

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

```
In [2]: pd.__version__
```

```
Out[2]: '2.2.2'
```

```
# pip install --upgrade openpyxl
```

```
In [3]: emp=pd.read_excel(r'D:\EDA\Rawdata.xlsx')
emp
```

	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 [4]: id(emp)
```

```
Out[4]: 2025608588096
```

```
In [5]: emp.columns
```

```
Out[5]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')
```

```
In [6]: emp.shape
```

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

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

	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 [8]: emp.tail()
```

Out[8]:

	Name	Domain	Age	Location	Salary	Exp
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 [9]:

emp.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 [10]:

emp.isnull()

Out[10]:

	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 [11]:

emp.isnull().sum()

Out[11]:

```

Name          0
Domain        0
Age           2
Location       2
Salary         0
Exp            1
dtype: int64

```

## Data Cleaning or Data Cleansing

In [12]: emp

	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 [13]: emp['Name']

```
Out[13]: 0      Mike
1      Teddy^
2      Uma#r
3      Jane
4      Uttam*
5      Kim
Name: Name, dtype: object
```

In [14]: emp['Name']=emp['Name'].str.replace(r'\W',' ',regex=True) # \W= non word character ,

In [15]: emp['Name']

```
Out[15]: 0      Mike
1      Teddy
2      Umar
3      Jane
4      Uttam
5      Kim
Name: Name, dtype: object
```

In [16]: emp

	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 [17]: emp['Domain']=emp['Domain'].str.replace(r'\W',' ',regex=True)
```

```
In [18]: emp['Domain']
```

```
Out[18]: 0    Datascienc
          1        Testing
          2   Dataanalyst
          3     Analytics
          4   Statistics
          5         NLP
Name: Domain, dtype: object
```

```
In [19]: emp
```

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascienc	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	Hyderabad	2000^0	NaN
4	Uttam	Statistics	67-yr	NaN	30000-	5+ year
5	Kim	NLP	55yr	Delhi	6000^\$0	10+

```
In [20]: emp['Age']=emp['Age'].str.replace(r'\W',' ',regex=True)
emp['Age']
```

```
Out[20]: 0    34years
          1      45yr
          2      NaN
          3      NaN
          4      67yr
          5      55yr
Name: Age, dtype: object
```

```
In [21]: emp['Age']=emp['Age'].str.extract('(\d+)') # r(r'(\d+)') -if we get error , ext
emp['Age']
```

```
Out[21]: 0    34
          1    45
          2    NaN
          3    NaN
          4    67
          5    55
Name: Age, dtype: object
```

```
In [22]: emp
```

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascienc	34	Mumbai	5^00#0	2+
1	Teddy	Testing	45	Bangalore	10%0000	<3
2	Umar	Dataanalyst	NaN	NaN	1\$5%0000	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 [23]: emp['Location']=emp['Location'].str.replace(r'\W',' ',regex=True)
emp['Location']
```

```
Out[23]: 0      Mumbai
         1      Bangalore
         2        NaN
         3      Hyderbad
         4        NaN
         5      Delhi
Name: Location, dtype: object
```

```
In [24]: emp['Salary']=emp['Salary'].str.replace(r'\W',' ',regex=True) # .replace-only number
emp['Salary']
```

```
Out[24]: 0      5000
         1     10000
         2    15000
         3    20000
         4    30000
         5    60000
Name: Salary, dtype: object
```

```
In [25]: emp['Exp']=emp['Exp'].str.extract('(\d+)')
emp['Exp']
```

```
Out[25]: 0      2
         1      3
         2      4
         3    NaN
         4      5
         5     10
Name: Exp, dtype: object
```

```
In [26]: emp
```

Out[26]:

	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 [27]: `clean_data=emp.copy()`In [28]: `clean_data`

Out[28]:

	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

## EDA Techniques

In [29]: `clean_data.isnull().sum()`

Out[29]:

In [30]: `clean_data['Age']`

Out[30]:

```
In [31]: import numpy as np
```

```
In [32]: clean_data['Age']=clean_data['Age'].fillna(np.mean(pd.to_numeric(clean_data['Age'])))
```

```
In [33]: clean_data['Age']
```

```
Out[33]: 0      34  
1      45  
2    50.25  
3    50.25  
4      67  
5      55  
Name: Age, dtype: object
```

```
In [34]: clean_data['Exp']
```

```
Out[34]: 0      2  
1      3  
2      4  
3    NaN  
4      5  
5     10  
Name: Exp, dtype: object
```

```
In [35]: clean_data['Exp']=clean_data['Exp'].fillna(np.mean(pd.to_numeric(clean_data['Exp'])))  
clean_data['Exp']
```

```
Out[35]: 0      2  
1      3  
2      4  
3    4.8  
4      5  
5     10  
Name: Exp, dtype: object
```

```
In [36]: clean_data['Location']
```

```
Out[36]: 0      Mumbai  
1    Bangalore  
2      NaN  
3    Hyderabad  
4      NaN  
5      Delhi  
Name: Location, dtype: object
```

```
In [37]: clean_data['Location']=clean_data['Location'].fillna(clean_data['Location'].mode()[  
clean_data['Location']]
```

```
Out[37]: 0      Mumbai  
1    Bangalore  
2    Bangalore  
3    Hyderabad  
4    Bangalore  
5      Delhi  
Name: Location, dtype: object
```

In [38]: `clean_data`

	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	Hyderabad	20000	4.8
4	Uttam	Statistics	67	Bangalore	30000	5
5	Kim	NLP	55	Delhi	60000	10

In [39]: `emp.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 [40]: `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 [41]: `clean_data['Age']=clean_data['Age'].astype(int) # convert system build in data`In [42]: `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 [43]: clean_data['Salary']=clean_data['Salary'].astype(int)
```

```
In [44]: clean_data['Exp']=clean_data['Exp'].astype(int)
```

```
In [45]: 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 [55]: 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 [56]: 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 [49]: clean_data
```

Out[49]:

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascienc	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 [50]: `clean_data.to_csv('clean_data.csv')`

In [51]: `import os  
os.getcwd() # from the os give the saved current working directly`

Out[51]: 'C:\\Users\\DELL'

## visualization

In [52]: `import matplotlib.pyplot as plt  
import seaborn as sns`

In [53]: `import warnings  
warnings.filterwarnings('ignore')`

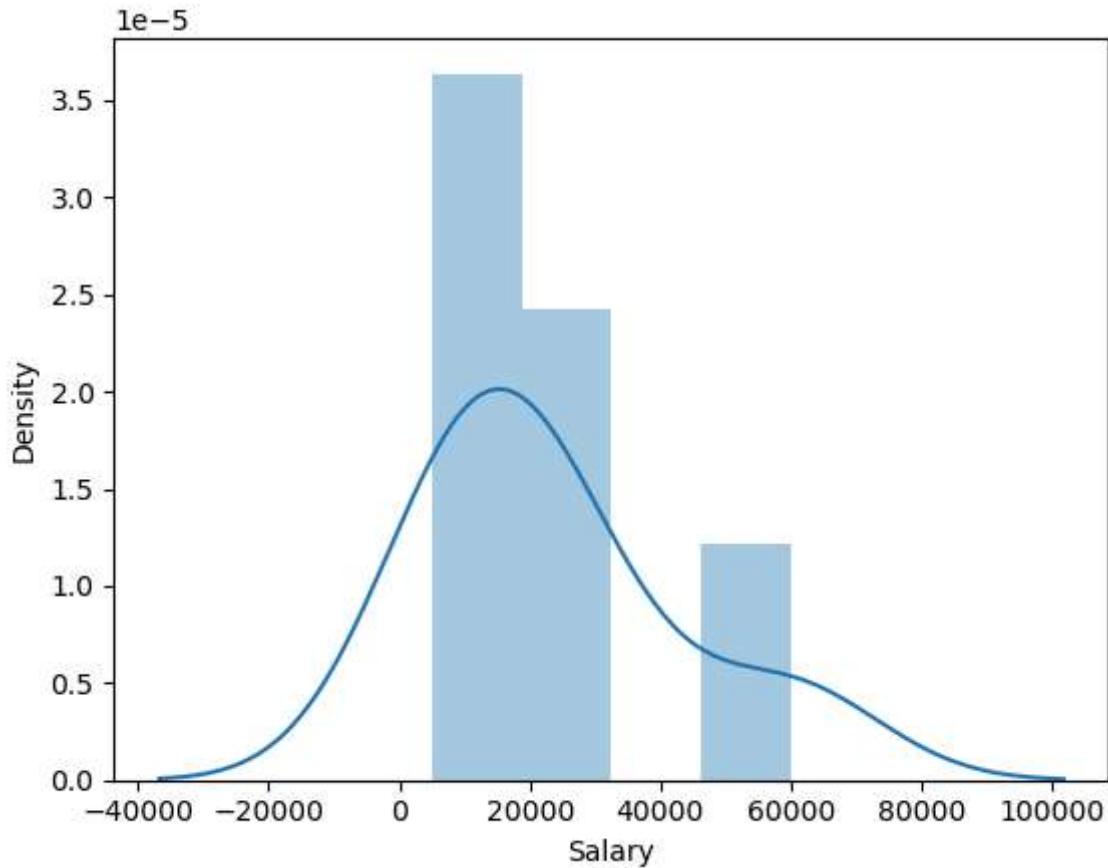
In [57]: `clean_data['Salary']`

Out[57]:

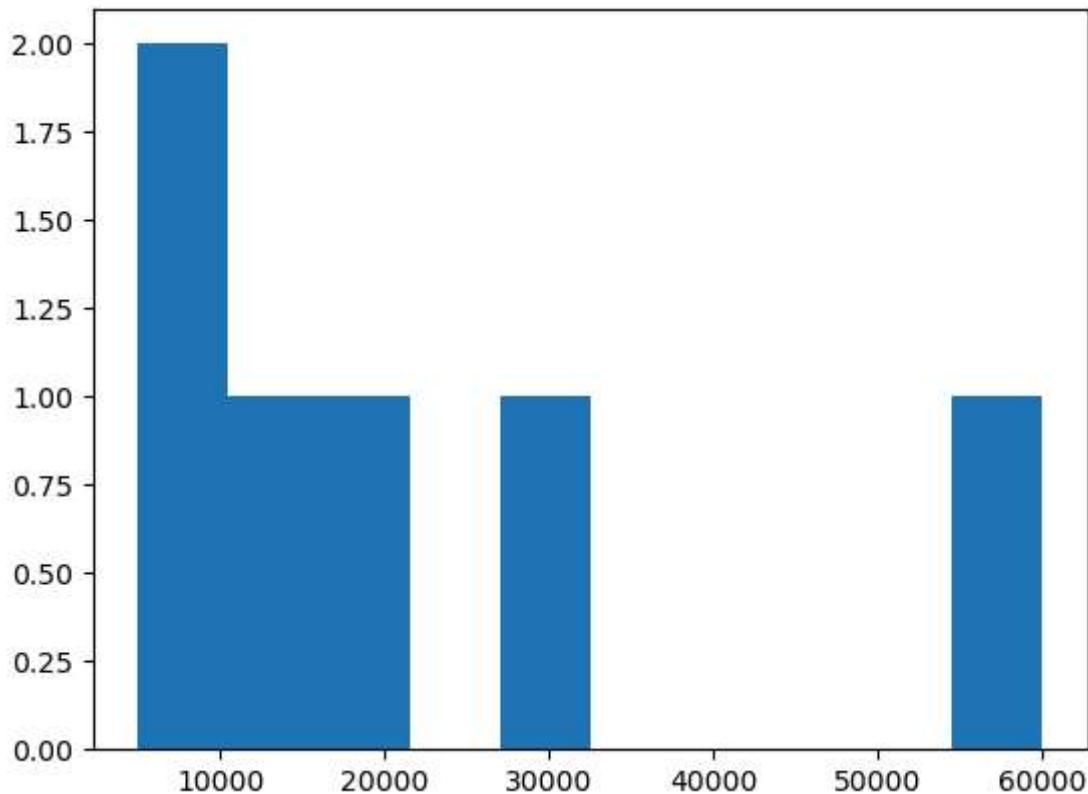
0	5000
1	10000
2	15000
3	20000
4	30000
5	60000

Name: Salary, dtype: int32

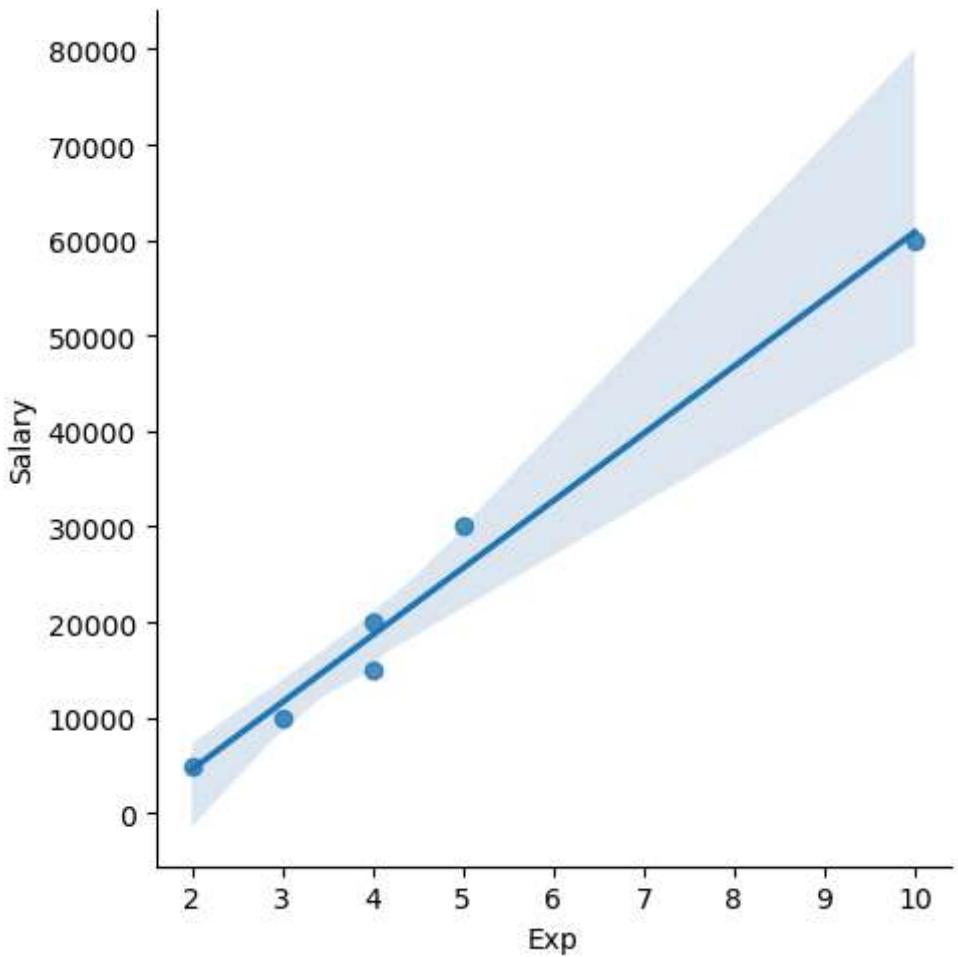
In [58]: `vis1=sns.distplot(clean_data['Salary'])`



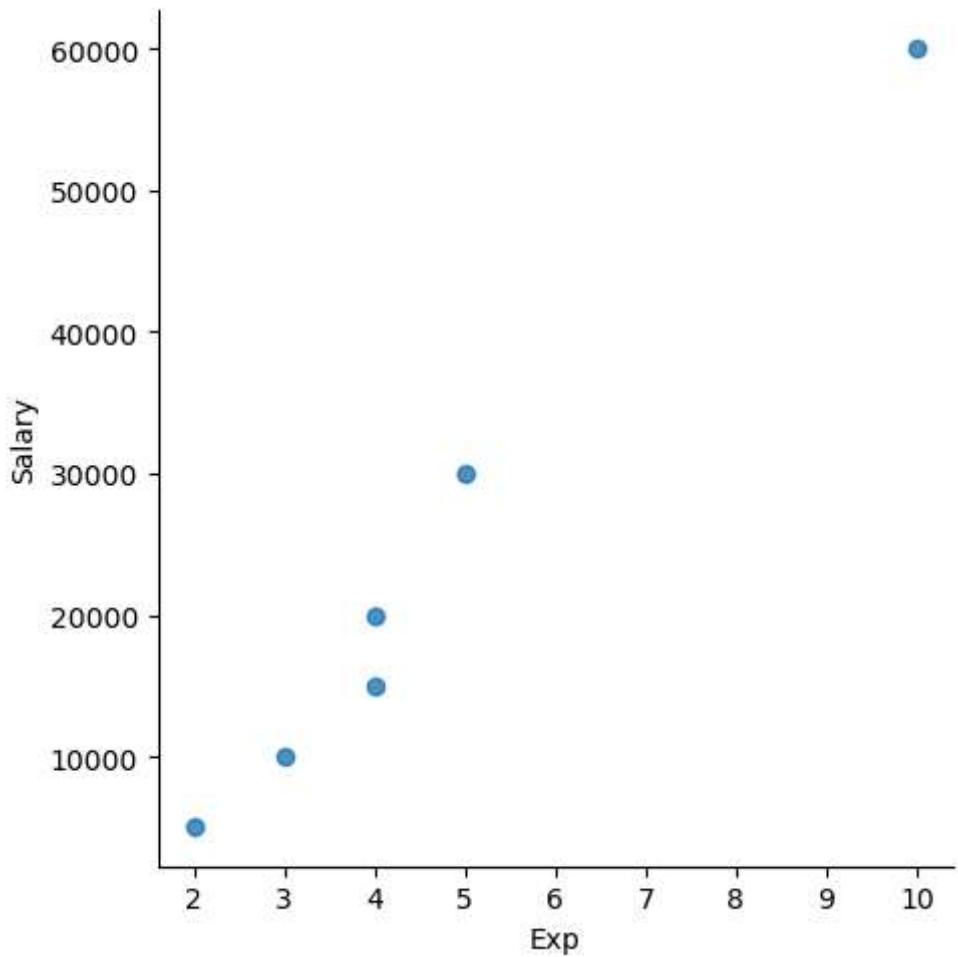
```
In [59]: vis2=plt.hist(clean_data['Salary'])
```



```
In [60]: vis3=sns.lmplot(data=clean_data,x='Exp',y='Salary')
```



```
In [62]: vis3=sns.lmplot(data=clean_data,x='Exp',y='Salary',fit_reg=False)
```



## splitting independent and dependent variable

```
In [63]: x_iv=clean_data[['Name','Domain','Age','Location','Exp']]  
y_iv
```

```
Out[63]:
```

	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 [70]: y_iv=clean_data[['Salary']]  
y_iv
```

Out[70]:

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

In [67]:

imputation=pd.get\_dummies(clean\_data,dtype=int)

In [68]:

imputation

Out[68]:

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



In [69]:

clean\_data

Out[69]:

	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 [71]:

len(clean\_data)

Out[71]:

6

In [76]:

imputation.columns

```
Out[76]: Index(['Age', 'Salary', 'Exp', 'Name_Jane', 'Name_Kim', 'Name_Mike',  
                 'Name_Teddy', 'Name_Umar', 'Name_Uttam', 'Domain_Analytics',  
                 'Domain_Dataanalyst', 'Domain_Datascience', 'Domain_NLP',  
                 'Domain_Statistics', 'Domain_Testing', 'Location_Bangalore',  
                 'Location_Delhi', 'Location_Hyderabad', 'Location_Mumbai'],  
                dtype='object')
```

```
In [77]: len(imputation.columns)
```

```
Out[77]: 19
```

```
In [ ]:
```