24.100000	33.200000	48.300000	91.900000	226.600000	374.800000	358.100000	234.100000	130.200000	32.300000
<b>max</b>	144.500000	229.600000	367.900000	554.400000	733.700000	1476.200000	1820.900000	1522.100000	826.300000	517.700000	475.100000

## data visulization

```
In [54]: ▶ sns.lmplot(x="Jun-Sep",y="ANNUAL",order=2,data=df1,ci=None)  
plt.show()
```



```
In [55]: #taking selected columns from dataset  
df1=df1[['Jun-Sep', 'ANNUAL']]  
df1
```

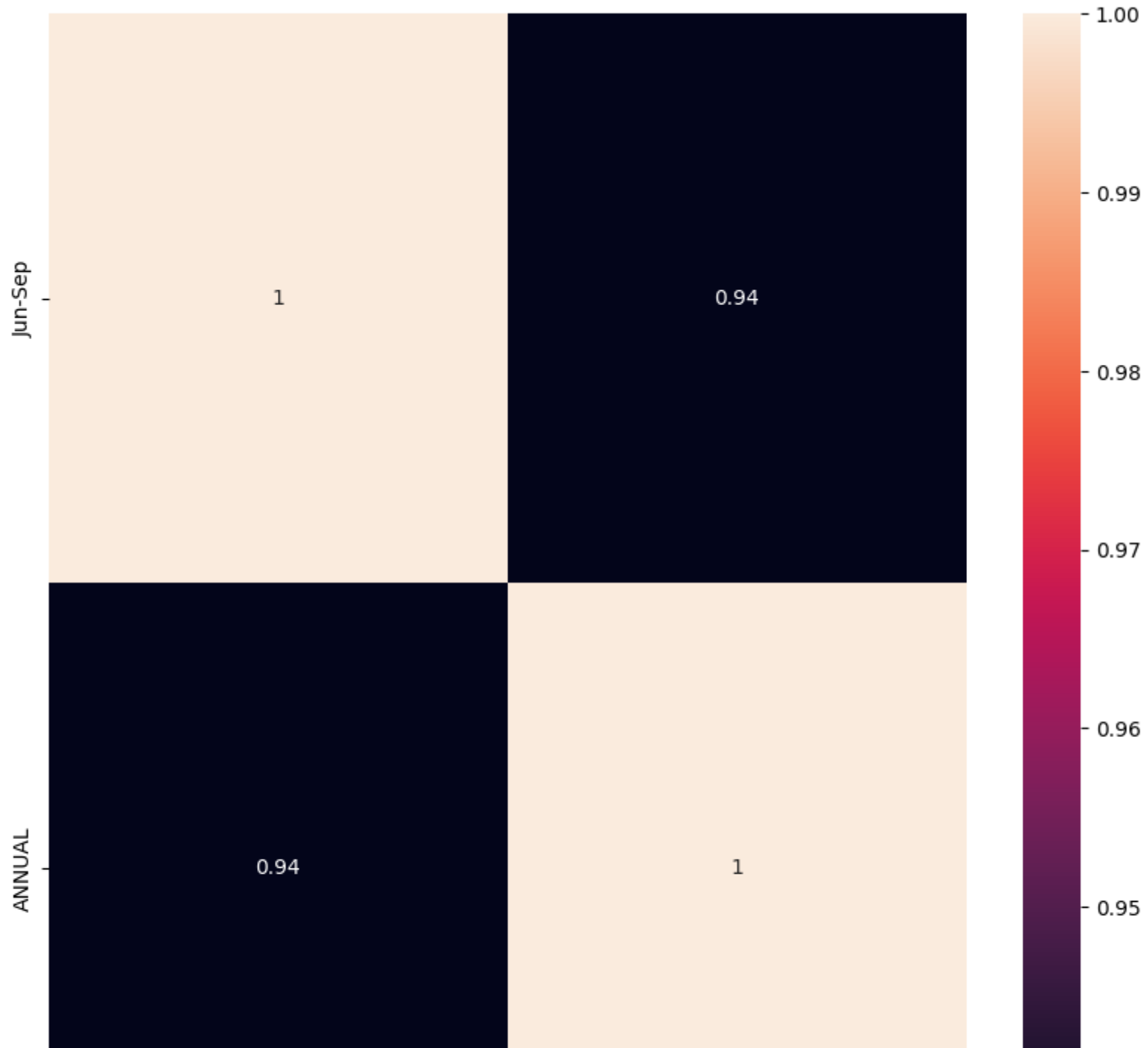
Out[55]:

	Jun-Sep	ANNUAL
0	1207.2	2805.2
1	1757.2	3015.7
2	1884.4	2913.3
3	1848.5	3043.8
4	3008.4	4034.7
...	...	...
636	2276.2	3302.5
637	3007.5	3621.6
638	1715.7	2958.4
639	2632.1	3253.1
640	998.5	1600.0

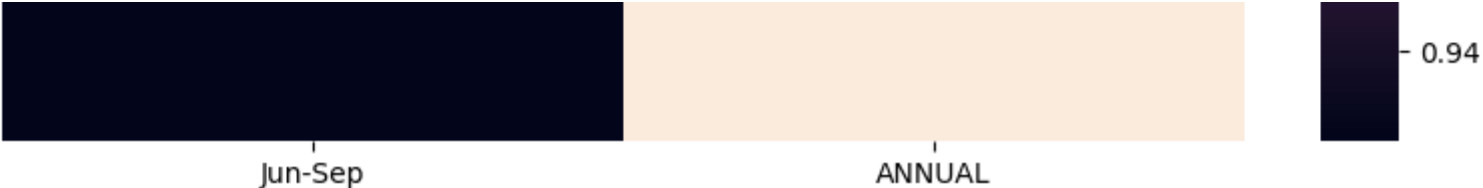
641 rows × 2 columns

```
In [56]: ▶ from sklearn.preprocessing import StandardScaler  
plt.figure(figsize=(10,10))  
sns.heatmap(df1.corr(),annot=True)  
plt.show()
```









```
In [58]: x=np.array(df1['Jun-Sep']).reshape(-1,1)
y=np.array(df1['ANNUAL']).reshape(-1,1)
#Dropping any rows with Nan values
df1.dropna(inplace=True)
df1
```

C:\Users\MY HOME\AppData\Local\Temp\ipykernel\_10092\178470992.py:4: 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) ([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))

```
df1.dropna(inplace=True)
```

Out[58]:

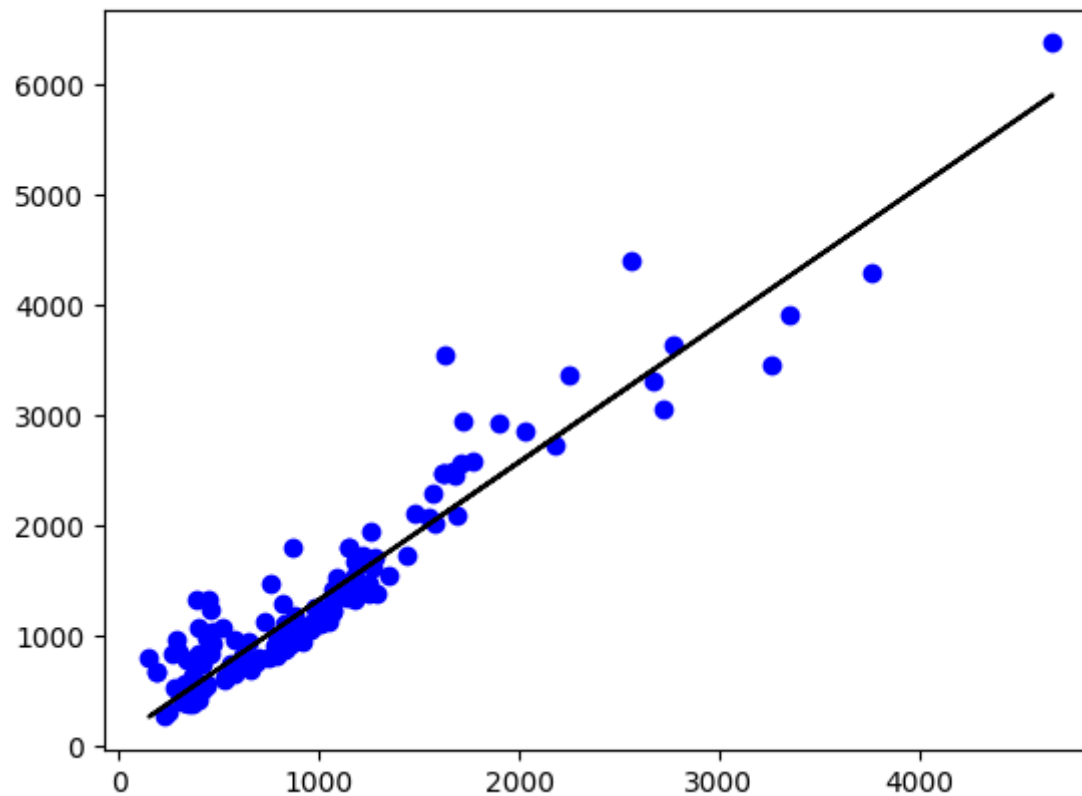
	Jun-Sep	ANNUAL
0	1207.2	2805.2
1	1757.2	3015.7
2	1884.4	2913.3
3	1848.5	3043.8
4	3008.4	4034.7
...	...	...
636	2276.2	3302.5
637	3007.5	3621.6
638	1715.7	2958.4
639	2632.1	3253.1
640	998.5	1600.0

641 rows × 2 columns

```
In [59]: 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(regr.score(x_test,y_test))
```

0.8931410883644894

```
In [60]: y_pred=regr.predict(x_test)
plt.scatter(x_test,y_test,color='b')
plt.plot(x_test,y_pred,color='k')
plt.show()
```



## ridge regression

```
In [61]: ▶ from sklearn.linear_model import Ridge, RidgeCV, Lasso
```

```
In [62]: ▶ ridge=Ridge(alpha=2)
ridge.fit(x_train,y_train)
train_score_ridge=ridge.score(x_train,y_train)
test_score_ridge=ridge.score(x_test,y_test)
print("\nRidgeRegression\n", (train_score_ridge))
print(test_score_ridge)
```

```
RidgeRegression
0.8698998531040942
0.8931410883440815
```

## Lasso Regression

```
In [63]: ▶ #Lasso regression model
print("\nLasso 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 0.8698998527066611

The test score for ls model is 0.8931410486981328

## CONCLUSION:

----->

Based on all models accuracies we conclude that the Linear REgression is the best model for this Dataset with the accuracy of 89%.