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

```python
In [2]: emp=pd.read_excel(r"C:\Users\MANISHA\Downloads\Rawdata.xlsx")
```

```python
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+ |

```python
In [4]: emp.shape
```

Out[4]: (6, 6)

```python
In [5]: len(emp)
```

Out[5]: 6

```python
In [6]: emp.columns
```

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

```python
In [7]: len(emp.columns)
```

Out[7]: 6

```python
In [8]: 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
```

```python
In [9]: emp['Name']
```

```
Out[9]:  0      Mike
         1    Teddy^
         2    Uma#r
         3      Jane
         4   Uttam*
         5       Kim
         Name: Name, dtype: object
```

```
In [10]:  emp['Domain']
```

```
Out[10]:  0    Datascience#$
          1          Testing
          2   Dataanalyst^^#
          3       Ana^^lytics
          4        Statistics
          5               NLP
          Name: Domain, dtype: object
```

```
In [11]:  emp['Salary']
```

```
Out[11]:  0     5^00#0
          1    10%%000
          2    1$5%000
          3     2000^0
          4     30000-
          5    6000^$0
          Name: Salary, dtype: object
```

```
In [12]:  emp[['Name','Domain','Age','Location','Salary','Exp']]
```

Out[12]:

|   | 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+ |

# Cleaning data

```
In [14]:  emp['Name']
```

```
Out[14]: 0      Mike
         1    Teddy^
         2     Uma#r
         3      Jane
         4    Uttam*
         5       Kim
         Name: Name, dtype: object
```

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

```
In [16]: emp['Name']
```

```
Out[16]: 0      Mike
         1     Teddy
         2      Umar
         3      Jane
         4     Uttam
         5       Kim
         Name: Name, dtype: object
```

```
In [17]: emp
```

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** | 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 [18]: emp['Domain']=emp['Domain'].str.replace(r'\W','',regex=True)
```

```
In [19]: emp['Domain']
```

```
Out[19]: 0    Datascience
         1        Testing
         2    Dataanalyst
         3      Analytics
         4     Statistics
         5            NLP
         Name: Domain, dtype: object
```

```
In [20]: emp
```

Out[20]:

| | 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 [21]:
```python
emp['Age']=emp['Age'].str.replace(r'\W','',regex=True)
```

In [22]:
```python
emp['Age']
```

Out[22]:
```
0    34years
1       45yr
2        NaN
3        NaN
4       67yr
5       55yr
Name: Age, dtype: object
```

In [23]:
```python
emp['Age']=emp['Age'].str.extract(r'(\d+)')
```

In [24]:
```python
emp['Age']
```

Out[24]:
```
0     34
1     45
2    NaN
3    NaN
4     67
5     55
Name: Age, dtype: object
```

In [25]:
```python
emp['Location']=emp['Location'].str.replace(r'\W','',regex=True)
```

In [26]:
```python
emp['Location']
```

Out[26]:
```
0       Mumbai
1    Bangalore
2          NaN
3     Hyderbad
4          NaN
5        Delhi
Name: Location, dtype: object
```

In [27]:
```python
emp['Salary']=emp['Salary'].str.replace(r'\W','',regex=True)
emp['Salary']
```

```
Out[27]: 0     5000
         1    10000
         2    15000
         3    20000
         4    30000
         5    60000
         Name: Salary, dtype: object
```

In [28]: `emp.head()`

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> yrs |
| **3** | Jane | Analytics | NaN | Hyderbad | 20000 | NaN |
| **4** | Uttam | Statistics | 67 | NaN | 30000 | 5+ year |

In [29]: `emp['Exp']=emp['Exp'].str.extract(r'(\d+)')`

In [30]: `emp['Exp']`

```
Out[30]: 0      2
         1      3
         2      4
         3    NaN
         4      5
         5     10
         Name: Exp, dtype: object
```

In [31]: `emp`

Out[31]:

| | 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 [32]: 
```
clean_data=emp.copy()
clean_data
```

Out[32]:

| | 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 [66]:
```python
import numpy as np
```

In [68]:
```python
clean_data['Age']= clean_data['Age'].fillna(np.mean(pd.to_numeric(clean_data['Age']
clean_data
```

Out[68]:

| | 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 | NaN |
| 4 | Uttam | Statistics | 67 | NaN | 30000 | 5 |
| 5 | Kim | NLP | 55 | Delhi | 60000 | 10 |

In [70]:
```python
clean_data['Exp']= clean_data['Exp'].fillna(np.mean(pd.to_numeric(clean_data['Exp']
clean_data
```

Out[70]:

| | 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 [72]:
```python
clean_data['Location']= clean_data['Location'].fillna(clean_data['Location'].mode()
clean_data
```

Out[72]:

| | 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 [74]:
```python
clean_data.isna().sum()
```

Out[74]:
```
Name        0
Domain      0
Age         0
Location    0
Salary      0
Exp         0
dtype: int64
```

In [76]:
```python
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 [78]:
```python
clean_data['Name']=clean_data['Name'].astype(str)
clean_data['Domain']=clean_data['Domain'].astype(str)
clean_data['Location']=clean_data['Location'].astype(str)
```

In [80]:
```python
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 [82]:
```python
clean_data['Age']=clean_data['Age'].astype(int)
clean_data['Salary']=clean_data['Salary'].astype(int)
clean_data['Exp']=clean_data['Exp'].astype(int)
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 [84]: `clean_data`

Out[84]:

|   | 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 [90]:
```python
clean_data.to_csv('clean_data.csv')
```

In [88]:
```python
clean_data.to_excel(r"C:\Users\MANISHA\Downloads\20th- EDA Practicle\Rawdata.xlsx")
```

In [92]:
```python
import os
os.getcwd()
```

Out[92]: 'C:\\Users\\MANISHA'

In [94]: `clean_data`

Out[94]:

| | 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 |

# EDA TECHNIQUE LETS APPLY

In [97]:
```python
import matplotlib.pyplot as plt
import seaborn as sns
```
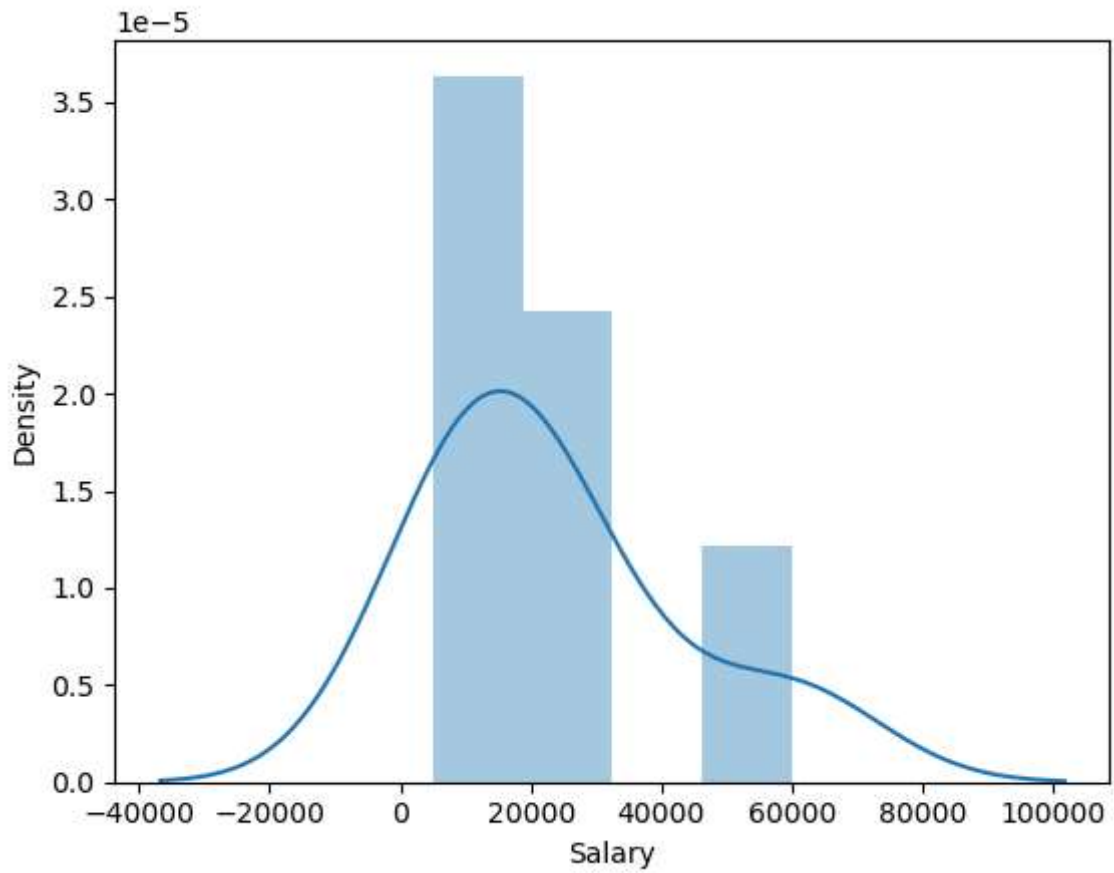
In [99]:
```python
import warnings
warnings.filterwarnings('ignore')
```
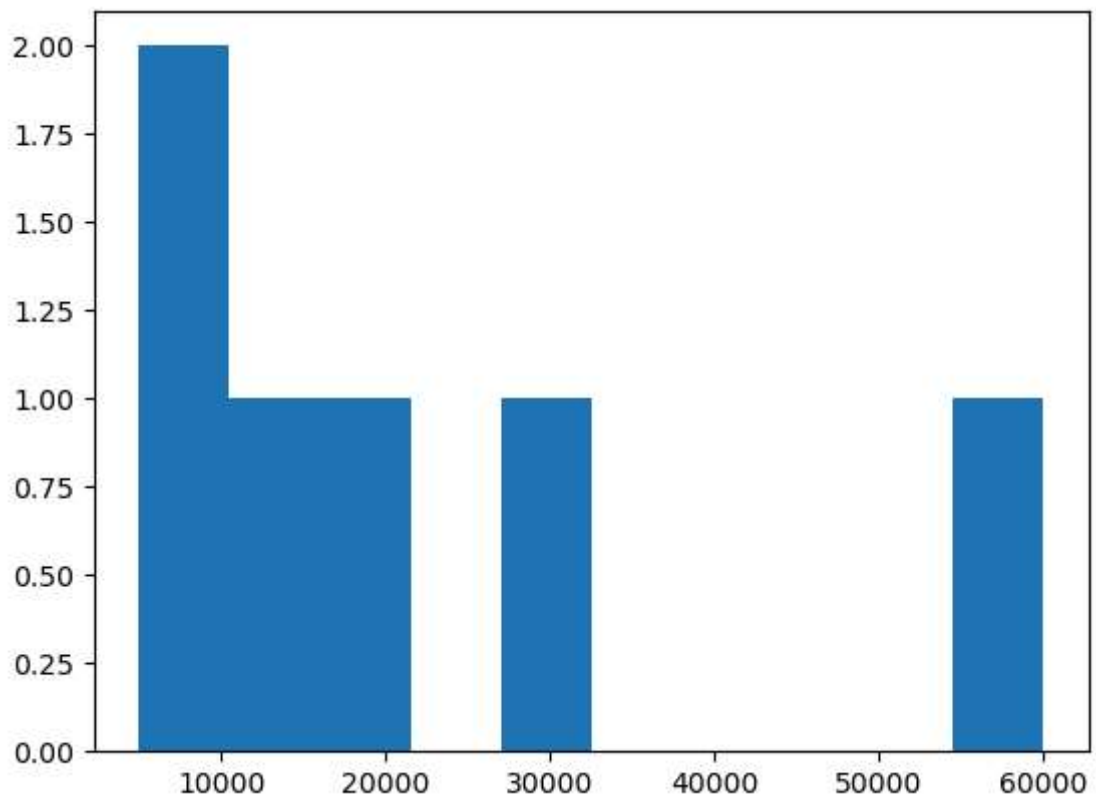
In [101...
```python
clean_data['Salary']
```

Out[101...
```
0     5000
1    10000
2    15000
3    20000
4    30000
5    60000
Name: Salary, dtype: int32
```
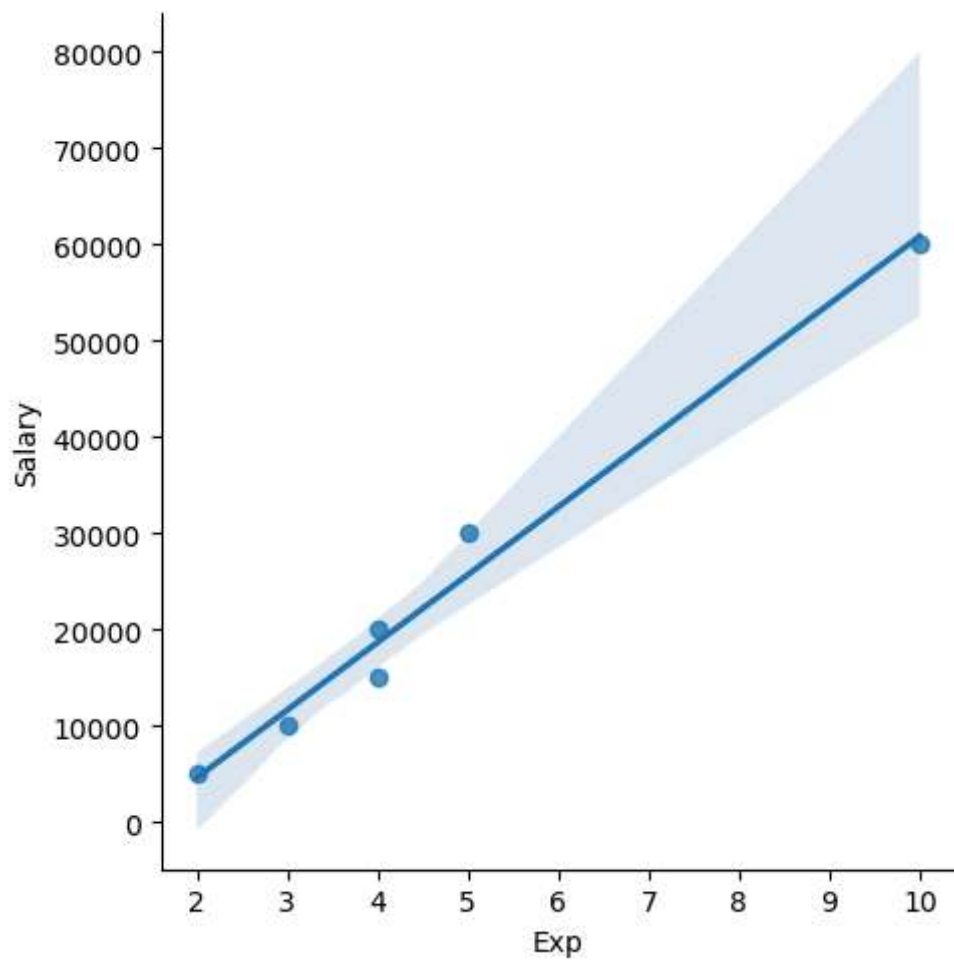
In [103...
```python
vis1=sns.distplot(clean_data['Salary'])
```
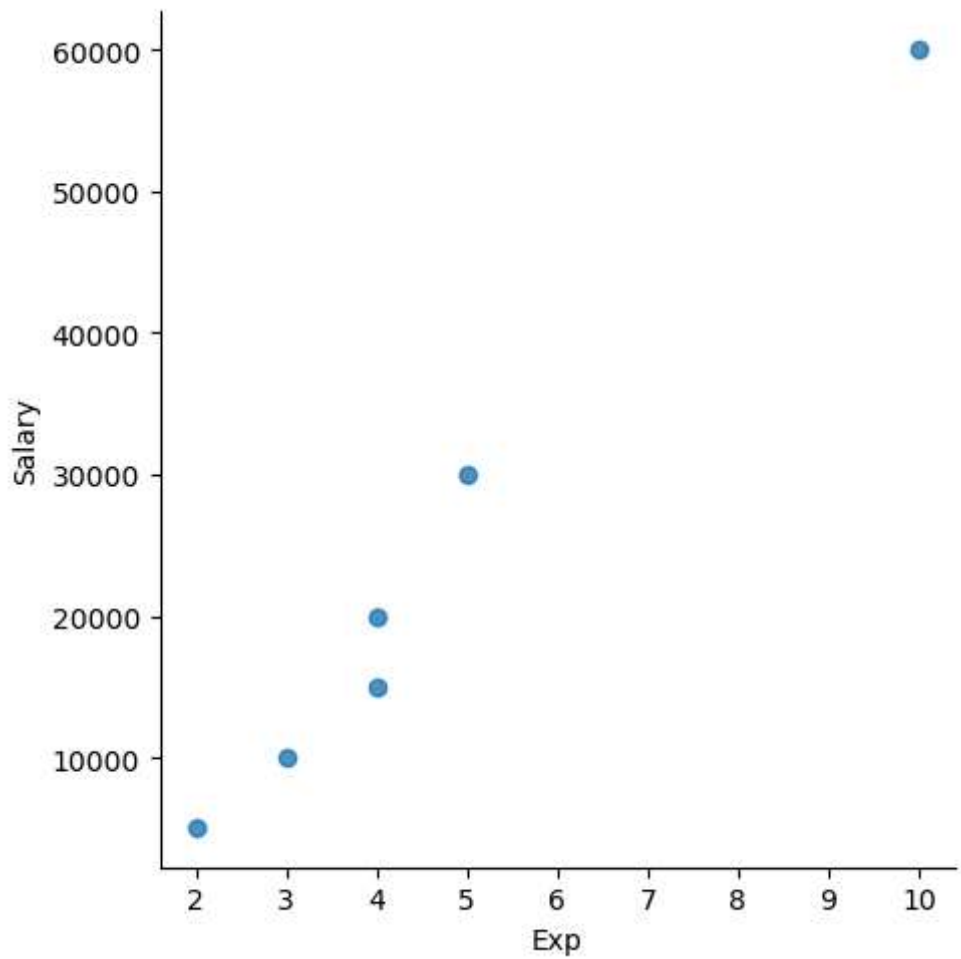
```python
vis2=plt.hist(clean_data['Salary'])
```

```python
vis3=sns.lmplot(data=clean_data,x='Exp',y='Salary')
```

```
In [109…   vis4=sns.lmplot(data=clean_data,x='Exp',y='Salary',fit_reg=False)
```

`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  | Bangalore | 15000  | 4   |
| 3 | Jane  | Analytics   | 50  | Hyderbad  | 20000  | 4   |
| 4 | Uttam | Statistics  | 67  | Bangalore | 30000  | 5   |
| 5 | Kim   | NLP         | 55  | Delhi     | 60000  | 10  |

`clean_data[0:6:2]`

|   | 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 [117…   clean_data[::-1]
```

Out[117…

|   | 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 [119…   clean_data.columns
```

Out[119…   Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')

```
In [121…   X_iv=clean_data[['Name','Domain','Age','Location','Exp']]
```

```
In [123…   X_iv
```

Out[123…

|   | 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 | Hyderbad | 4 |
| **4** | Uttam | Statistics | 67 | Bangalore | 5 |
| **5** | Kim | NLP | 55 | Delhi | 10 |

```
In [125…   Y_dv=clean_data[['Salary']]
```

```
In [127…   Y_dv
```

Out[127…

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

```
emp
```

|   | 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 |

```
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 | Bangalore | 15000 | 4 |
| 3 | Jane | Analytics | 50 | Hyderbad | 20000 | 4 |
| 4 | Uttam | Statistics | 67 | Bangalore | 30000 | 5 |
| 5 | Kim | NLP | 55 | Delhi | 60000 | 10 |

```
X_iv
```

|   | 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 | Hyderbad | 4 |
| 4 | Uttam | Statistics | 67 | Bangalore | 5 |
| 5 | Kim | NLP | 55 | Delhi | 10 |

```
Y_dv
```

Out[135...

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

In [137...

```
clean_data
```

Out[137...

| | 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 [145...

```
imputation=pd.get_dummies(clean_data,dtype=int)
```

In [147...

```
imputation
```

Out[147...

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

In [ ]: