

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

```
In [2]: emp = pd.read_excel(r"D:\sir gen Ai\EDA\Rawdata.xlsx")
```

```
In [3]: emp
```

```
Out[3]:
```

	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]: 2112457861472
```

```
In [5]: emp.isnull()
```

```
Out[5]:
```

	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 [6]: emp.columns
```

```
Out[6]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')
```

```
In [7]: emp.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 [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.isnull().sum()`

Out[9]:

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

In [10]: `emp.Name[0]`

Out[10]: 'Mike'

In [11]: `emp`

Out[11]:

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

In [13]: `emp['Name']`

Out[13]:

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

In [14]: `emp`

Out[14]:

	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 [15]: `emp['Domain'] = emp['Domain'].str.replace(r'\W', '', regex=True)`In [16]: `emp['Domain']`

Out[16]:

0	Datascience
1	Testing
2	Dataanalyst
3	Analytics
4	Statistics
5	NLP

Name: Domain, dtype: object

In [17]: `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	Analytics	NaN	Hyderbad	2000^0	NaN
4	Uttam	Statistics	67-yr	NaN	30000-	5+ year
5	Kim	NLP	55yr	Delhi	6000^\$0	10+

In [18]: `emp['Age'] = emp['Age'].str.replace(r'\W', '', regex=True)`In [19]: `emp`

	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 [20]: emp['Age']= emp['Age'].str.extract('(\d+)')
```

```
In [21]: emp
```

Out[21]:

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

Out[22]:

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

```
In [24]: 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.replace(r'\W',' ', regex=True)
```

```
In [26]: emp['Exp']
```

Out[26]:

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

```
In [27]: emp['Exp'] = emp['Exp'].str.extract('(\d+)')
emp['Exp']
```

```
Out[27]: 0      2
          1      3
          2      4
          3    NaN
          4      5
          5     10
Name: Exp, dtype: object
```

```
In [28]: clean_data = emp.copy()
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	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 [29]: clean_data.isnull().sum()
```

```
Out[29]: Name      0
          Domain    0
          Age       2
          Location  2
          Salary    0
          Exp       1
          dtype: int64
```

```
In [30]: import numpy as np
```

```
In [31]: clean_data['Age'] = clean_data['Age'].fillna(np.mean(pd.to_numeric(clean_data['A
```

```
In [32]: clean_data['Age']
```

```
Out[32]: 0      34
          1      45
          2    50.25
          3    50.25
          4      67
          5      55
Name: Age, dtype: object
```

```
In [33]: clean_data['Exp']=clean_data['Exp'].fillna(np.mean(pd.to_numeric(clean_data['Exp
```

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

```
Out[34]: 0      2
         1      3
         2      4
         3    4.8
         4      5
         5     10
Name: Exp, dtype: object
```

```
In [35]: clean_data['Location'] = clean_data['Location'].fillna(clean_data['Location'].mode().iloc[0])
```

```
In [36]: emp
```

```
Out[36]:   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 [37]: clean_data
```

```
Out[37]:   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 [38]: 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 [39]: clean_data['Age'] = clean_data['Age'].astype(int)
```

```
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      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 [41]: clean_data['Salary']= clean_data['Salary'].astype(int)
```

```
In [42]: clean_data['Exp']= clean_data['Exp'].astype(int)
```

```
In [43]: 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 [44]: clean_data['Name']= clean_data['Name'].astype('category')
```

```
In [45]: clean_data['Domain']= clean_data['Domain'].astype('category')
```

```
In [46]: clean_data['Location']= clean_data['Location'].astype('category')
```

```
In [47]: 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 [48]: clean_data.to_csv('clean_data.csv')
```

```
In [49]: import os
```

```
In [50]: os.getcwd()
```

```
Out[50]: 'C:\\Users\\komme'
```

```
In [ ]: # EDA Techniques
```

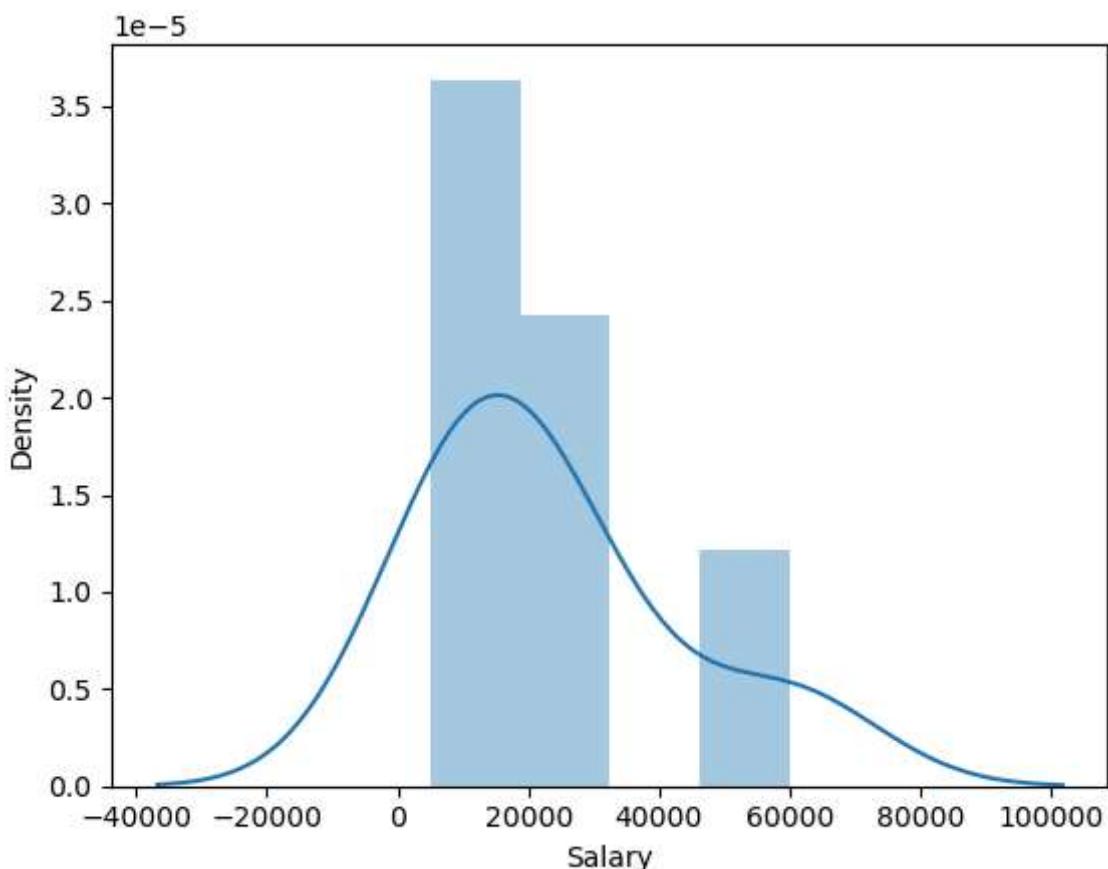
```
In [51]: import matplotlib.pyplot as plt
import seaborn as sns
```

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

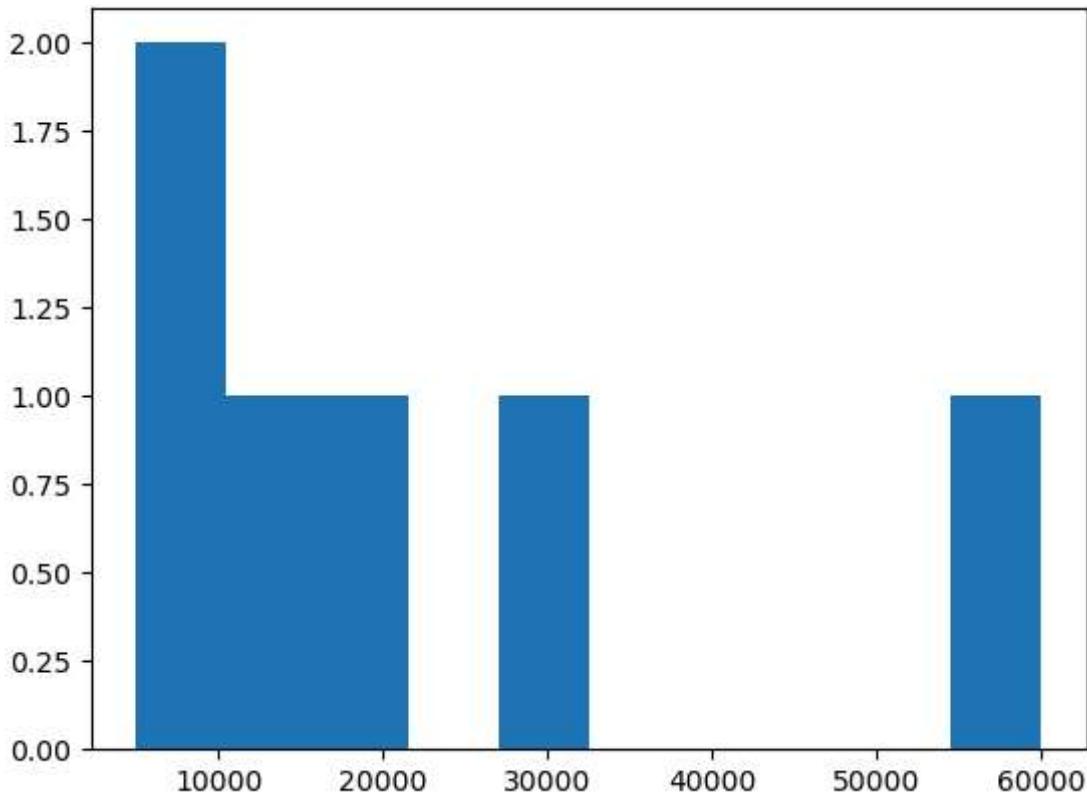
```
In [54]: clean_data['Salary']
```

```
Out[54]: 0      5000
1     10000
2     15000
3     20000
4     30000
5     60000
Name: Salary, dtype: int32
```

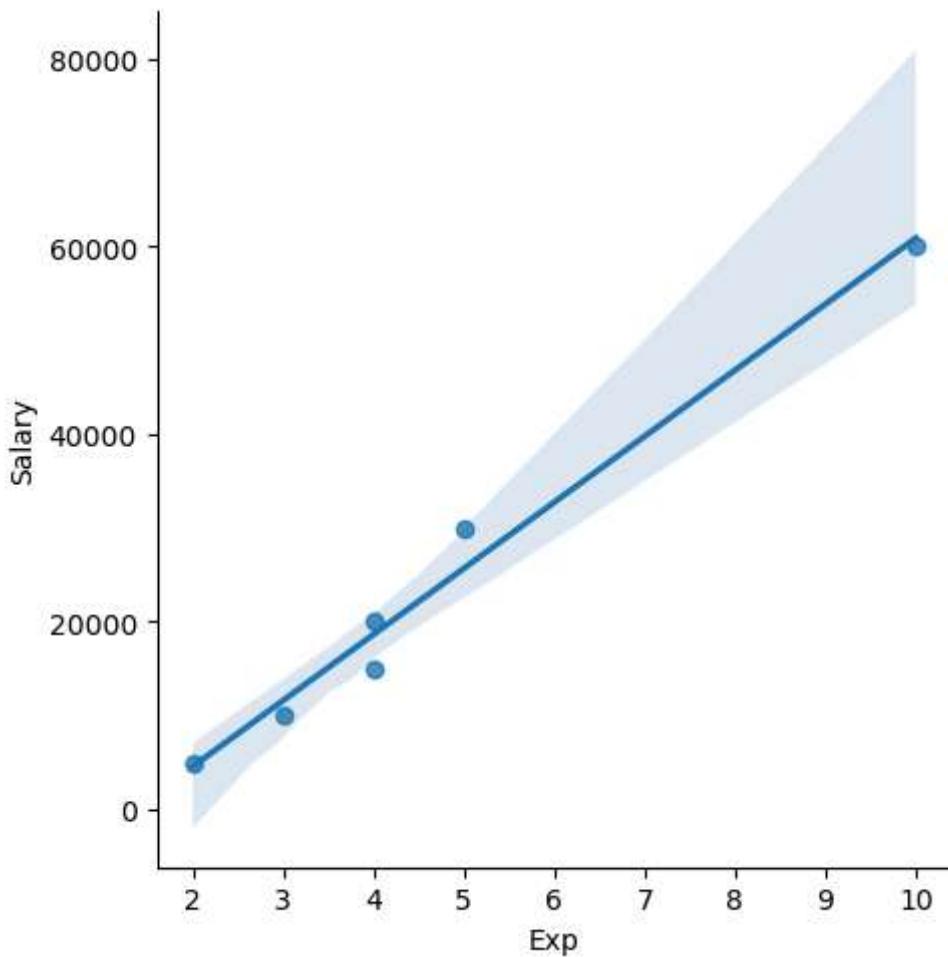
```
In [57]: vis1 = sns.distplot(clean_data['Salary'])    # distplot func using 1 parameter
```



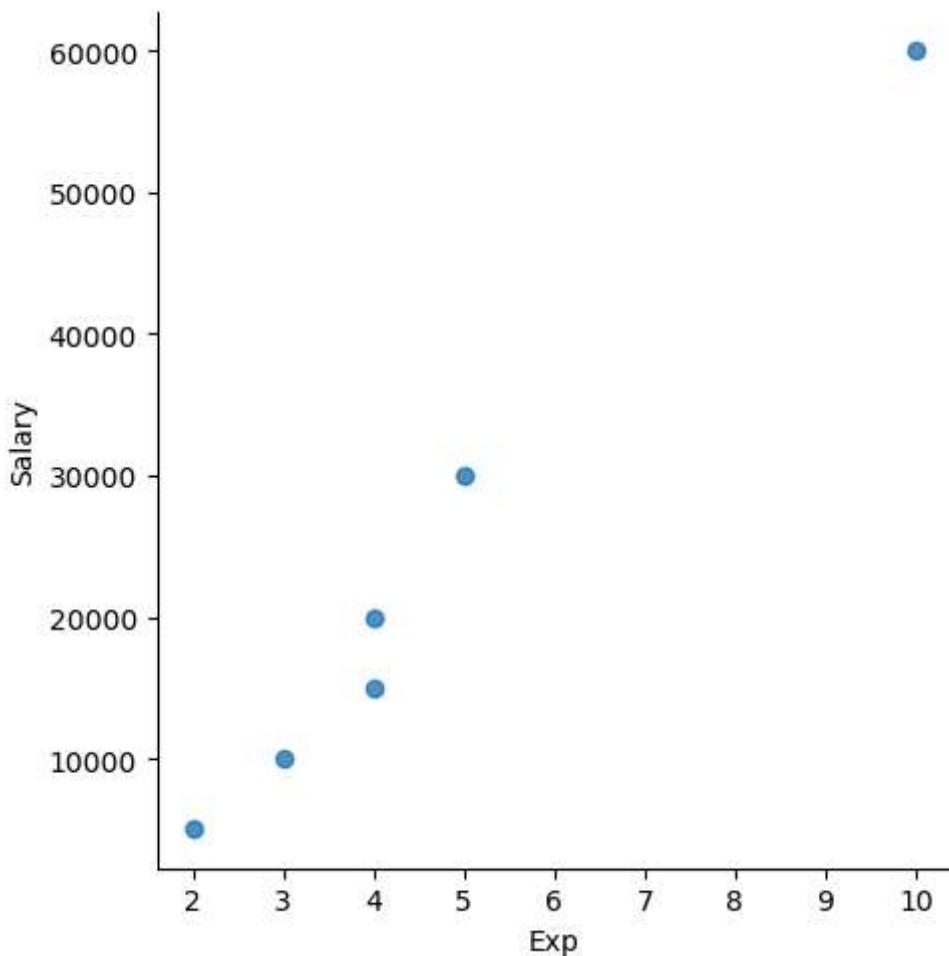
```
In [58]: v2 = plt.hist(clean_data['Salary'])    # histogram
```



```
In [59]: v3 = sns.lmplot(data=clean_data,x = 'Exp', y = 'Salary')      # Lmpplot is used pl
```



```
In [60]: v4 = sns.lmplot(data = clean_data,x = 'Exp',y = 'Salary',fit_reg=False)    # fig
```



```
In [61]: clean_data
```

```
Out[61]:
```

	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

SLICING

```
In [62]: clean_data[:]
```

Out[62]:

	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 [63]:

clean_data[::-1] # Reverse print

Out[63]:

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

In [64]:

clean_data[1:4]

Out[64]:

	Name	Domain	Age	Location	Salary	Exp
1	Teddy	Testing	45	Bangalore	10000	3
2	Umar	Dataanalyst	50	Bangalore	15000	4
3	Jane	Analytics	50	Hyderbad	20000	4

In [66]:

clean_data[0:4:2]

Out[66]:

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascience	34	Mumbai	5000	2
2	Umar	Dataanalyst	50	Bangalore	15000	4

SEPERATE CATEGIRICAL DATA AND NUMERICAL DATA OR SEPERATE DEPENDENT VARIABLE(Y) AND INDIPENDENT VARIABLE(X) COLUMNS

In [70]:

clean_data.columns

```
Out[70]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')
```

```
In [75]: xiv= clean_data[['Name','Age','Location','Exp']]  
xiv
```

```
Out[75]:   Name  Age  Location  Exp
```

0	Mike	34	Mumbai	2
1	Teddy	45	Bangalore	3
2	Umar	50	Bangalore	4
3	Jane	50	Hyderbad	4
4	Uttam	67	Bangalore	5
5	Kim	55	Delhi	10

```
In [78]: yiv = clean_data[['Salary']]  
yiv
```

```
Out[78]:  Salary
```

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

```
In [79]: clean_data
```

```
Out[79]:   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 [84]: imputation = pd.get_dummies(clean_data)
```

```
In [85]: imputation
```

Out[85]:

	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 []:

In []:

In []:

In []:

In []: