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

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

Out[2]:  '2.1.4'

```
In [3]:  emp=pd.read_excel(r"D:\DS_NIT\Oct_24\EDA.xlsx")
```

```
In [4]:  emp
```

Out[4]:

|   | 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 [5]:  emp.shape
```

Out[5]:  (6, 6)

```
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.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['Domain']`

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

In [11]: `emp.isnull()`

Out[11]:

| | 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 [12]: `emp.isnull().sum()`

```
Out[12]:  Name         0
          Domain       0
          Age          2
          Location     2
          Salary       0
          Exp          1
          dtype: int64
```

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','')`

In [15]: `emp['Name']`

```
Out[15]:  0       Mike
          1      Teddy^
          2      Uma#r
          3       Jane
          4     Uttam*
          5        Kim
          Name: Name, dtype: object
```

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

In [17]: `emp['Name']`

```
Out[17]:  0       Mike
          1      Teddy
          2       Umar
          3       Jane
          4      Uttam
          5        Kim
          Name: Name, dtype: object
```

In [18]: `emp.columns`

Out[18]: `Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')`

In [19]: `emp.Domain`

```
Out[19]:  0     Datascience#$
          1          Testing
          2    Dataanalyst^^#
          3       Ana^^lytics
          4       Statistics
          5              NLP
          Name: Domain, dtype: object
```

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

```
In [21]: emp['Domain']
```

```
Out[21]: 0    Datascience
         1        Testing
         2     Dataanalyst
         3      Analytics
         4     Statistics
         5            NLP
         Name: Domain, dtype: object
```

```
In [22]: emp.columns
```

```
Out[22]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')
```

```
In [23]: emp.Age
```

```
Out[23]: 0     34 years
         1       45' yr
         2          NaN
         3          NaN
         4        67-yr
         5         55yr
         Name: Age, dtype: object
```

```
In [24]: emp['Age']=emp['Age'].str.extract(r'(\d+)')
```

```
In [25]: emp['Age']
```

```
Out[25]: 0      34
         1      45
         2     NaN
         3     NaN
         4      67
         5      55
         Name: Age, dtype: object
```

```
In [26]: emp.columns
```

```
Out[26]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')
```

```
In [27]: emp.Location
```

```
Out[27]: 0        Mumbai
         1     Bangalore
         2           NaN
         3      Hyderbad
         4           NaN
         5         Delhi
         Name: Location, dtype: object
```

```
In [28]: emp.Salary
```

```
Out[28]:  0      5^00#0
          1     10%%000
          2     1$5%000
          3      2000^0
          4      30000-
          5     6000^$0
          Name: Salary, dtype: object
```

```
In [29]:  emp['Salary']=emp['Salary'].str.replace(r'\W','',regex=True)
```

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

```
Out[30]:  0      5000
          1     10000
          2     15000
          3     20000
          4     30000
          5     60000
          Name: Salary, dtype: object
```

```
In [31]:  emp.Exp
```

```
Out[31]:  0          2+
          1          <3
          2      4> yrs
          3         NaN
          4     5+ year
          5         10+
          Name: Exp, dtype: object
```

```
In [32]:  emp['Exp']=emp['Exp'].str.extract(r'(\d+)')
```

```
In [33]:  emp['Exp']
```

```
Out[33]:  0       2
          1       3
          2       4
          3     NaN
          4       5
          5      10
          Name: Exp, dtype: object
```

```
In [34]:  emp.head()
```

Out[34]:

| | 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 |

```
In [35]:  clean_data=emp.copy()
```

```
In [36]:  clean_data
```

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]:  import numpy as np
```

```
In [38]:  clean_data['Age']=clean_data['Age'].fillna(np.mean(pd.to_numeric(clean_data['Age'])
```

```
In [39]:  clean_data['Age']
```

```
Out[39]:  0        34
          1        45
          2     50.25
          3     50.25
          4        67
          5        55
          Name: Age, dtype: object
```

```
In [40]:  clean_data['Location']
```

```
Out[40]:  0        Mumbai
          1     Bangalore
          2           NaN
          3      Hyderbad
          4           NaN
          5         Delhi
          Name: Location, dtype: object
```

```
In [41]:  clean_data['Location']=clean_data['Location'].fillna(clean_data['Location'].mode()[
```

```
In [42]:  clean_data['Location']
```

```
Out[42]:  0        Mumbai
          1     Bangalore
          2     Bangalore
          3      Hyderbad
          4     Bangalore
          5         Delhi
          Name: Location, dtype: object
```

```
In [43]: clean_data
```

Out[43]:

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

```
In [44]: clean_data['Exp']=clean_data['Exp'].fillna(np.mean(pd.to_numeric(clean_data['Exp'])
```

```
In [45]: clean_data['Exp']
```

```
Out[45]: 0      2
         1      3
         2      4
         3    4.8
         4      5
         5     10
         Name: Exp, dtype: object
```

```
In [46]: clean_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 | 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 [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      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 [48]:
```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)
```

In [49]:
```python
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 [50]:
```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      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 [51]:
```python
clean_data
```

Out[51]:

|   | 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 [52]:
```python
clean_data.to_csv('clean_datapr.csv')
```

```
In [53]: import os
         os.getcwd()
```

Out[53]: "C:\\Users\\evenk\\OneDrive\\Desktop\\DS_NIT\\Oct'24"
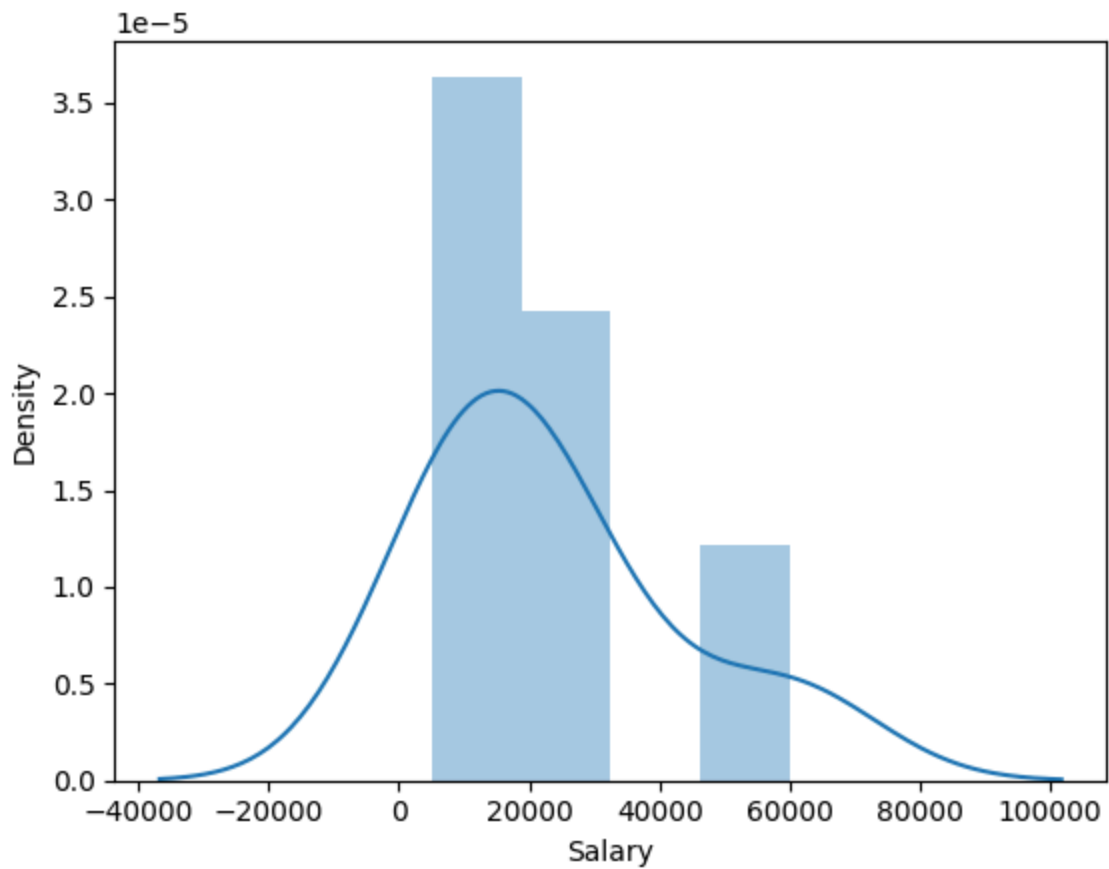
```
In [54]: clean_data
```

Out[54]:

|   | 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  |

## LETS APPLY EDA TECHNIQUES

```
In [55]: import matplotlib.pyplot as plt
         import seaborn as sns
```

```
In [56]: import warnings
         warnings.filterwarnings('ignore')
```

```
In [57]: clean_data
```

Out[57]:

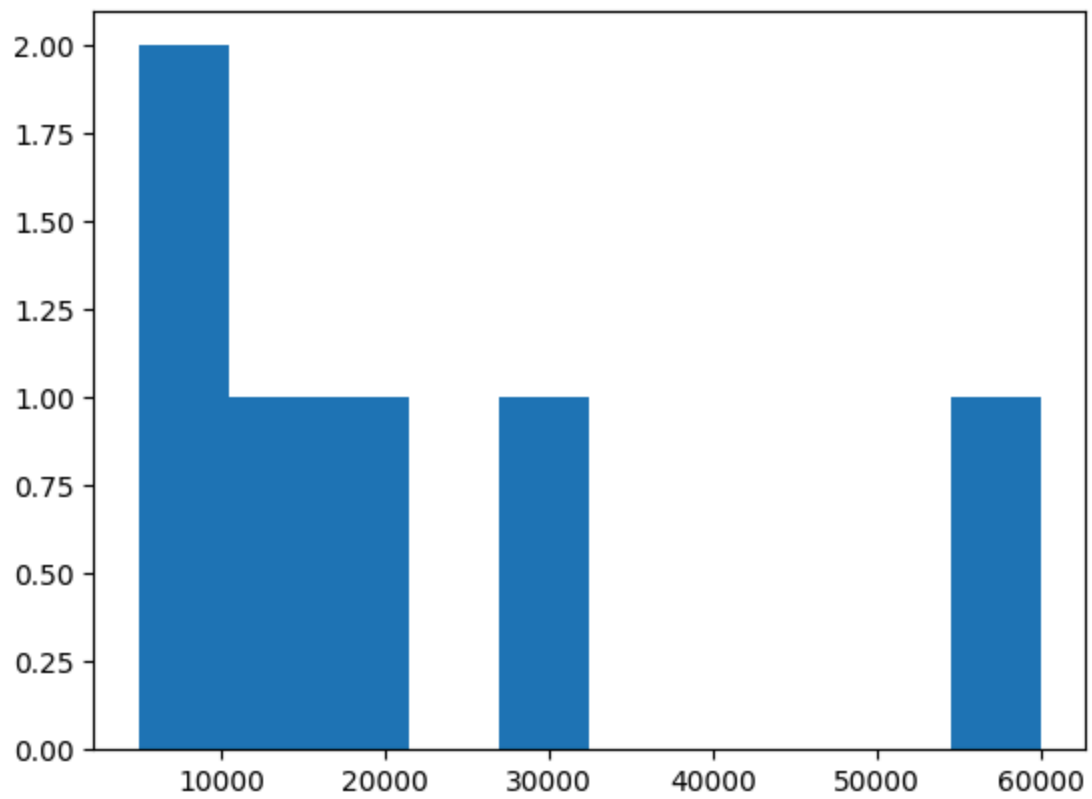|   | Name  | Domain      | Age | Location  | Salary | Exp |
|---|-------|-------------|-----|-----------|--------|-----|
| 0 | Mike  | Datascience | 34  | Mumbai    | 5000   | 2   |
| 1 | Teddy | Testing     | 45  | Bangalore | 10000  | 3   |
| 2 | Umar  | Dataanalyst | 50  | Bangalore | 15000  | 4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

```
In [58]: clean_data['Salary']
```

0     5000
        1    10000
        2    15000
        3    20000
        4    30000
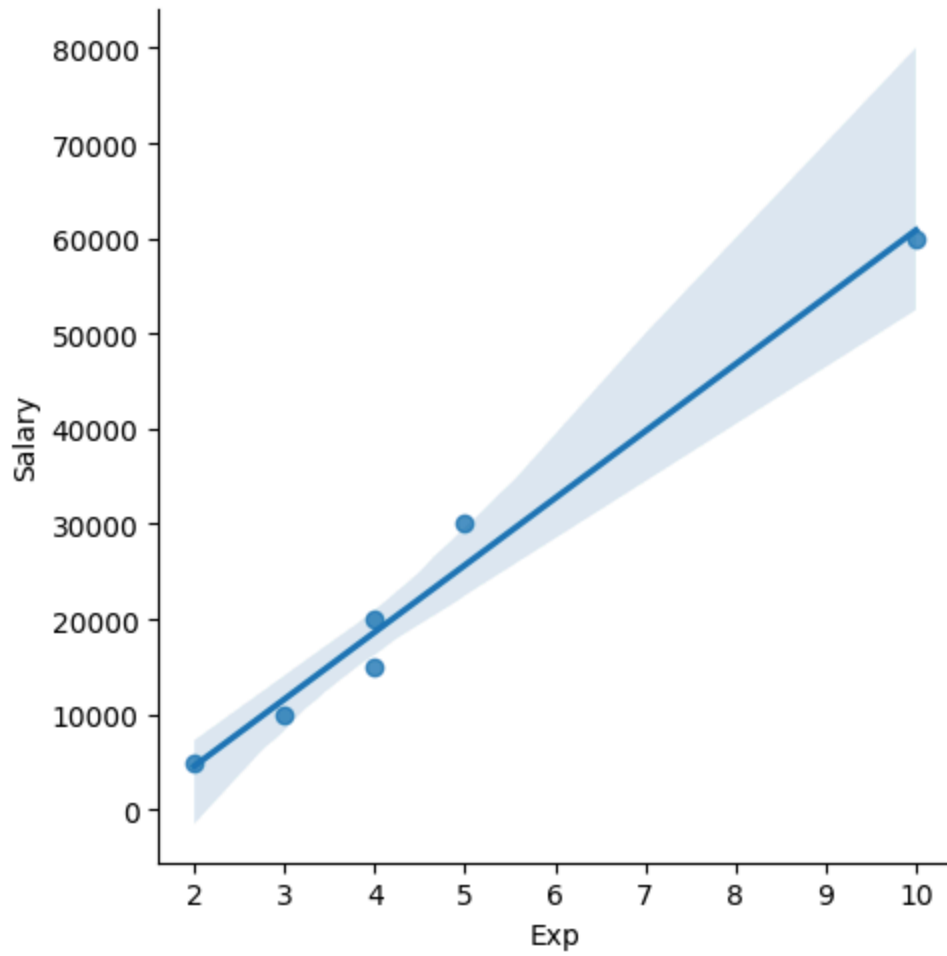        5    60000
        Name: Salary, dtype: int32

# Visualization

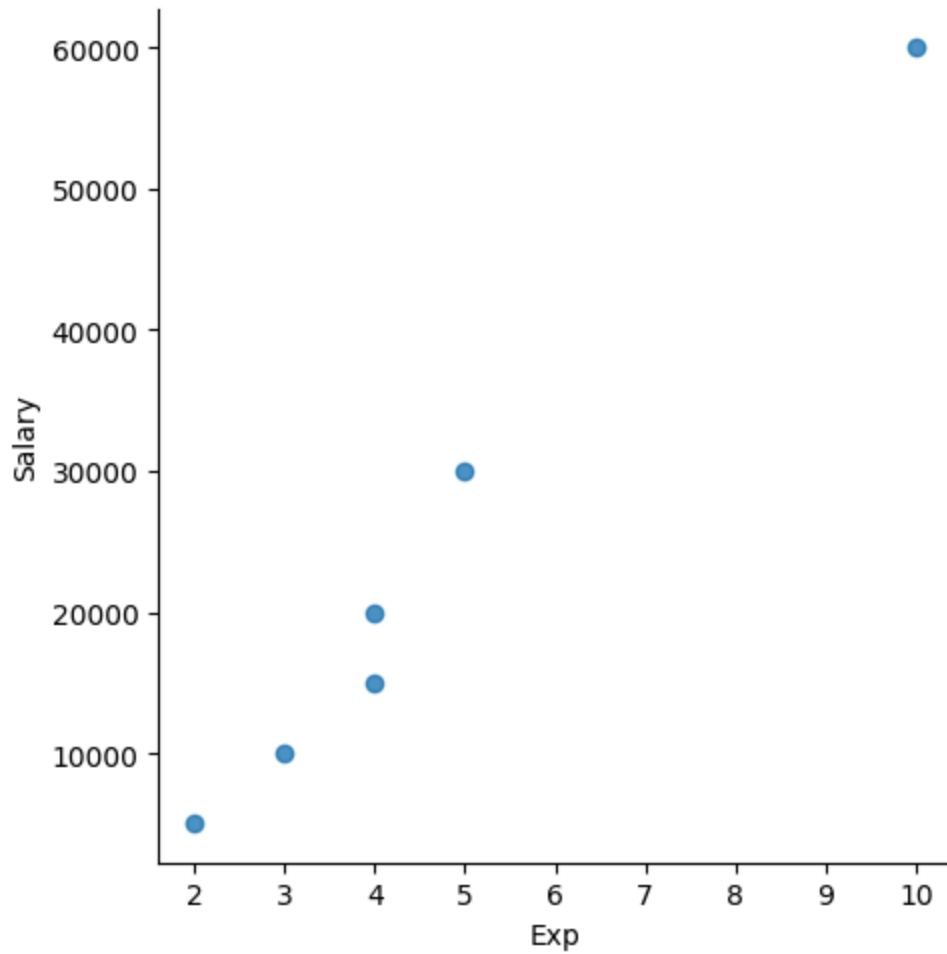In [59]: ```python
#univariate
vis1=sns.distplot(clean_data['Salary'])
```



In [60]: ```python
#outlier
vis2=plt.hist(clean_data['Salary'])
```

In [61]: 
```python
#bivariate
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)
```

`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[:]`

Out[64]:

| | 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 [65]: `clean_data[1:6:2]`

Out[65]:

| | Name | Domain | Age | Location | Salary | Exp |
|---|---|---|---|---|---|---|
| **1** | Teddy | Testing | 45 | Bangalore | 10000 | 3 |
| **3** | Jane | Analytics | 50 | Hyderbad | 20000 | 4 |
| **5** | Kim | NLP | 55 | Delhi | 60000 | 10 |

In [66]: `clean_data[2:6]`

Out[66]:

| | Name | Domain | Age | Location | Salary | Exp |
|---|---|---|---|---|---|---|
| **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 [67]: `clean_data[:]`

Out[67]:

| | 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 [68]: `clean_data[::-1]`

Out[68]:

| | 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 [69]: `clean_data.columns`

Out[69]: `Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')`

In [70]:
```
#variable identification(independent variable)
x_iv=clean_data[['Name', 'Domain', 'Age', 'Location', 'Exp']]
```

In [71]: `x_iv`

Out[71]:

| | 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 [72]:
```
#Variable identification(Dependent variable)
y_iv=clean_data[['Salary']]
```

In [73]: `y_iv`

Out[73]:

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

In [74]: `clean_data`

Out[74]:

| | 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 [75]: `x_iv`

Out[75]:

| | 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 [76]: `y_iv`

Out[76]:

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

In [77]: `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 |

In [78]:
```
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 |

In [83]:
```
imputation=pd.get_dummies(clean_data)
```

In [84]:
```
imputation
```

| | Age | Salary | Exp | Name_Jane | Name_Kim | Name_Mike | Name_Teddy | Name_Umar | Nan |
|---|---|---|---|---|---|---|---|---|---|
| **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 [85]:
```
imputation=pd.get_dummies(clean_data,dtype=int)
```

In [86]:
```
imputation
```

| | 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 | |