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

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

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

```python
In [3]:  emp = pd.read_excel(r'C:\Users\UMA SESHA KUMARI\Downloads\rawdata.xlsx ')
```

```python
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+ |

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

```
Out[5]:  2633150196496
```

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

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

```python
In [7]:  emp.shape
```

```
Out[7]:  (6, 6)
```

```python
In [8]:  emp.head()
```

Out[8]:

|   | 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 |

```python
In [9]:  emp.tail()
```

Out[9]:

| | 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 [10]: `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 [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.isna()`

Out[12]:

| | 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 [13]:
```python
emp.isnull().sum()
```

Out[13]:
```
Name        0
Domain      0
Age         2
Location    2
Salary      0
Exp         1
dtype: int64
```

DATA CLEANSING or DATA CLEANING

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

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

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

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

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

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

localhost:8888/doc/tree/Project 5 EDA Exploratory DATA.ipynb?

3/13

```
Out[18]:  0        Datascience#$
          1             Testing
          2        Dataanalyst^^#
          3          Ana^^lytics
          4           Statistics
          5                 NLP
          Name: Domain, dtype: object
```

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

```
In [20]:  emp['Domain']
```

```
Out[20]:  0     Datascience
          1         Testing
          2      Dataanalyst
          3        Analytics
          4       Statistics
          5              NLP
          Name: Domain, dtype: object
```

```
In [21]:  emp.isnull().sum()
```

```
Out[21]:  Name        0
          Domain      0
          Age         2
          Location    2
          Salary      0
          Exp         1
          dtype: int64
```

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

```
Out[22]:  0      34 years
          1       45' yr
          2          NaN
          3          NaN
          4        67-yr
          5         55yr
          Name: Age, 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.replace(r'\W','',regex=True)
```

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

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

```
In [26]:  emp['Age'] = emp['Age'].str.extract('(\d+)')
```

```
<>:1: SyntaxWarning: invalid escape sequence '\d'
<>:1: SyntaxWarning: invalid escape sequence '\d'
C:\Users\UMA SESHA KUMARI\AppData\Local\Temp\ipykernel_9160\1884116463.py:1: SyntaxW
arning: invalid escape sequence '\d'
  emp['Age'] = emp['Age'].str.extract('(\d+)')
```

```
In [27]:  emp['Age']
```

```
Out[27]:  0      34
          1      45
          2     NaN
          3     NaN
          4      67
          5      55
          Name: Age, dtype: object
```

```
In [28]:  emp['Location']
```

```
Out[28]:  0       Mumbai
          1    Bangalore
          2          NaN
          3     Hyderbad
          4          NaN
          5        Delhi
          Name: Location, dtype: object
```

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

```
In [30]:  emp['Domain'] = emp['Domain'].str.replace(r'\W','',regex=True)
```

```
In [31]:  emp['Salary']
```

```
Out[31]:  0      5^00#0
          1     10%%000
          2     1$5%000
          3      2000^0
          4      30000-
          5     6000^$0
          Name: Salary, dtype: object
```

```
In [32]:  emp['Salary'] = emp['Salary'].str.replace(r'\W','',regex=True)
```

```
In [33]:  emp
```

Out[33]:

| | 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 [34]:
```python
emp['Exp'] = emp['Exp'].str.extract('(\d+)')
```

```
<>:1: SyntaxWarning: invalid escape sequence '\d'
<>:1: SyntaxWarning: invalid escape sequence '\d'
C:\Users\UMA SESHA KUMARI\AppData\Local\Temp\ipykernel_9160\3836251810.py:1: SyntaxW
arning: invalid escape sequence '\d'
  emp['Exp'] = emp['Exp'].str.extract('(\d+)')
```

In [35]:
```python
emp['Exp']
```

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

In [36]:
```python
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]:
```python
clean_data = emp
```

In [38]:
```python
clean_data
```

Out[38]:

| | 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 [39]:
```
id(clean_data)
```

Out[39]:  2633150196496

In [40]:
```
id(emp)
```

Out[40]:  2633150196496

### EDA TECHNIQUES

In [42]:
```
clean_data
```

Out[42]:

| | 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 [43]:
```
clean_data.isnull()
```

Out[43]:

| | 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 [44]:  clean_data.isnull().sum()
```

```
Out[44]:  Name         0
          Domain       0
          Age          2
          Location     2
          Salary       0
          Exp          1
          dtype: int64
```

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

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

```
In [46]:  import numpy as np
```

```
In [47]:  clean_data['Age'] = clean_data['Age'].fillna(np.mean(pd.to_numeric(clean_data['Age'
```

```
In [95]:  clean_data['Age']
```

```
Out[95]:  0        34
          1        45
          2     50.25
          3     50.25
          4        67
          5        55
          Name: Age, dtype: object
```

```
In [99]:  clean_data['Exp'] = clean_data['Exp'].fillna(np.mean(pd.to_numeric(clean_data['Exp'
```

```
In [101…  clean_data['Exp']
```

```
Out[101…  0      2
          1      3
          2      4
          3    4.8
          4      5
          5     10
          Name: Exp, dtype: object
```

```
In [103…  clean_data['Location']
```

Out[103...    0        Mumbai
             1     Bangalore
             2          NaN
             3     Hyderbad
             4          NaN
             5        Delhi
             Name: Location, dtype: object

In [105...  ```python
clean_data['Location'] = clean_data['Location'].fillna(clean_data['Location'].mode(
```

In [107...  ```python
clean_data['Location']
```

Out[107...    0        Mumbai
             1     Bangalore
             2     Bangalore
             3     Hyderbad
             4     Bangalore
             5        Delhi
             Name: Location, dtype: object

In [110...  ```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 [116...  ```python
clean_data['Age'] = clean_data['Age'].astype(int)
```

In [118...  ```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      object
 5   Exp       6 non-null      object
dtypes: int32(1), object(5)
memory usage: 396.0+ bytes
```

In [124...  ```python
clean_data['Salary'] = clean_data['Salary'].astype(int)
```

In [126...
```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      object
dtypes: int32(2), object(4)
memory usage: 372.0+ bytes
```

In [128...
```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 [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      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      object
dtypes: category(3), int32(2), object(1)
memory usage: 890.0+ bytes
```

In [134...
```python
clean_data
```

Out[134...

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

In [136...
```python
type(clean_data)
```

Out[136...
```
pandas.core.frame.DataFrame
```

In [138… 
```python
clean_data.to_csv('clean_data.csv')
```

In [140… 
```python
import os
os.getcwd()
```

Out[140… 
```
'C:\\Users\\UMA SESHA KUMARI'
```

In [150… 
```python
import matplotlib.pyplot as plt
import seaborn as sns
```

In [154… 
```python
vis1 = sns.distplot(clean_data['Salary'])
```



In [158… 
```python
vis2 = plt.hist(clean_data['Salary'])
```
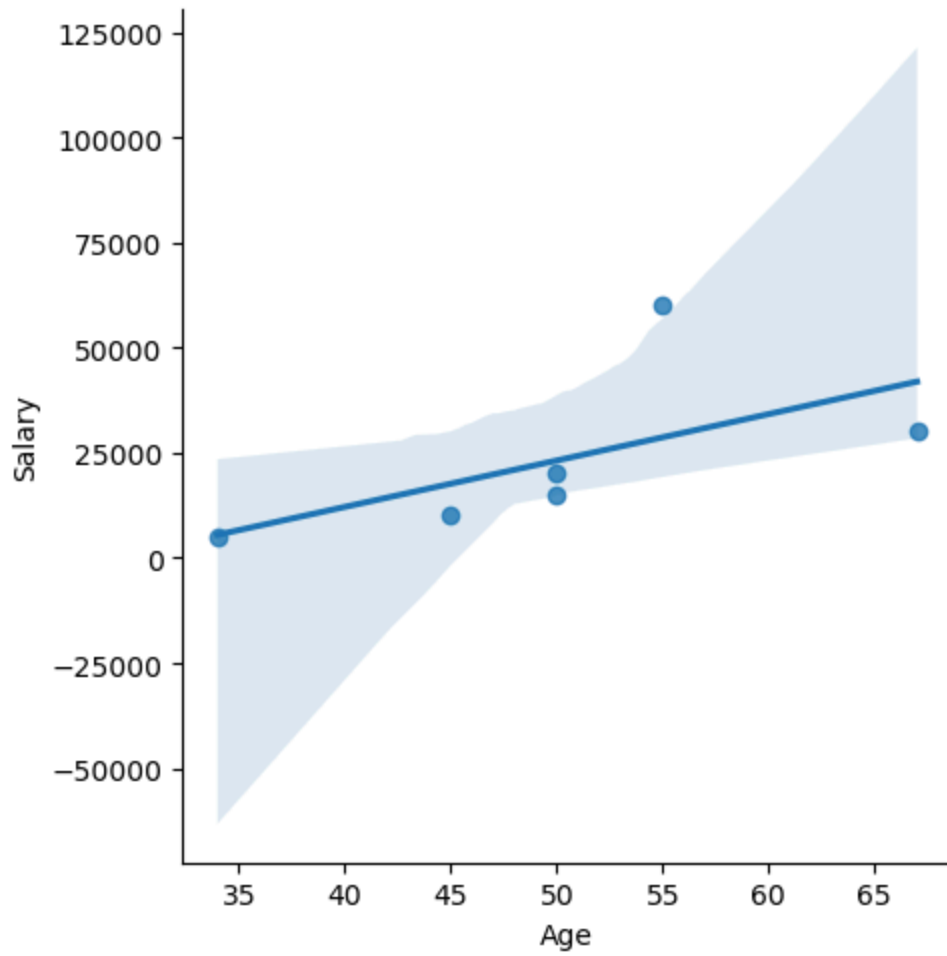
In [160…  `clean_data`

Out[160…

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

In [170…  `clean_data.columns`

Out[170…  `Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')`

In [182…  `vix3 = sns.lmplot(data = clean_data , x = 'Age', y = 'Salary')`

In [ ]: