

# Importing Libraries

In [1]:

```
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
```

# Importing Datasets

In [2]:

```
df=pd.read_csv("rainfall_gujarat region.csv")
df
```

Out[2]:

	index	SUBDIVISION	YEAR	JAN	FEB	MAR	APR	MAY	JUN	JUL	AUG	SEP	OCT	NOV	12
0	2277	GUJARAT REGION	1901	4.2	0.0	0.6	1.6	7.0	60.3	240.2	205.4	18.1	16.6	0.0	
1	2278	GUJARAT REGION	1902	3.9	0.0	0.0	0.6	1.0	32.8	229.8	299.0	281.2	2.3	1.5	
2	2279	GUJARAT REGION	1903	0.3	0.1	1.4	0.0	12.3	30.1	452.9	202.0	183.2	5.4	0.0	
3	2280	GUJARAT REGION	1904	0.8	10.6	16.8	0.2	3.9	48.3	194.8	71.8	138.0	6.1	0.1	
4	2281	GUJARAT REGION	1905	0.1	0.7	1.1	0.3	0.0	20.1	668.3	37.9	81.3	1.4	0.2	
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	
110	2387	GUJARAT REGION	2011	0.0	0.2	0.0	0.0	0.0	16.3	259.2	451.7	162.5	0.4	0.0	
111	2388	GUJARAT REGION	2012	0.1	0.0	0.0	0.0	0.0	34.4	178.2	230.3	263.8	7.1	0.0	
112	2389	GUJARAT REGION	2013	0.0	0.9	0.1	4.6	0.0	155.7	405.4	211.1	287.3	53.2	0.1	
113	2390	GUJARAT REGION	2014	5.7	0.1	0.2	1.0	1.3	11.6	307.5	138.6	235.1	3.3	1.3	
114	2391	GUJARAT REGION	2015	1.8	0.0	6.1	5.5	0.9	120.7	354.7	37.4	93.4	2.2	0.3	

115 rows × 20 columns



## head

In [3]:

```
df.head(5)
df
```

Out[3]:

		index	SUBDIVISION	YEAR	JAN	FEB	MAR	APR	MAY	JUN	JUL	AUG	SEP	OCT	NOV	I
0	2277		GUJARAT REGION	1901	4.2	0.0	0.6	1.6	7.0	60.3	240.2	205.4	18.1	16.6	0.0	
1	2278		GUJARAT REGION	1902	3.9	0.0	0.0	0.6	1.0	32.8	229.8	299.0	281.2	2.3	1.5	
2	2279		GUJARAT REGION	1903	0.3	0.1	1.4	0.0	12.3	30.1	452.9	202.0	183.2	5.4	0.0	
3	2280		GUJARAT REGION	1904	0.8	10.6	16.8	0.2	3.9	48.3	194.8	71.8	138.0	6.1	0.1	
4	2281		GUJARAT REGION	1905	0.1	0.7	1.1	0.3	0.0	20.1	668.3	37.9	81.3	1.4	0.2	
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
110	2387		GUJARAT REGION	2011	0.0	0.2	0.0	0.0	0.0	16.3	259.2	451.7	162.5	0.4	0.0	
111	2388		GUJARAT REGION	2012	0.1	0.0	0.0	0.0	0.0	34.4	178.2	230.3	263.8	7.1	0.0	
112	2389		GUJARAT REGION	2013	0.0	0.9	0.1	4.6	0.0	155.7	405.4	211.1	287.3	53.2	0.1	
113	2390		GUJARAT REGION	2014	5.7	0.1	0.2	1.0	1.3	11.6	307.5	138.6	235.1	3.3	1.3	
114	2391		GUJARAT REGION	2015	1.8	0.0	6.1	5.5	0.9	120.7	354.7	37.4	93.4	2.2	0.3	

115 rows × 20 columns



## tail

In [4]:

```
df.tail(5)
df
```

Out[4]:

		index	SUBDIVISION	YEAR	JAN	FEB	MAR	APR	MAY	JUN	JUL	AUG	SEP	OCT	NOV	I
0	2277		GUJARAT REGION	1901	4.2	0.0	0.6	1.6	7.0	60.3	240.2	205.4	18.1	16.6	0.0	
1	2278		GUJARAT REGION	1902	3.9	0.0	0.0	0.6	1.0	32.8	229.8	299.0	281.2	2.3	1.5	
2	2279		GUJARAT REGION	1903	0.3	0.1	1.4	0.0	12.3	30.1	452.9	202.0	183.2	5.4	0.0	

	index	SUBDIVISION	YEAR	JAN	FEB	MAR	APR	MAY	JUN	JUL	AUG	SEP	OCT	NOV	12
3	2280	GUJARAT REGION	1904	0.8	10.6	16.8	0.2	3.9	48.3	194.8	71.8	138.0	6.1	0.1	
4	2281	GUJARAT REGION	1905	0.1	0.7	1.1	0.3	0.0	20.1	668.3	37.9	81.3	1.4	0.2	
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	
110	2387	GUJARAT REGION	2011	0.0	0.2	0.0	0.0	0.0	16.3	259.2	451.7	162.5	0.4	0.0	
111	2388	GUJARAT REGION	2012	0.1	0.0	0.0	0.0	0.0	34.4	178.2	230.3	263.8	7.1	0.0	
112	2389	GUJARAT REGION	2013	0.0	0.9	0.1	4.6	0.0	155.7	405.4	211.1	287.3	53.2	0.1	
113	2390	GUJARAT REGION	2014	5.7	0.1	0.2	1.0	1.3	11.6	307.5	138.6	235.1	3.3	1.3	
114	2391	GUJARAT REGION	2015	1.8	0.0	6.1	5.5	0.9	120.7	354.7	37.4	93.4	2.2	0.3	

115 rows × 20 columns

## Data Cleaning and Data Preprocessing

### describe()

In [5]:

```
df.describe()
```

Out[5]:

	index	YEAR	JAN	FEB	MAR	APR	MAY	JUN
count	115.000000	115.000000	115.000000	115.000000	115.000000	115.000000	115.000000	115.000000
mean	2334.000000	1958.000000	1.786087	1.191304	1.220870	1.116522	5.809565	121.284348
std	33.341666	33.341666	4.762590	2.870710	4.784102	3.980389	13.981353	84.287119
min	2277.000000	1901.000000	0.000000	0.000000	0.000000	0.000000	0.000000	2.600000
25%	2305.500000	1929.500000	0.000000	0.000000	0.000000	0.000000	0.100000	58.750000
50%	2334.000000	1958.000000	0.100000	0.000000	0.000000	0.100000	0.900000	112.500000
75%	2362.500000	1986.500000	1.500000	0.650000	0.250000	0.750000	4.100000	155.850000
max	2391.000000	2015.000000	44.100000	14.600000	42.100000	40.400000	98.300000	367.300000

### shape

```
In [6]: np.shape(df)
```

```
Out[6]: (115, 20)
```

## size

```
In [7]: np.size(df)
```

```
Out[7]: 2300
```

## dropna

```
In [8]: df=df.dropna()
```

## columns

```
In [9]: df.columns
```

```
Out[9]: Index(['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')
```

## info()

```
In [10]: df.info()
```

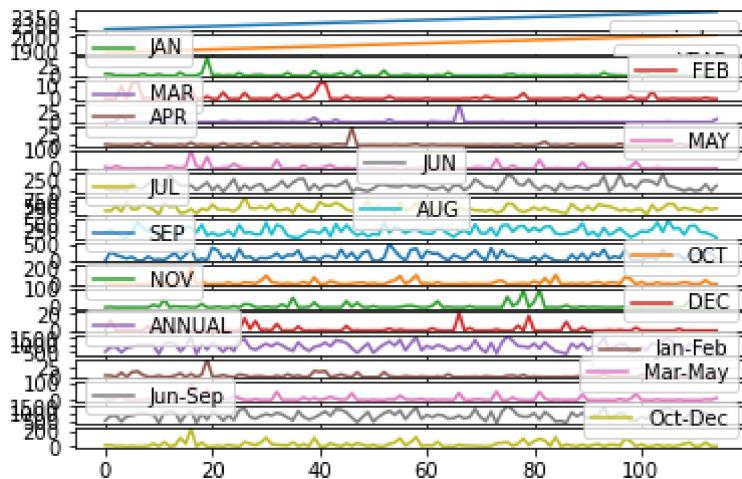
```
<class 'pandas.core.frame.DataFrame'>  
Int64Index: 115 entries, 0 to 114  
Data columns (total 20 columns):  
 #   Column      Non-Null Count  Dtype     
---  --          --          --  
 0   index       115 non-null    int64    
 1   SUBDIVISION 115 non-null    object   
 2   YEAR        115 non-null    int64    
 3   JAN         115 non-null    float64  
 4   FEB         115 non-null    float64  
 5   MAR         115 non-null    float64  
 6   APR         115 non-null    float64  
 7   MAY         115 non-null    float64  
 8   JUN         115 non-null    float64  
 9   JUL         115 non-null    float64  
 10  AUG         115 non-null    float64  
 11  SEP         115 non-null    float64  
 12  OCT         115 non-null    float64  
 13  NOV         115 non-null    float64  
 14  DEC         115 non-null    float64  
 15  ANNUAL      115 non-null    float64  
 16  Jan-Feb     115 non-null    float64
```

```
17 Mar-May      115 non-null    float64
18 Jun-Sep      115 non-null    float64
19 Oct-Dec      115 non-null    float64
dtypes: float64(17), int64(2), object(1)
memory usage: 18.9+ KB
```

## Line chart

```
In [11]: df.plot.line(subplots=True)
```

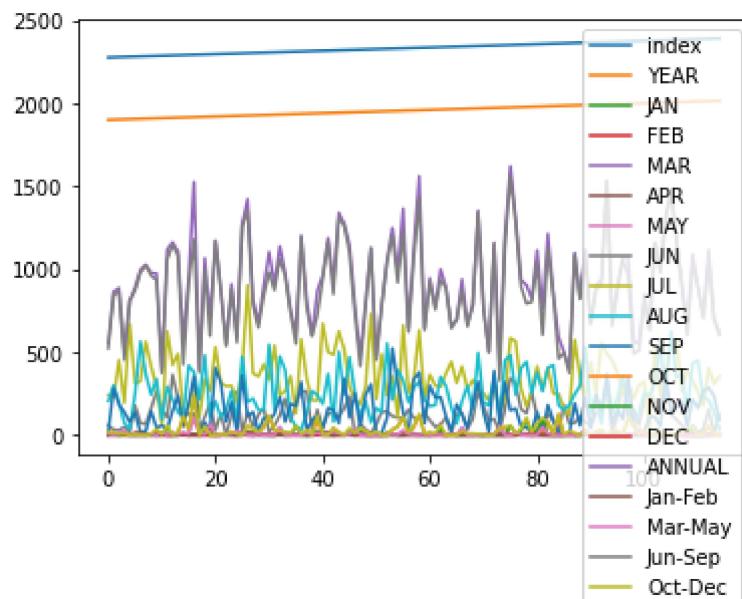
```
Out[11]: array([<AxesSubplot:>, <AxesSubplot:>, <AxesSubplot:>,
   <AxesSubplot:>, <AxesSubplot:>, <AxesSubplot:>, <AxesSubplot:>,
   <AxesSubplot:>, <AxesSubplot:>, <AxesSubplot:>, <AxesSubplot:>,
   <AxesSubplot:>, <AxesSubplot:>, <AxesSubplot:>, <AxesSubplot:>,
   <AxesSubplot:>], dtype=object)
```



## Line chart

```
In [12]: df.plot.line()
```

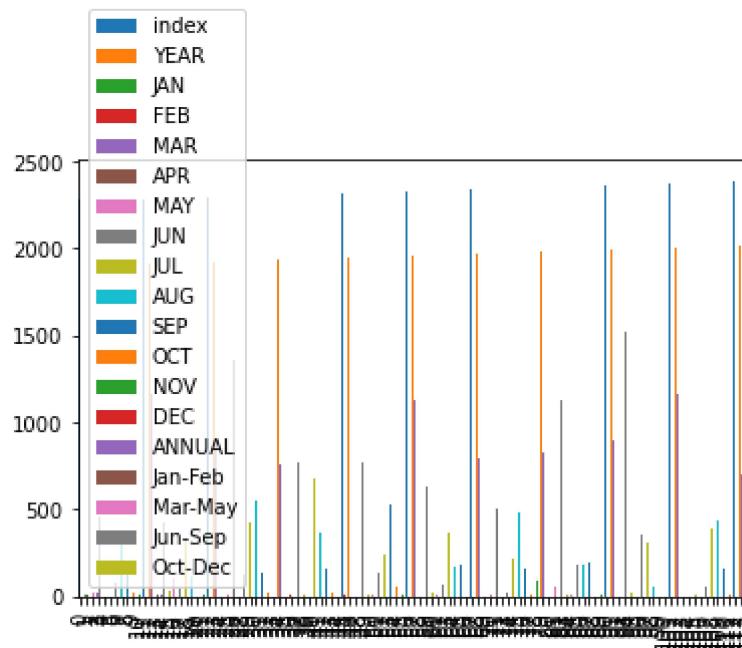
```
Out[12]: <AxesSubplot:>
```



## Bar chart

```
In [13]: df.plot.bar()
```

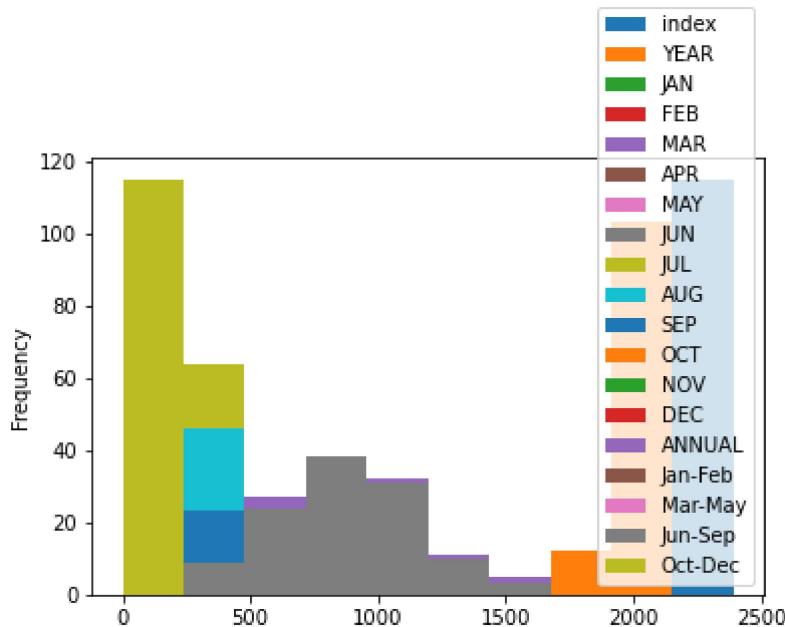
```
Out[13]: <AxesSubplot:>
```



## Histogram

```
In [14]: df.plot.hist()
```

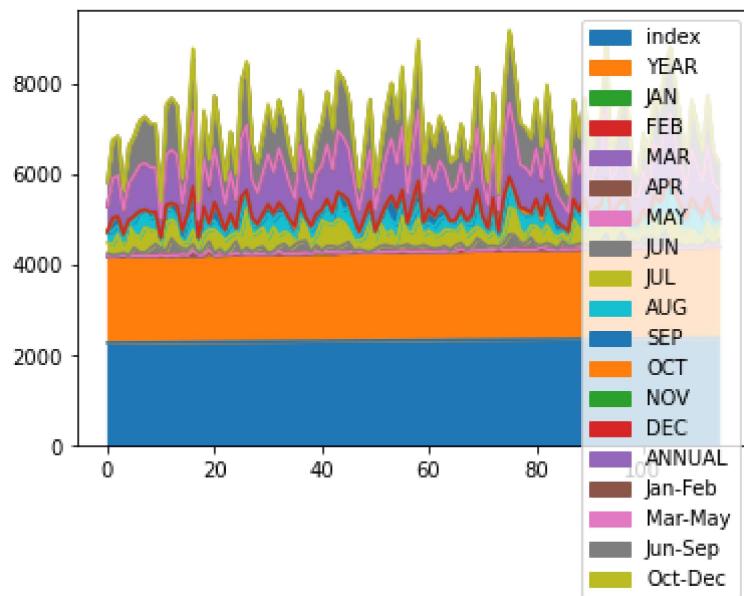
```
Out[14]: <AxesSubplot:ylabel='Frequency'>
```



## Area chart

In [15]: `df.plot.area()`

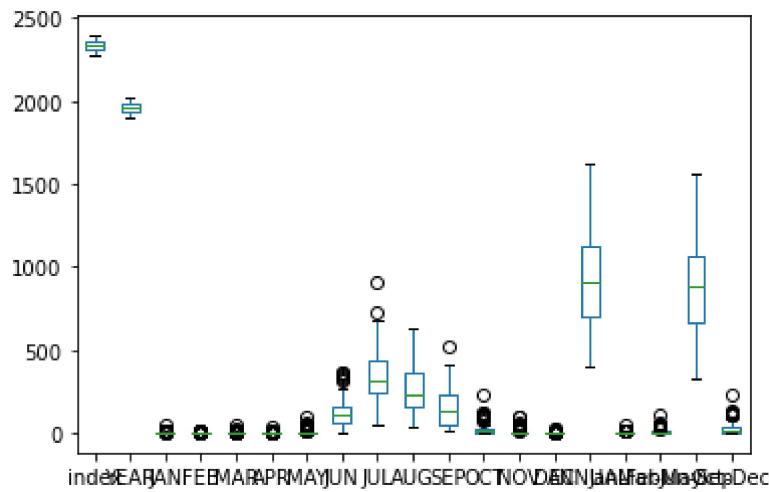
Out[15]: <AxesSubplot:>



## Box chart

In [16]: `df.plot.box()`

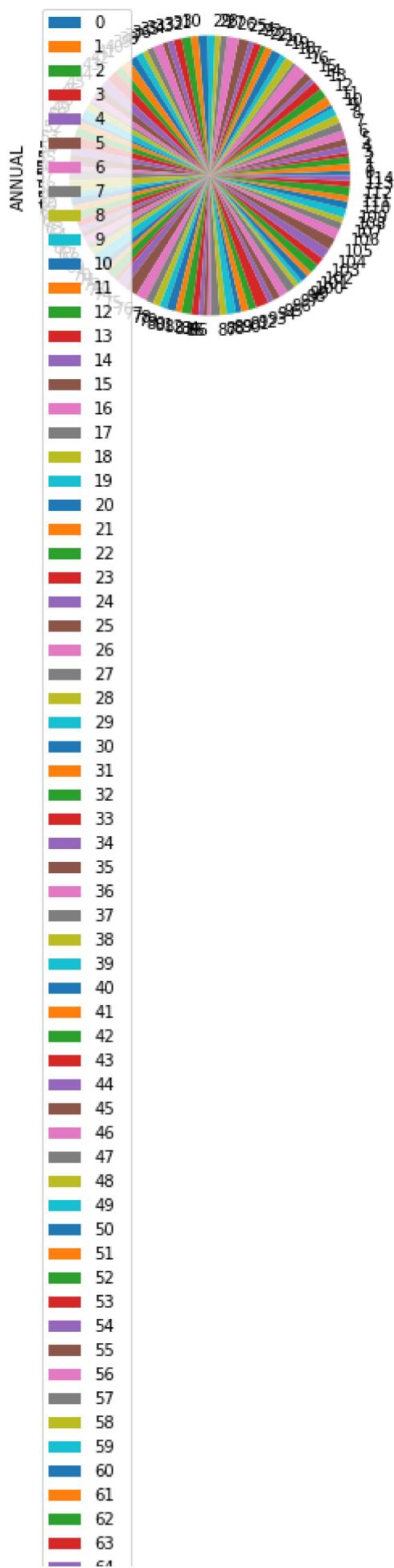
Out[16]: <AxesSubplot:>

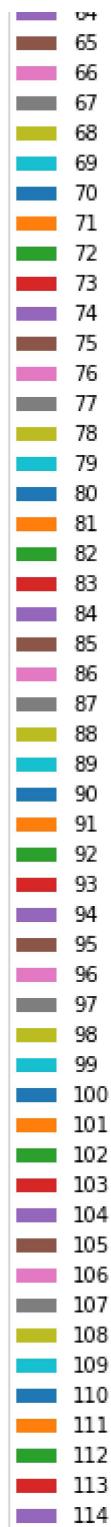


## Pie chart

```
In [17]: df.plot.pie(y='ANNUAL')
```

```
Out[17]: <AxesSubplot:ylabel='ANNUAL'>
```

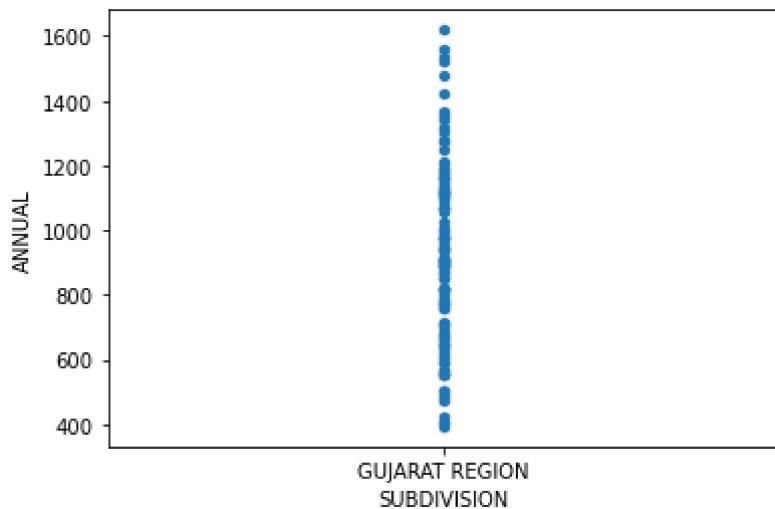




## Scatter chart

```
In [18]: df.plot.scatter(x='SUBDIVISION' ,y='ANNUAL')
```

```
Out[18]: <AxesSubplot:xlabel='SUBDIVISION', ylabel='ANNUAL'>
```



In [19]:

`df.info()`

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 115 entries, 0 to 114
Data columns (total 20 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   index       115 non-null    int64  
 1   SUBDIVISION 115 non-null    object  
 2   YEAR        115 non-null    int64  
 3   JAN         115 non-null    float64 
 4   FEB         115 non-null    float64 
 5   MAR         115 non-null    float64 
 6   APR         115 non-null    float64 
 7   MAY         115 non-null    float64 
 8   JUN         115 non-null    float64 
 9   JUL         115 non-null    float64 
 10  AUG         115 non-null    float64 
 11  SEP         115 non-null    float64 
 12  OCT         115 non-null    float64 
 13  NOV         115 non-null    float64 
 14  DEC         115 non-null    float64 
 15  ANNUAL      115 non-null    float64 
 16  Jan-Feb     115 non-null    float64 
 17  Mar-May     115 non-null    float64 
 18  Jun-Sep     115 non-null    float64 
 19  Oct-Dec     115 non-null    float64 
dtypes: float64(17), int64(2), object(1)
memory usage: 18.9+ KB
```

In [20]:

`df.describe()`

Out[20]:

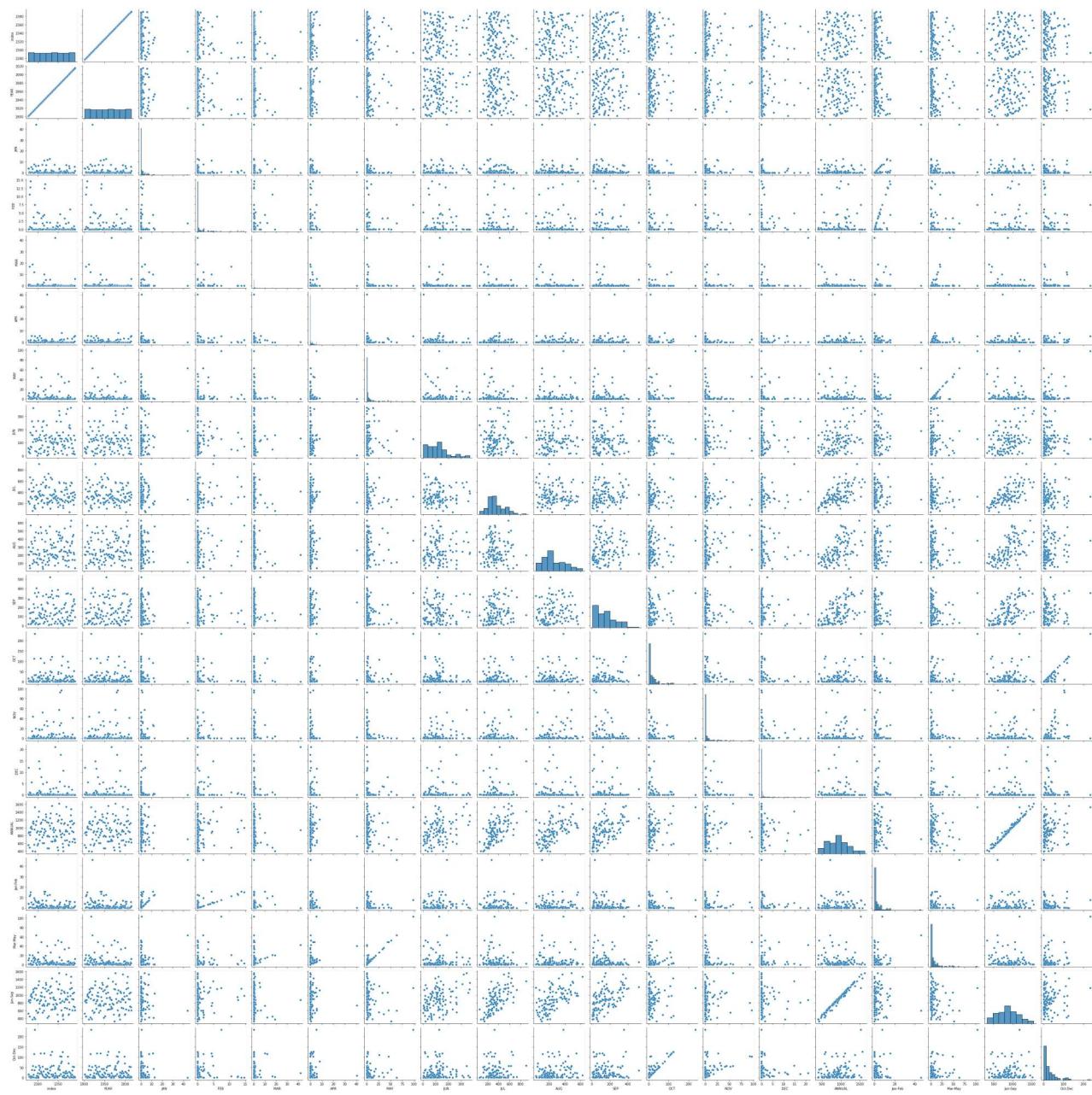
	index	YEAR	JAN	FEB	MAR	APR	MAY	JUN
<b>count</b>	115.000000	115.000000	115.000000	115.000000	115.000000	115.000000	115.000000	115.000000
<b>mean</b>	2334.000000	1958.000000	1.786087	1.191304	1.220870	1.116522	5.809565	121.284348
<b>std</b>	33.341666	33.341666	4.762590	2.870710	4.784102	3.980389	13.981353	84.287119
<b>min</b>	2277.000000	1901.000000	0.000000	0.000000	0.000000	0.000000	0.000000	2.600000
<b>25%</b>	2305.500000	1929.500000	0.000000	0.000000	0.000000	0.000000	0.100000	58.750000

	index	YEAR	JAN	FEB	MAR	APR	MAY	JUN
<b>50%</b>	2334.000000	1958.000000	0.100000	0.000000	0.000000	0.100000	0.900000	112.500000
<b>75%</b>	2362.500000	1986.500000	1.500000	0.650000	0.250000	0.750000	4.100000	155.850000
<b>max</b>	2391.000000	2015.000000	44.100000	14.600000	42.100000	40.400000	98.300000	367.300000

## EDA AND VISUALIZATION

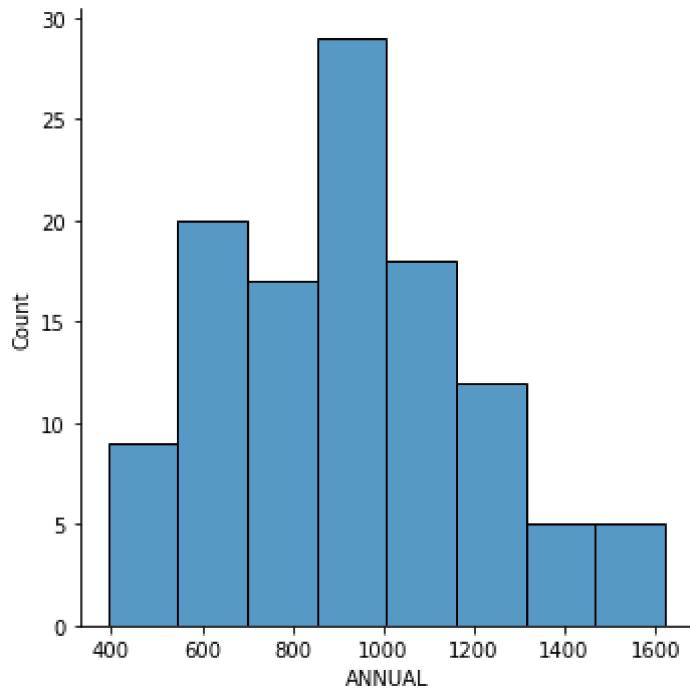
In [21]: `sns.pairplot(df)`

Out[21]: <seaborn.axisgrid.PairGrid at 0x1da997f8ac0>



In [22]: `sns.displot(df['ANNUAL'])`

Out[22]: &lt;seaborn.axisgrid.FacetGrid at 0x1daa52a2700&gt;

In [23]: 

```
sns.heatmap(df.corr())
```

Out[23]: &lt;AxesSubplot:&gt;

