In [27]:
```python
import pandas as pd
```

In [29]:
```python
emp = pd.read_excel(r'C:\Users\samik\Downloads\Rawdata.xlsx')
```

In [31]:
```python
emp
```

Out[31]:

| | 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 [33]:
```python
emp.columns
```

Out[33]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')

In [35]:
```python
emp.shape
```

Out[35]: (6, 6)

In [37]:
```python
emp.head()
```

Out[37]:

| | 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 [39]:
```python
emp.tail()
```

Out[39]:

| | 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 [41]:
```python
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 [43]: `emp['Domain']`

Out[43]:
```
0       Datascience#$
1             Testing
2       Dataanalyst^^#
3          Ana^^lytics
4          Statistics
5                 NLP
Name: Domain, dtype: object
```

In [45]: `emp.isnull()`

Out[45]:

|   | 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 [47]: `emp.isnull().any()`

Out[47]:
```
Name        False
Domain      False
Age          True
Location     True
Salary      False
Exp          True
dtype: bool
```

In [51]: `emp.isnull().sum()`

Out[51]:
```
Name        0
Domain      0
Age         2
Location    2
Salary      0
Exp         1
dtype: int64
```

```
In [53]: emp['Name']
```

```
Out[53]: 0      Mike
         1    Teddy^
         2     Uma#r
         3      Jane
         4    Uttam*
         5       Kim
         Name: Name, dtype: object
```

# We will use regex to clean the data and removed all noise charactered from the dataset.

```
In [57]: emp['Name'] = emp['Name'].str.replace(r'\W','',regex = True)
```

```
In [59]: emp
```

Out[59]:

|   | 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 [61]: emp['Domain'] = emp['Domain'].str.replace(r'\W','',regex = True)
```

```
In [63]: emp
```

Out[63]:

|   | 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 [65]: emp['Age'] = emp['Age'].str.replace(r'\W','',regex = True)
```

```
In [69]: emp['Age']
```

```
Out[69]:  0      34years
          1        45yr
          2         NaN
          3         NaN
          4        67yr
          5        55yr
          Name: Age, dtype: object
```

```
In [71]:  emp['Age'] = emp['Age'].str.extract('(\d+)')
```

```
<>:1: SyntaxWarning: invalid escape sequence '\d'
<>:1: SyntaxWarning: invalid escape sequence '\d'
C:\Users\samik\AppData\Local\Temp\ipykernel_10224\1884116463.py:1: SyntaxWarning:
invalid escape sequence '\d'
  emp['Age'] = emp['Age'].str.extract('(\d+)')
```

```
In [73]:  emp
```

Out[73]:

| | 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 [75]:  emp['Salary'] = emp['Salary'].str.replace(r'\W','',regex = True)
```

```
In [77]:  emp
```

Out[77]:

| | 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 |
| 5 | Kim | NLP | 55 | Delhi | 60000 | 10+ |

```
In [79]:  emp['Exp'] = emp['Exp'].str.replace(r'\W','',regex = True)
```

```
In [81]:  emp
```

Out[81]:

| | 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 | 4yrs |
| 3 | Jane | Analytics | NaN | Hyderbad | 20000 | NaN |
| 4 | Uttam | Statistics | 67 | NaN | 30000 | 5year |
| 5 | Kim | NLP | 55 | Delhi | 60000 | 10 |

In [83]:
```python
import warnings
warnings.filterwarnings('ignore')
```

In [85]:
```python
emp['Exp'] = emp['Exp'].str.extract('(\d+)')
```

In [87]:
```python
emp
```

Out[87]:

| | 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 [89]:
```python
clean_data = emp.copy()
```

In [91]:
```python
clean_data
```

Out[91]:

| | 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 |

# Missing values treatment for numerical data

In [94]:
```python
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 | 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 [96]:
```python
clean_data['Age']
```

Out[96]:
```
0     34
1     45
2     NaN
3     NaN
4     67
5     55
Name: Age, dtype: object
```

In [98]:
```python
import numpy as np
```

In [102…
```python
clean_data['Age'] = clean_data['Age'].fillna(np.mean(pd.to_numeric(clean_data['A
```

In [104…
```python
clean_data['Age']
```

Out[104…
```
0       34
1       45
2     50.25
3     50.25
4       67
5       55
Name: Age, dtype: object
```

In [106…
```python
clean_data['Exp']
```

Out[106…
```
0      2
1      3
2      4
3     NaN
4      5
5     10
Name: Exp, dtype: object
```

In [108…
```python
clean_data['Exp'] = clean_data['Exp'].fillna(np.mean(pd.to_numeric(clean_data['E
```

In [110…
```python
clean_data['Exp']
```

```
Out[110…   0      2
           1      3
           2      4
           3    4.8
           4      5
           5     10
           Name: Exp, dtype: object
```

In [112…
```
clean_data
```

Out[112…

|   | 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 |

# We will do missing value traetment for categorical data with mode .

In [119…
```
clean_data
```

Out[119…

|   | 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 [116…
```
clean_data['Location']
```

Out[116…
```
0        Mumbai
1     Bangalore
2           NaN
3      Hyderbad
4           NaN
5         Delhi
Name: Location, dtype: object
```

In [121…
```
clean_data['Location'] = clean_data['Location'].fillna(clean_data['Location'].mo
```

In [123…
```
clean_data
```

Out[123…

| | 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 [125…

```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
```

# We will change numerical data type to int and categorical to category.

In [128…

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

In [130…

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

In [132…

```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      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 [134…

```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 [136…  `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 [138…  `clean_data`

Out[138…

|   | 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 [140…  `clean_data.to_csv('clean_data.csv')`

In [142…
```python
import os
os.getcwd()
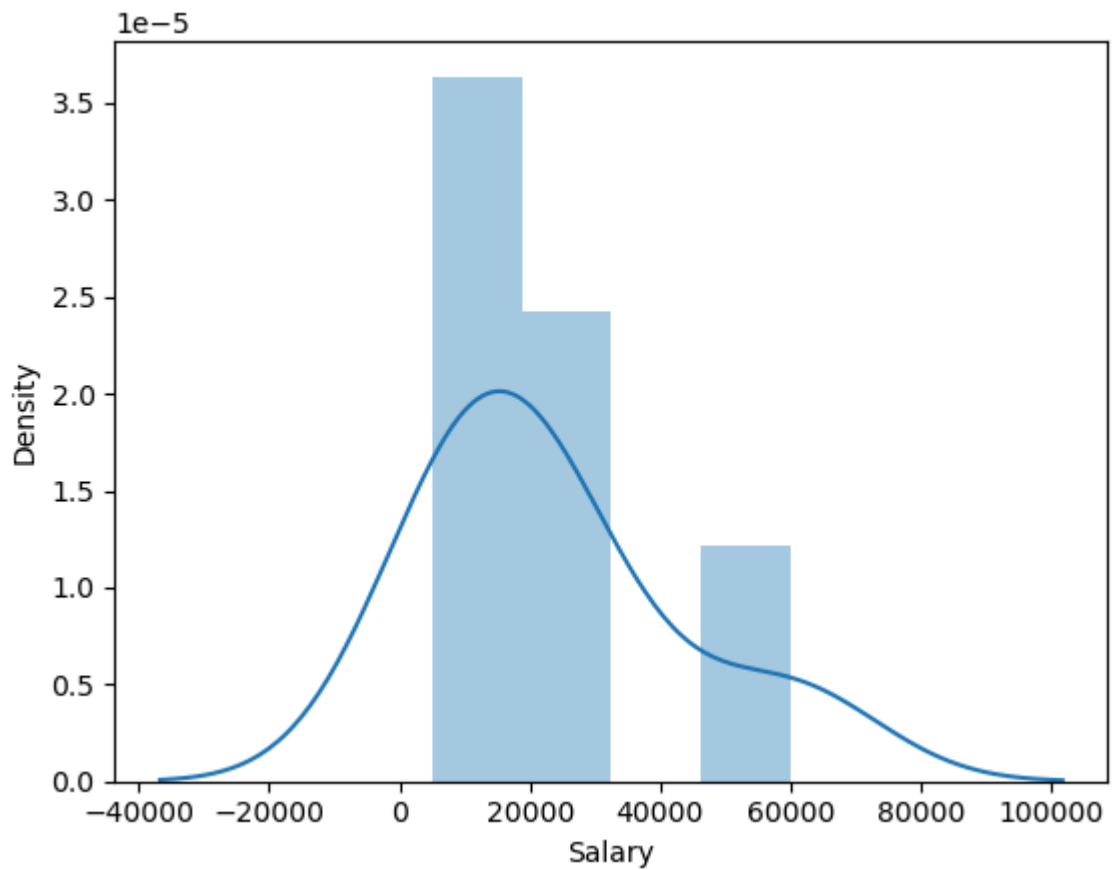```

Out[142…   `'C:\\Users\\samik\\python project'`

# EDA technologies

In [145…
```python
import matplotlib.pyplot as plt
import seaborn as sns
```

In [147…
```python
import warnings
warnings.filterwarnings('ignore')
```
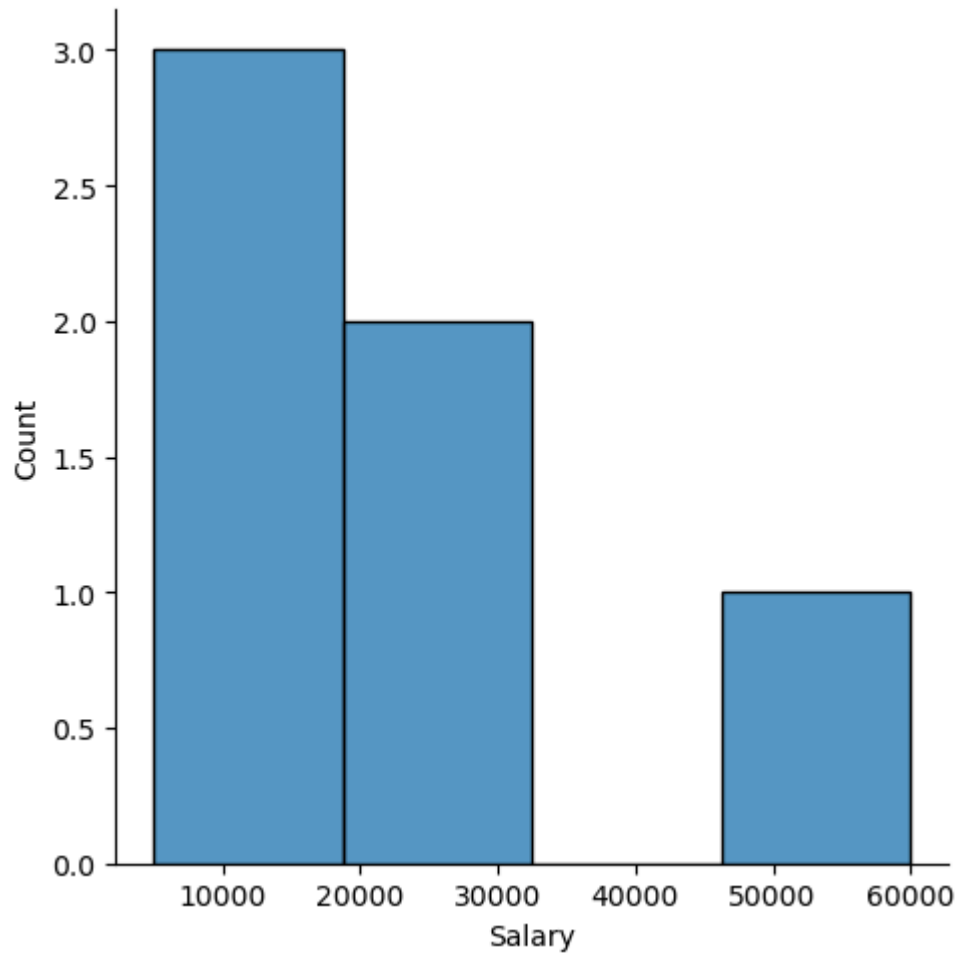
In [149…  `clean_data['Salary']`

Out[149…
```
0     5000
1    10000
2    15000
3    20000
4    30000
5    60000
Name: Salary, dtype: int32
```
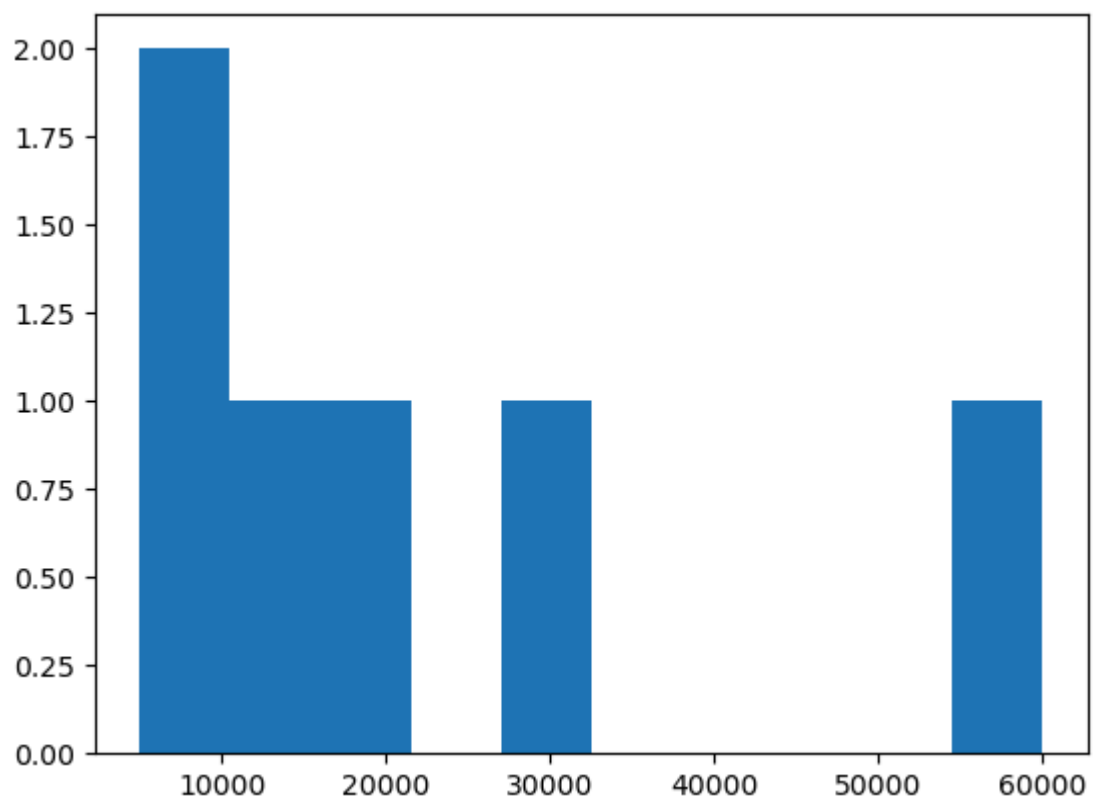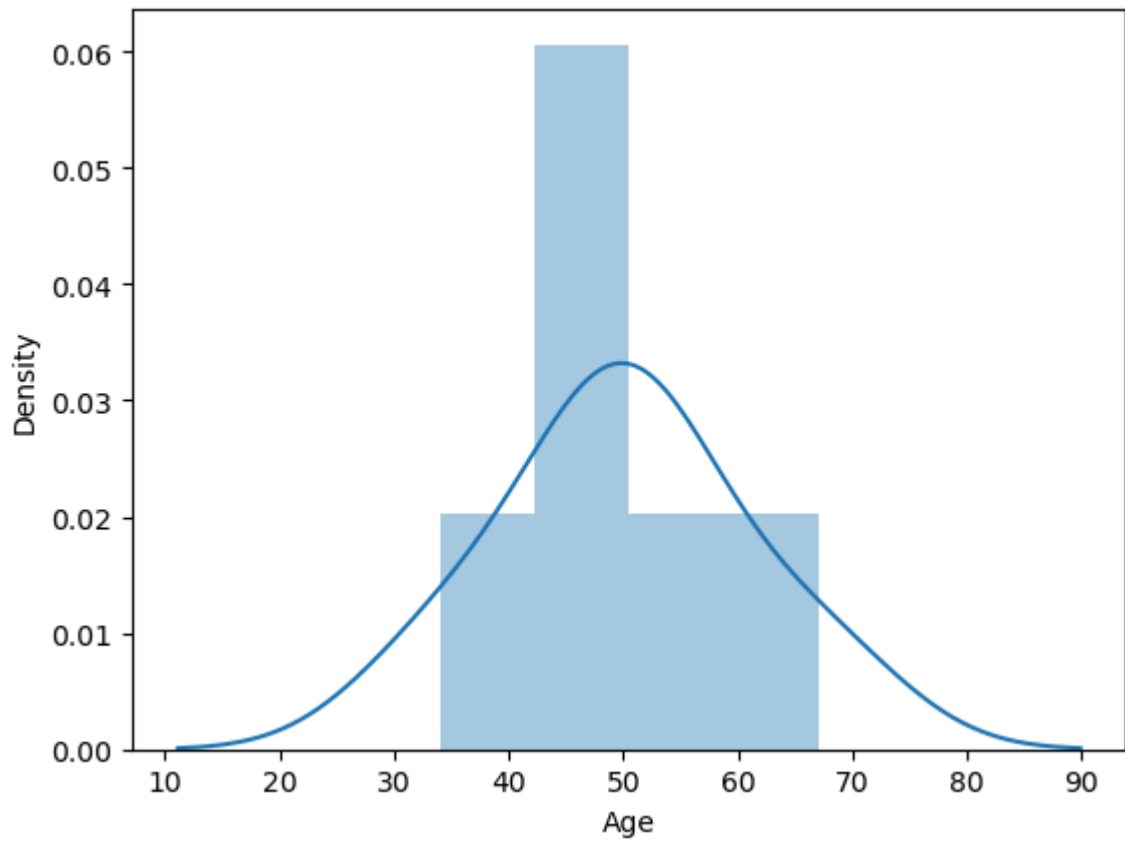
In [151…  `vis1 = sns.distplot(clean_data['Salary'])`
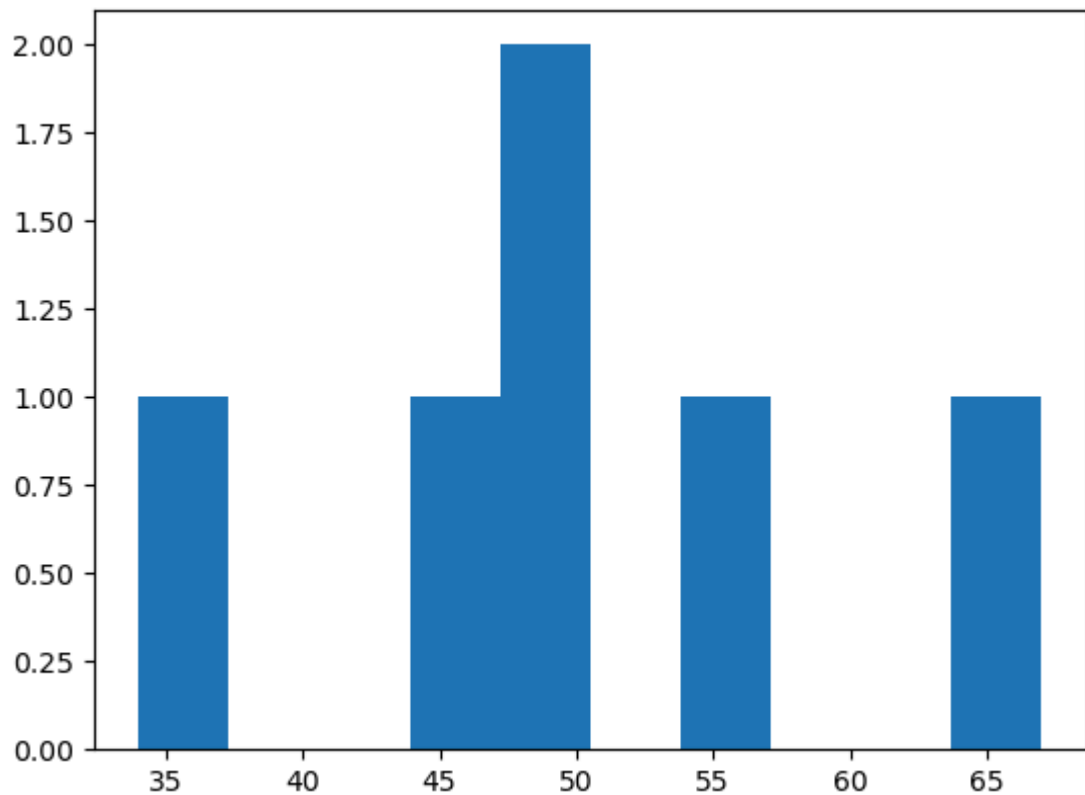


In [153…  `vis2 = sns.displot(clean_data['Salary'])`
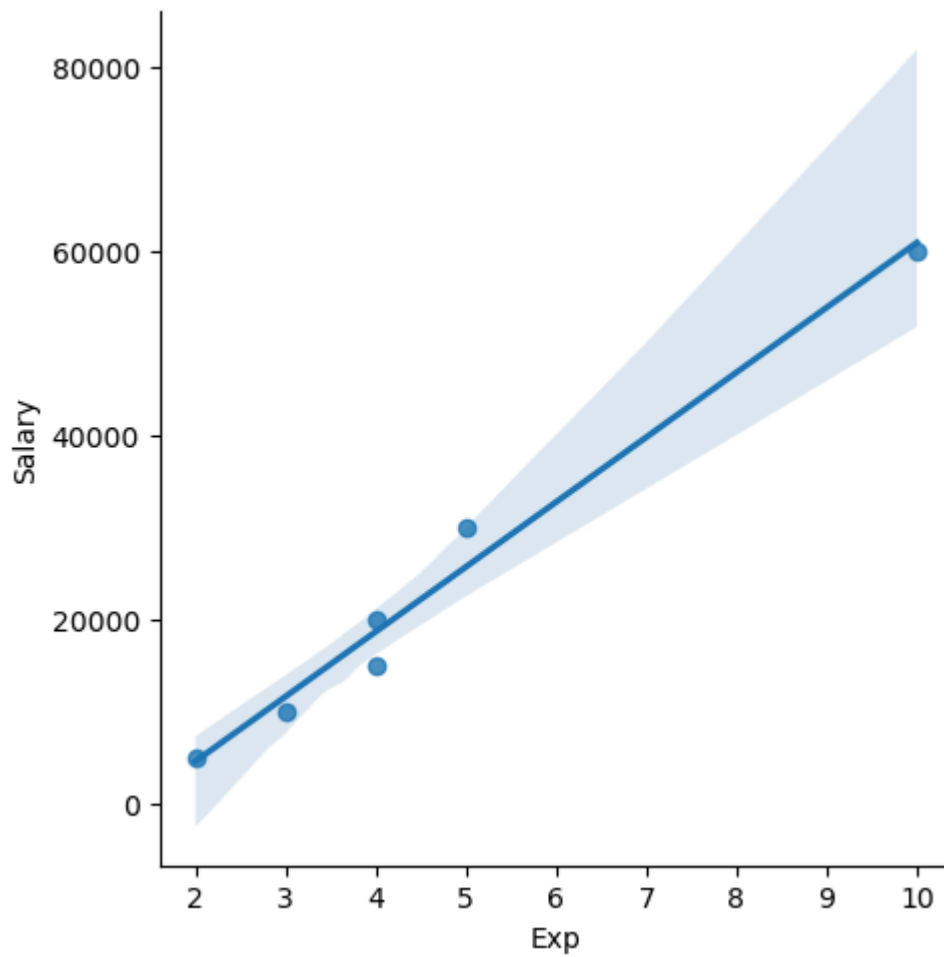
```
In [155…   vis3 = plt.hist(clean_data['Salary'])
```
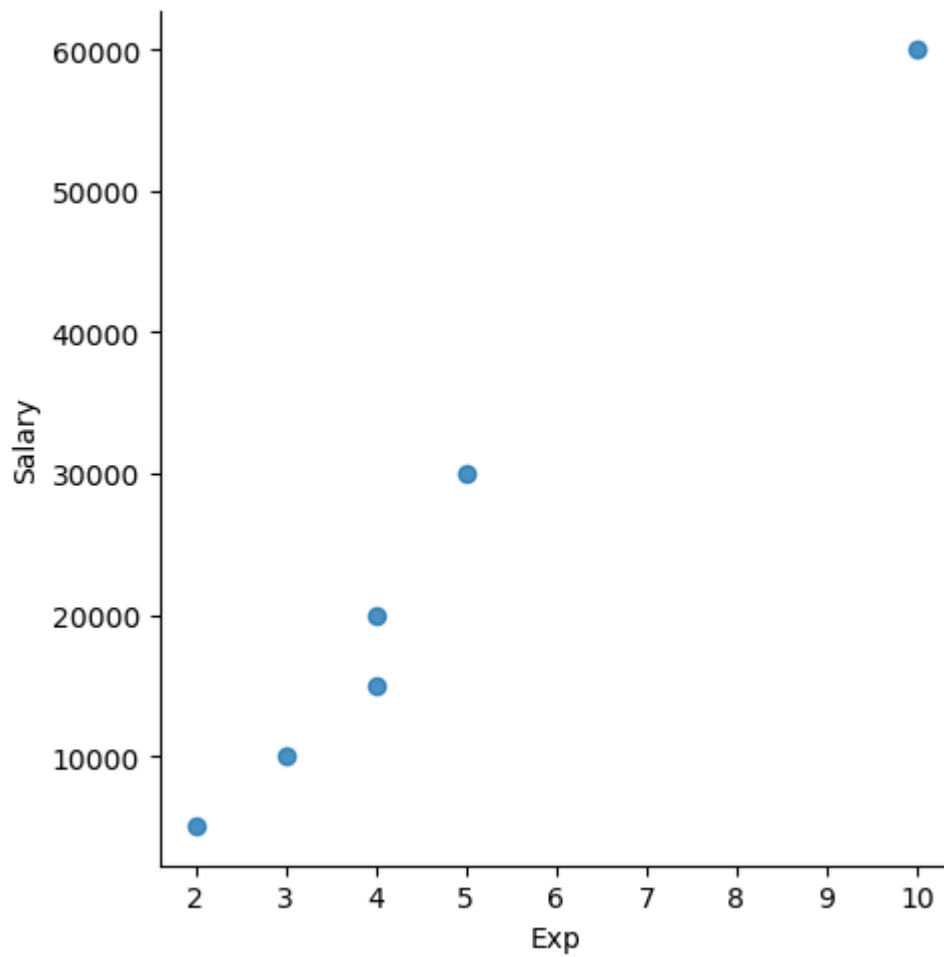


```
In [157…   vis4 = sns.distplot(clean_data['Age'])
```

In [159… `vis5= plt.hist(clean_data['Age'])`
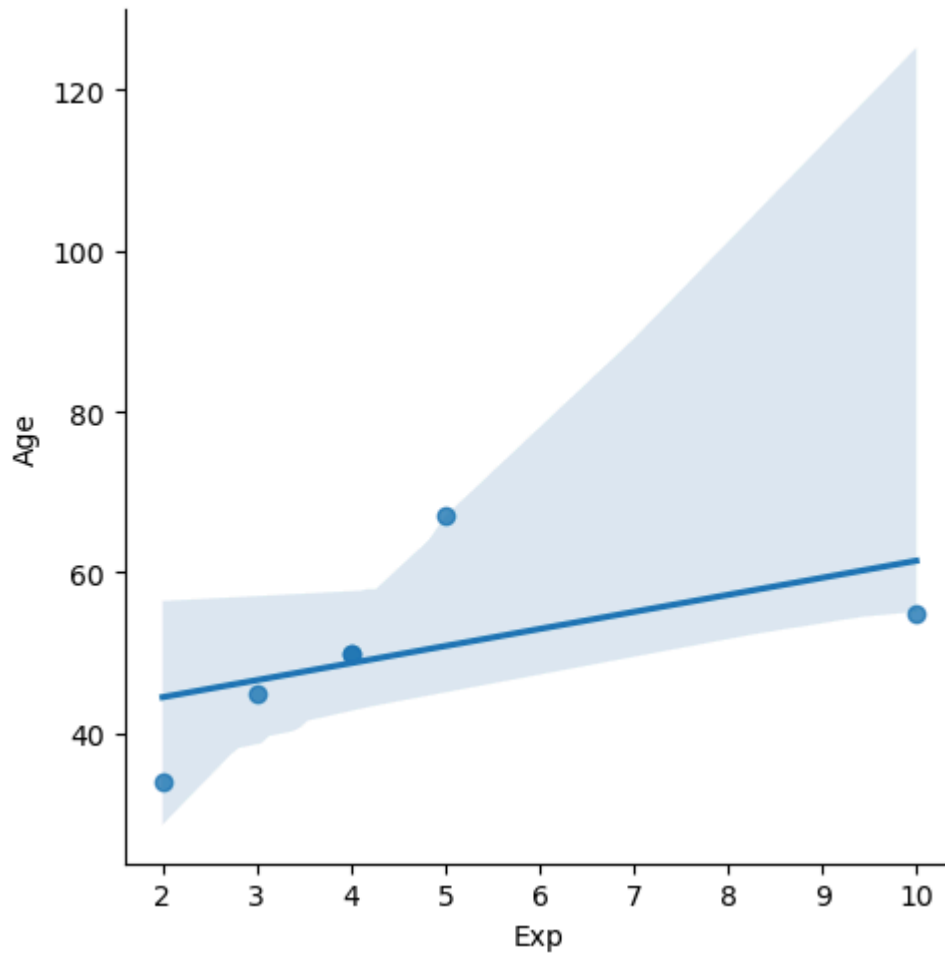


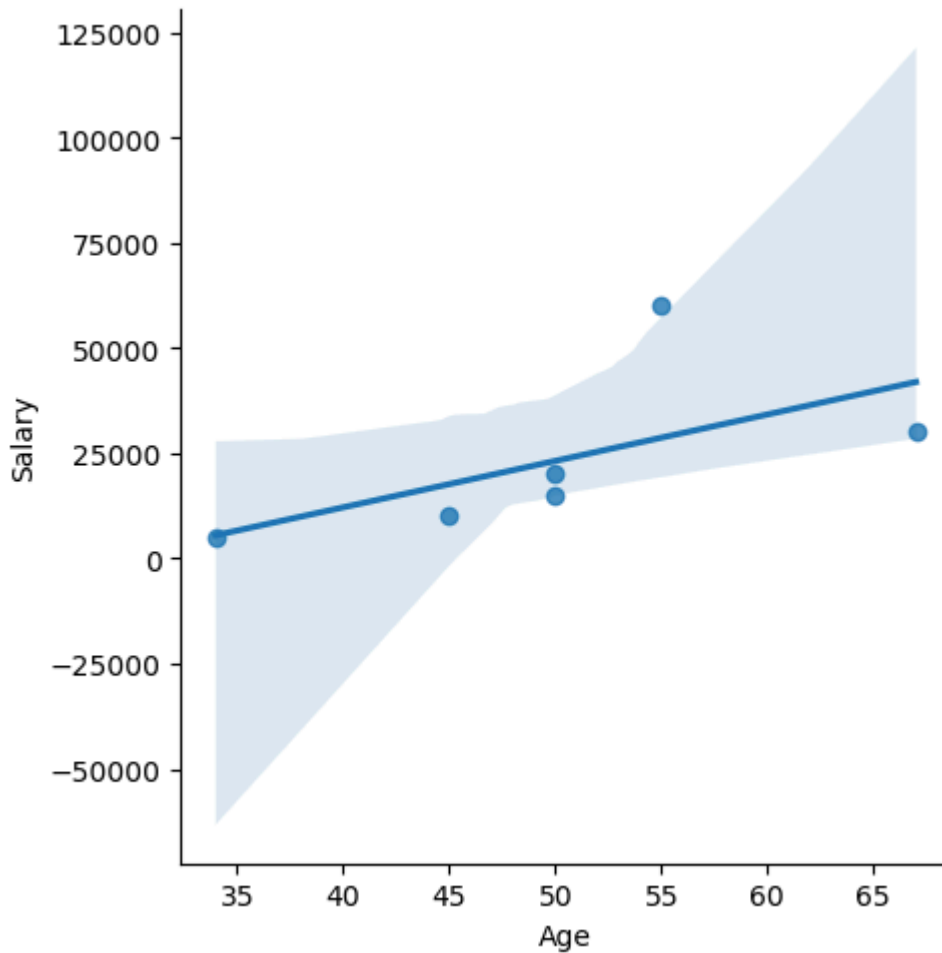In [161… `vis6 = sns.lmplot(data = clean_data , x = 'Exp', y ='Salary')`

```
In [163...    vis7 = sns.lmplot(data = clean_data , x = 'Exp', y ='Salary',fit_reg = False)
```

```
In [165…   vis8 = sns.lmplot(data = clean_data , x = 'Exp', y ='Age')
```

```
In [167…   vis9= sns.lmplot(data = clean_data , x = 'Age', y ='Salary')
```

```
In [169…   clean_data[:]
```

Out[169…

|   | 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 [171…   clean_data[0:6:2]
```

Out[171…

|   | 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 [173…   clean_data[0:8:3]
```

Out[173...

| | Name | Domain | Age | Location | Salary | Exp |
|---|---|---|---|---|---|---|
| **0** | Mike | Datascience | 34 | Mumbai | 5000 | 2 |
| **3** | Jane | Analytics | 50 | Hyderbad | 20000 | 4 |

In [175...
```python
clean_data[::-1]
```

Out[175...

| | 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 [177...
```python
clean_data.columns
```

Out[177...    Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')

In [179...
```python
X_iv = clean_data[['Name', 'Domain', 'Age', 'Location', 'Exp']]
```

In [181...
```python
X_iv
```

Out[181...

| | 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 [183...
```python
Y_dv = clean_data[['Salary']]
```

In [185...
```python
Y_dv
```

Out[185…

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

In [187…    `emp`

Out[187…

|   | 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 [189…    `clean_data`

Out[189…

|   | 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 [191…
```python
imputation = pd.get_dummies(clean_data)
```

In [193…
```python
imputation
```

Out[193…

| | 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 [195… `imputation.astype(int)`

Out[195…

| | Age | Salary | Exp | Name_Jane | Name_Kim | Name_Mike | Name_Teddy | Name_Umar |
|---|---|---|---|---|---|---|---|---|
| 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 [197… `imputation.columns`

Out[197…
```
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_Hyderbad', 'Location_Mumbai'],
      dtype='object')
```

In [199… `len(imputation)`

Out[199… `6`

In [201… `imputation.shape`

Out[201… `(6, 19)`

In [203… `len(imputation.columns)`

Out[203… `19`

In [ ]: