

```
In [1]: import pandas as pd
import numpy as np
```

```
In [3]: data=pd.read_excel(r"D:\Sid 17-03-2025\SIDDHARTH BOSE\FSDS & GEN AI\March\27th -
```

```
In [5]: data
```

```
Out[5]:
```

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascience#\$	34 years	Mumbai	5^00#0	2+
1	Teddy^	Testing	45' yr	Bangalore	10%%000	<3
2	Uma#r	Dataanalyst^^#	NaN	NaN	1\$5%000	4> yrs
3	Jane	Ana^^lytics	NaN	Hyderbad	2000^0	NaN
4	Uttam*	Statistics	67-yr	NaN	30000-	5+ year
5	Kim	NLP	55yr	Delhi	6000^\$0	10+

```
In [7]: data.head()
```

```
Out[7]:
```

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascience#\$	34 years	Mumbai	5^00#0	2+
1	Teddy^	Testing	45' yr	Bangalore	10%%000	<3
2	Uma#r	Dataanalyst^^#	NaN	NaN	1\$5%000	4> yrs
3	Jane	Ana^^lytics	NaN	Hyderbad	2000^0	NaN
4	Uttam*	Statistics	67-yr	NaN	30000-	5+ year

```
In [9]: id(data)
```

```
Out[9]: 1802901236400
```

```
In [11]: data.columns
```

```
Out[11]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')
```

```
In [13]: data.shape
```

```
Out[13]: (6, 6)
```

```
In [15]: data.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6 entries, 0 to 5
Data columns (total 6 columns):
 #   Column      Non-Null Count  Dtype
---  -
 0   Name        6 non-null     object
 1   Domain      6 non-null     object
 2   Age         4 non-null     object
 3   Location    4 non-null     object
 4   Salary      6 non-null     object
 5   Exp         5 non-null     object
dtypes: object(6)
memory usage: 420.0+ bytes

```

In [17]: data

Out[17]:

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascience#\$	34 years	Mumbai	5^00#0	2+
1	Teddy^	Testing	45' yr	Bangalore	10%%000	<3
2	Uma#r	Dataanalyst^^#	NaN	NaN	1\$5%000	4> yrs
3	Jane	Ana^^lytics	NaN	Hyderbad	2000^0	NaN
4	Uttam*	Statistics	67-yr	NaN	30000-	5+ year
5	Kim	NLP	55yr	Delhi	6000^\$0	10+

In [19]: data.isnull()

Out[19]:

	Name	Domain	Age	Location	Salary	Exp
0	False	False	False	False	False	False
1	False	False	False	False	False	False
2	False	False	True	True	False	False
3	False	False	True	False	False	True
4	False	False	False	True	False	False
5	False	False	False	False	False	False

In [21]: data.isna()

Out[21]:

	Name	Domain	Age	Location	Salary	Exp
0	False	False	False	False	False	False
1	False	False	False	False	False	False
2	False	False	True	True	False	False
3	False	False	True	False	False	True
4	False	False	False	True	False	False
5	False	False	False	False	False	False

In [23]: `data.isnull().sum()`

Out[23]:

Name	0
Domain	0
Age	2
Location	2
Salary	0
Exp	1

dtype: int64

## Data Cleaning or Data Cleansing

In [26]: `data['Name']`

Out[26]:

0	Mike
1	Teddy^
2	Uma#r
3	Jane
4	Uttam*
5	Kim

Name: Name, dtype: object

In [30]: `data['Name']=data['Name'].str.replace(r'\W','',regex=True)`In [32]: `data`

Out[32]:

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascience#\$	34 years	Mumbai	5^00#0	2+
1	Teddy	Testing	45' yr	Bangalore	10%%000	<3
2	Umar	Dataanalyst^^#	NaN	NaN	1\$5%000	4> yrs
3	Jane	Ana^^lytics	NaN	Hyderbad	2000^0	NaN
4	Uttam	Statistics	67-yr	NaN	30000-	5+ year
5	Kim	NLP	55yr	Delhi	6000^\$0	10+

In [34]: `data['Domain']=data['Domain'].str.replace(r'\W','',regex=True)`  
`data`

Out[34]:

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascience	34 years	Mumbai	5^00#0	2+
1	Teddy	Testing	45' yr	Bangalore	10%%000	<3
2	Umar	Dataanalyst	NaN	NaN	1\$5%000	4> yrs
3	Jane	Analytics	NaN	Hyderbad	2000^0	NaN
4	Uttam	Statistics	67-yr	NaN	30000-	5+ year
5	Kim	NLP	55yr	Delhi	6000^\$0	10+

In [36]: `data['Age']=data['Age'].str.replace(r'\W','',regex=True)`  
`data`

Out[36]:

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascience	34years	Mumbai	5^00#0	2+
1	Teddy	Testing	45yr	Bangalore	10%%000	<3
2	Umar	Dataanalyst	NaN	NaN	1\$5%000	4> yrs
3	Jane	Analytics	NaN	Hyderbad	2000^0	NaN
4	Uttam	Statistics	67yr	NaN	30000-	5+ year
5	Kim	NLP	55yr	Delhi	6000^\$0	10+

In [38]: `data['Age']=data['Age'].str.extract(r'(\d+)')`

In [40]: `data`

Out[40]:

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascience	34	Mumbai	5^00#0	2+
1	Teddy	Testing	45	Bangalore	10%%000	<3
2	Umar	Dataanalyst	NaN	NaN	1\$5%000	4> yrs
3	Jane	Analytics	NaN	Hyderbad	2000^0	NaN
4	Uttam	Statistics	67	NaN	30000-	5+ year
5	Kim	NLP	55	Delhi	6000^\$0	10+

In [42]: `data['Salary']=data['Salary'].str.replace(r"W","",regex=True)`  
`data`

Out[42]:

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascience	34	Mumbai	5000	2+
1	Teddy	Testing	45	Bangalore	10000	<3
2	Umar	Dataanalyst	NaN	NaN	15000	4> yrs
3	Jane	Analytics	NaN	Hyderabad	20000	NaN
4	Uttam	Statistics	67	NaN	30000	5+ year
5	Kim	NLP	55	Delhi	60000	10+

In [44]: `data['Location']=data['Location'].str.replace(r'\W','',regex=True)`

In [46]: `data`

Out[46]:

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascience	34	Mumbai	5000	2+
1	Teddy	Testing	45	Bangalore	10000	<3
2	Umar	Dataanalyst	NaN	NaN	15000	4> yrs
3	Jane	Analytics	NaN	Hyderabad	20000	NaN
4	Uttam	Statistics	67	NaN	30000	5+ year
5	Kim	NLP	55	Delhi	60000	10+

In [48]: `data['Exp']=data['Exp'].str.extract(r'(\d+)')`  
`data`

Out[48]:

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascience	34	Mumbai	5000	2
1	Teddy	Testing	45	Bangalore	10000	3
2	Umar	Dataanalyst	NaN	NaN	15000	4
3	Jane	Analytics	NaN	Hyderabad	20000	NaN
4	Uttam	Statistics	67	NaN	30000	5
5	Kim	NLP	55	Delhi	60000	10

till now we have raw data we use regex to clean all nosiy data ad cancel from main data

## Apply EDA Technique

In [102...]: `clean_data=data.copy()`

In [104...

clean\_data

Out[104...

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascience	34	Mumbai	5000	2
1	Teddy	Testing	45	Bangalore	10000	3
2	Umar	Dataanalyst	NaN	NaN	15000	4
3	Jane	Analytics	NaN	Hyderbad	20000	NaN
4	Uttam	Statistics	67	NaN	30000	5
5	Kim	NLP	55	Delhi	60000	10

In [106...

```
clean_data['Age']=clean_data['Age'].fillna(np.mean(pd.to_numeric(clean_data['Age
```

In [108...

clean\_data['Age']

Out[108...

```
0      34
1      45
2    50.25
3    50.25
4      67
5      55
Name: Age, dtype: object
```

In [110...

clean\_data['Exp']

Out[110...

```
0      2
1      3
2      4
3    NaN
4      5
5     10
Name: Exp, dtype: object
```

In [112...

```
clean_data['Exp']=clean_data['Exp'].fillna(np.mean(pd.to_numeric(clean_data['Exp
```

In [114...

clean\_data['Exp']

Out[114...

```
0      2
1      3
2      4
3    4.8
4      5
5     10
Name: Exp, dtype: object
```

In [116...

clean\_data

Out[116...

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascience	34	Mumbai	5000	2
1	Teddy	Testing	45	Bangalore	10000	3
2	Umar	Dataanalyst	50.25	NaN	15000	4
3	Jane	Analytics	50.25	Hyderbad	20000	4.8
4	Uttam	Statistics	67	NaN	30000	5
5	Kim	NLP	55	Delhi	60000	10

In [118...

```
clean_data['Location'].isnull().sum()
```

Out[118...

2

In [120...

```
clean_data['Location']=clean_data['Location'].fillna(clean_data['Location'].mode
```

In [122...

```
clean_data['Location']
```

Out[122...

```
0      Mumbai
1    Bangalore
2    Bangalore
3    Hyderbad
4    Bangalore
5         Delhi
Name: Location, dtype: object
```

In [124...

```
clean_data
```

Out[124...

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascience	34	Mumbai	5000	2
1	Teddy	Testing	45	Bangalore	10000	3
2	Umar	Dataanalyst	50.25	Bangalore	15000	4
3	Jane	Analytics	50.25	Hyderbad	20000	4.8
4	Uttam	Statistics	67	Bangalore	30000	5
5	Kim	NLP	55	Delhi	60000	10

In [126...

```
clean_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6 entries, 0 to 5
Data columns (total 6 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   Name        6 non-null      object
1   Domain      6 non-null      object
2   Age         6 non-null      object
3   Location    6 non-null      object
4   Salary      6 non-null      object
5   Exp         6 non-null      object
dtypes: object(6)
memory usage: 420.0+ bytes
```

```
In [128... clean_data['Age']=clean_data['Age'].astype(int)
```

```
In [130... clean_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6 entries, 0 to 5
Data columns (total 6 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   Name        6 non-null      object
1   Domain      6 non-null      object
2   Age         6 non-null      int32
3   Location    6 non-null      object
4   Salary      6 non-null      object
5   Exp         6 non-null      object
dtypes: int32(1), object(5)
memory usage: 396.0+ bytes
```

```
In [131... clean_data['Salary']=clean_data['Salary'].astype(int)
clean_data['Exp']=clean_data['Exp'].astype(int)
```

```
In [133... clean_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6 entries, 0 to 5
Data columns (total 6 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   Name        6 non-null      object
1   Domain      6 non-null      object
2   Age         6 non-null      int32
3   Location    6 non-null      object
4   Salary      6 non-null      int32
5   Exp         6 non-null      int32
dtypes: int32(3), object(3)
memory usage: 348.0+ bytes
```

```
In [136... clean_data['Name']= clean_data['Name'].astype('category')
clean_data['Domain']= clean_data['Domain'].astype('category')
clean_data['Location']= clean_data['Location'].astype('category')
```

```
In [138... clean_data.info()
```



```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6 entries, 0 to 5
Data columns (total 6 columns):
#   Column      Non-Null Count  Dtype
---  -
0   Name        6 non-null     category
1   Domain      6 non-null     category
2   Age         6 non-null     int32
3   Location    6 non-null     category
4   Salary      6 non-null     int32
5   Exp         6 non-null     int32
dtypes: category(3), int32(3)
memory usage: 866.0 bytes
```

```
In [146... clean_data.to_csv('Clean_data.csv')
```

```
In [150... import os
os.getcwd()
```

```
Out[150... 'C:\\Users\\siddharth.bose'
```

## Visualizing the Data using Matplotlib and Seaborn

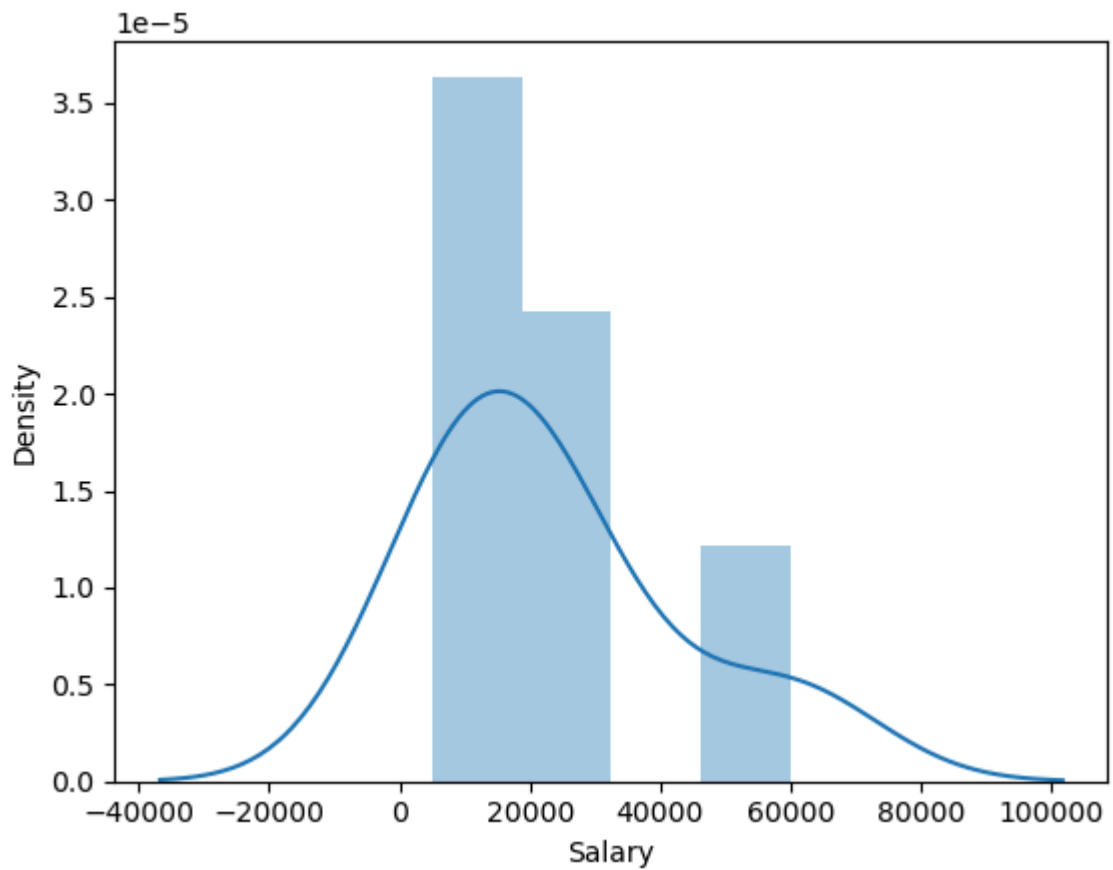
```
In [153... import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [157... import warnings
warnings.filterwarnings('ignore')
```

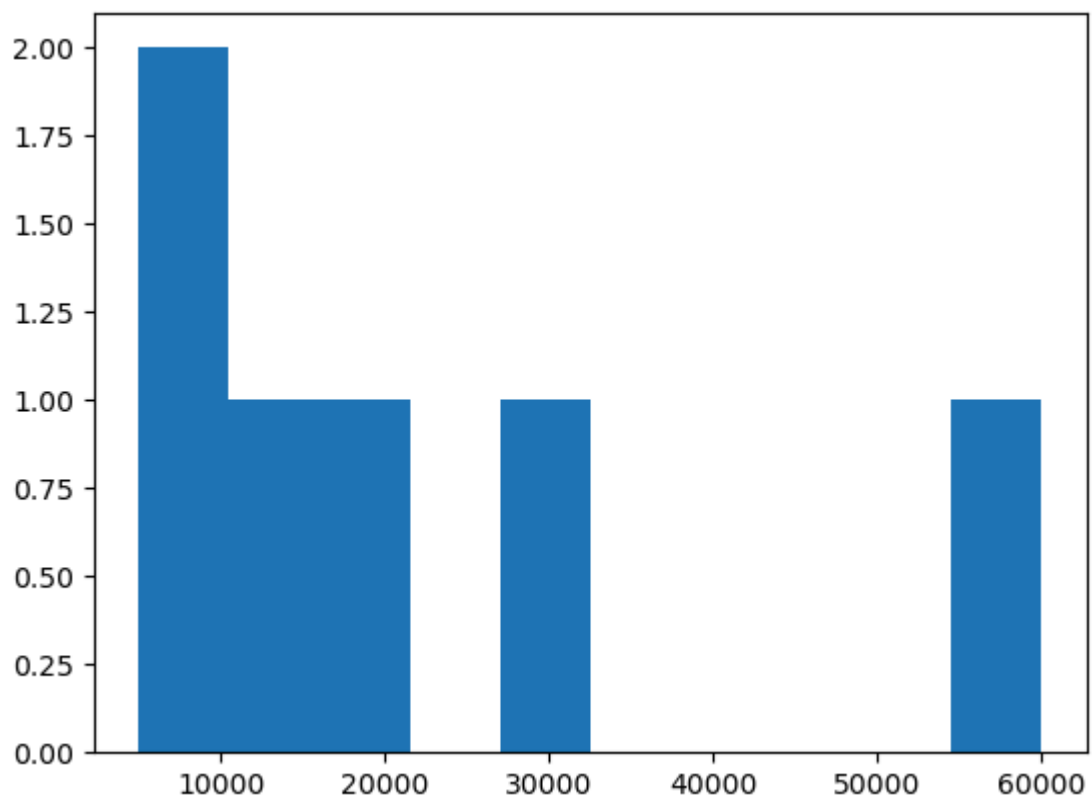
```
In [159... clean_data['Salary']
```

```
Out[159... 0      5000
1     10000
2     15000
3     20000
4     30000
5     60000
Name: Salary, dtype: int32
```

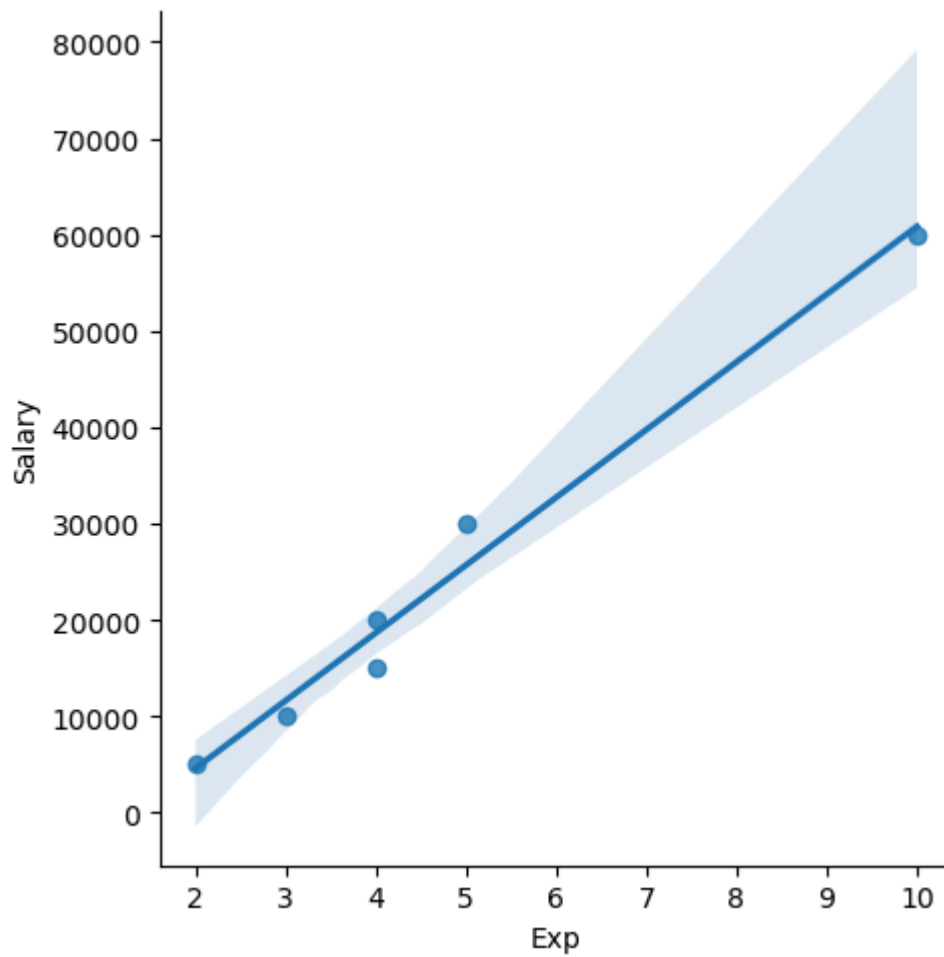
```
In [163... vis1= sns.distplot(clean_data['Salary']) #Univariate Analysis
plt.show()
```



```
In [165... vis2 = plt.hist(clean_data['Salary']) #univariate Analysis
plt.show()
```

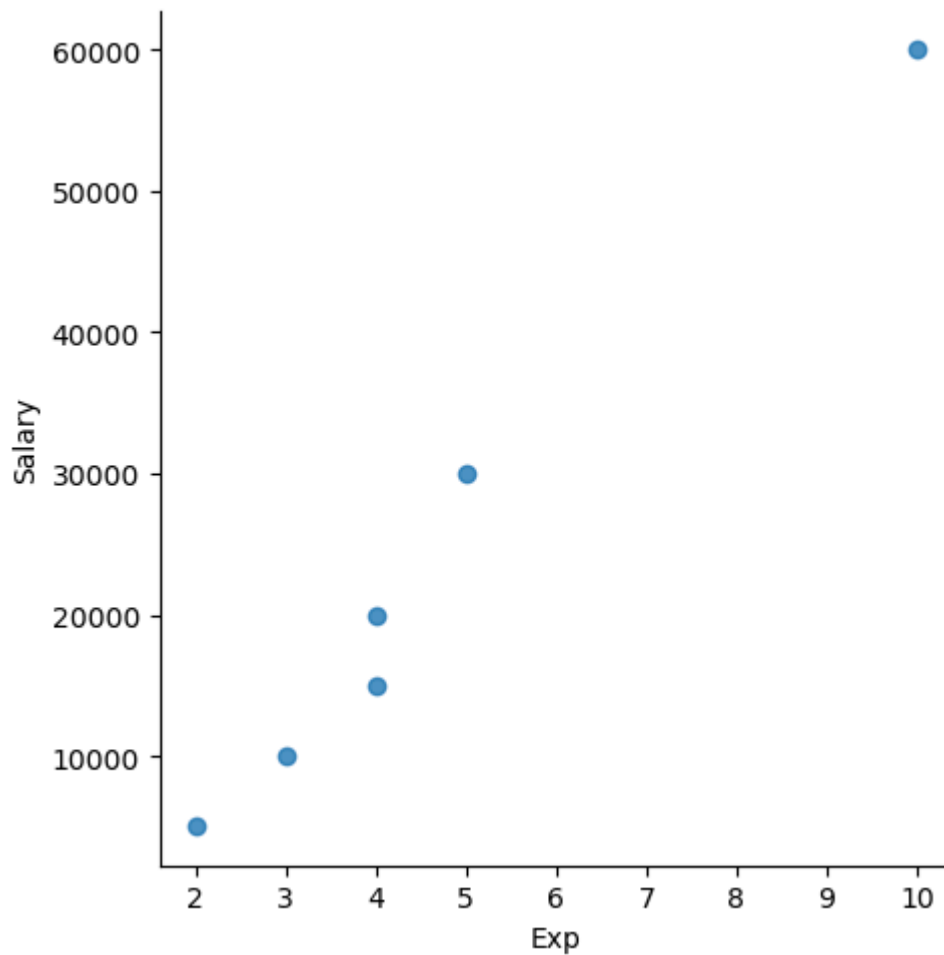


```
In [167... vis3 = sns.lmplot(data = clean_data,x = 'Exp',y = 'Salary') #Bivariate Analysis
plt.show()
```



In [169...

```
vis4 = sns.lmplot(data= clean_data,x='Exp',y = 'Salary',fit_reg=False) #Bivariate  
plt.show()
```



In [171...] `clean_data[:]`

Out[171...] 

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascience	34	Mumbai	5000	2
1	Teddy	Testing	45	Bangalore	10000	3
2	Umar	Dataanalyst	50	Bangalore	15000	4
3	Jane	Analytics	50	Hyderbad	20000	4
4	Uttam	Statistics	67	Bangalore	30000	5
5	Kim	NLP	55	Delhi	60000	10

In [173...] `clean_data[0:6:2]`

Out[173...] 

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascience	34	Mumbai	5000	2
2	Umar	Dataanalyst	50	Bangalore	15000	4
4	Uttam	Statistics	67	Bangalore	30000	5

In [175...] `clean_data.columns`

Out[175...] `Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')`

```
In [177... X_iv = clean_data[['Name', 'Domain', 'Age', 'Location', 'Exp']] #Extracting the
```

```
In [179... X_iv
```

```
Out[179...      Name    Domain  Age  Location  Exp
0   Mike  Datascience   34   Mumbai    2
1  Teddy    Testing   45  Bangalore    3
2   Umar  Dataanalyst   50  Bangalore    4
3   Jane    Analytics   50  Hyderabad    4
4  Uttam   Statistics   67  Bangalore    5
5    Kim        NLP    55    Delhi   10
```

```
In [181... y_dv = clean_data[['Salary']]
y_dv
```

```
Out[181...      Salary
0     5000
1    10000
2    15000
3    20000
4    30000
5    60000
```

```
In [183... clean_data
```

```
Out[183...      Name    Domain  Age  Location  Salary  Exp
0   Mike  Datascience   34   Mumbai    5000    2
1  Teddy    Testing   45  Bangalore   10000    3
2   Umar  Dataanalyst   50  Bangalore   15000    4
3   Jane    Analytics   50  Hyderabad   20000    4
4  Uttam   Statistics   67  Bangalore   30000    5
5    Kim        NLP    55    Delhi   60000   10
```

```
In [185... X_iv
```

Out[185...

	Name	Domain	Age	Location	Exp
0	Mike	Datascience	34	Mumbai	2
1	Teddy	Testing	45	Bangalore	3
2	Umar	Dataanalyst	50	Bangalore	4
3	Jane	Analytics	50	Hyderabad	4
4	Uttam	Statistics	67	Bangalore	5
5	Kim	NLP	55	Delhi	10

In [187...

y\_dv

Out[187...

	Salary
0	5000
1	10000
2	15000
3	20000
4	30000
5	60000

In [189...

```
imputation = pd.get_dummies(clean_data)
imputation
```

Out[189...

	Age	Salary	Exp	Name_Jane	Name_Kim	Name_Mike	Name_Teddy	Name_Umar
0	34	5000	2	False	False	True	False	False
1	45	10000	3	False	False	False	True	False
2	50	15000	4	False	False	False	False	True
3	50	20000	4	True	False	False	False	False
4	67	30000	5	False	False	False	False	False
5	55	60000	10	False	True	False	False	False



In [ ]: