

Raw data to Clean data conversion using Python EDA

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

```
In [3]: emp = pd.read_excel(r"C:\Users\ANITHA\OneDrive\Desktop\DATASCIENCE NOTES\Rawdata (1)
```

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

Out[4]:

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

```
Out[5]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')
```

```
In [6]: emp.shape
```

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

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

Out[8]:

| | 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 |

```
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

Out[11]:

| | 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 [12]:

emp['Domain']

Out[12]:

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

In [13]:

emp.isnull()

| | 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 [14]: `emp.isnull().sum()`

Out[14]:

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

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

Out[15]:

| | |
|---------------------------|--------|
| 0 | Mike |
| 1 | Teddy^ |
| 2 | Uma#r |
| 3 | Jane |
| 4 | Uttam* |
| 5 | Kim |
| Name: Name, dtype: object | |

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

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

Out[26]:

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

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

Out[18]:

| | |
|-----------------------------|----------------|
| 0 | Datascience#\$ |
| 1 | Testing |
| 2 | Dataanalyst^^# |
| 3 | Ana^^lytics |
| 4 | Statistics |
| 5 | NLP |
| Name: Domain, dtype: object | |

```
In [28]: emp['Domain'] = emp['Domain'].str.replace(r'\W+', ' ', regex=True)
```

```
In [29]: emp['Domain']
```

```
Out[29]: 0    Datascienc
          1        Testing
          2   Dataanalyst
          3     Analytics
          4   Statistics
          5         NLP
Name: Domain, dtype: object
```

```
In [30]: emp
```

| | 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 | Umar | Dataanalyst | NaN | NaN | 1\$5%000 | 4> yrs |
| 3 | Jane | Analytics | NaN | Hyderabad | 2000^0 | NaN |
| 4 | Uttam | Statistics | 67-yr | NaN | 30000- | 5+ year |
| 5 | Kim | NLP | 55yr | Delhi | 6000^\$0 | 10+ |

```
In [33]: emp['Age'] = emp['Age'].str.replace(r'\W+', ' ', regex = True)
```

```
In [34]: emp['Age']
```

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

```
In [35]: emp['Age'] = emp['Age'].str.extract('(\d+)')
```

```
In [36]: emp['Age']
```

```
Out[36]: 0    34
          1    45
          2    NaN
          3    NaN
          4    67
          5    55
Name: Age, dtype: object
```

```
In [37]: emp
```

Out[37]:

| | Name | Domain | Age | Location | Salary | Exp |
|---|-------|-------------|-----|-----------|----------|---------|
| 0 | Mike | Datascienc | 34 | Mumbai | 5^00#0 | 2+ |
| 1 | Teddy | Testing | 45 | Bangalore | 10%0000 | <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 [38]: `emp['Location'] = emp['Location'].str.replace(r'\W+', '')`In [39]: `emp['Location']`

Out[39]:

| | |
|---|-----------|
| 0 | Mumbai |
| 1 | Bangalore |
| 2 | NaN |
| 3 | Hyderbad |
| 4 | NaN |
| 5 | Delhi |

Name: Location, dtype: object

In [40]: `emp['Salary']`

Out[40]:

| | |
|---|----------|
| 0 | 5^00#0 |
| 1 | 10%0000 |
| 2 | 1\$5%000 |
| 3 | 2000^0 |
| 4 | 30000- |
| 5 | 6000^\$0 |

Name: Salary, dtype: object

In [43]: `emp['Salary'] = emp['Salary'].str.replace(r'\w+', '', regex = True)`In [44]: `emp['Salary']`

Out[44]:

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

Name: Salary, dtype: object

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

```
Out[45]: 0      2+
         1      <3
         2    4> yrs
         3      NaN
         4    5+ year
         5    10+
Name: Exp, dtype: object
```

```
In [46]: emp['Exp'] = emp['Exp'].str.extract('(\d+)')
```

```
In [47]: emp['Exp']
```

```
Out[47]: 0      2
         1      3
         2      4
         3      NaN
         4      5
         5     10
Name: Exp, dtype: object
```

```
In [48]: 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 [49]: clean_data = emp.copy()
```

```
In [50]: 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 [51]: clean_data['Age']
```

```
Out[51]: 0    34
         1    45
         2    NaN
         3    NaN
         4    67
         5    55
Name: Age, dtype: object
```

```
In [52]: import numpy as np
```

```
In [53]: clean_data['Age'] = clean_data['Age'].fillna(np.mean(pd.to_numeric(clean_data['Age'])))
```

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

```
Out[54]: 0    34
         1    45
         2    50.25
         3    50.25
         4    67
         5    55
Name: Age, dtype: object
```

```
In [55]: clean_data['Exp']
```

```
Out[55]: 0    2
         1    3
         2    4
         3    NaN
         4    5
         5    10
Name: Exp, dtype: object
```

```
In [57]: clean_data['Exp'] = clean_data['Exp'].fillna(np.mean(pd.to_numeric(clean_data['Exp'])))
```

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

```
Out[58]: 0    2
         1    3
         2    4
         3    4.8
         4    5
         5    10
Name: Exp, dtype: object
```

```
In [59]: clean_data
```

Out[59]:

| | Name | Domain | Age | Location | Salary | Exp |
|---|-------|-------------|-------|-----------|--------|-----|
| 0 | Mike | Datascienc | 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 [60]: `clean_data['Location'].isnull().sum()`Out[60]: `np.int64(2)`In [61]: `clean_data['Location']`

```
Out[61]: 0      Mumbai
         1      Bangalore
         2      NaN
         3      Hyderbad
         4      NaN
         5      Delhi
Name: Location, dtype: object
```

In [62]: `clean_data['Location'] = clean_data['Location'].fillna(clean_data['Location'].mode()`In [63]: `clean_data['Location']`

```
Out[63]: 0      Mumbai
         1      Bangalore
         2      Bangalore
         3      Hyderbad
         4      Bangalore
         5      Delhi
Name: Location, dtype: object
```

In [64]: `clean_data`

| | Name | Domain | Age | Location | Salary | Exp |
|---|-------|-------------|-------|-----------|--------|-----|
| 0 | Mike | Datascienc | 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 [65]: 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 [66]: clean_data['Age'] = clean_data['Age'].astype(int)
```

```
In [67]: 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    int64  
 3   Location   6 non-null    object  
 4   Salary     6 non-null    object  
 5   Exp        6 non-null    object  
dtypes: int64(1), object(5)
memory usage: 420.0+ bytes
```

```
In [69]: clean_data['Salary'] = clean_data['Salary'].astype(int)
clean_data['Exp'] = clean_data['Exp'].astype(int)
```

```
In [70]: 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    int64  
 3   Location   6 non-null    object  
 4   Salary     6 non-null    int64  
 5   Exp        6 non-null    int64  
dtypes: int64(3), object(3)
memory usage: 420.0+ bytes
```

```
In [71]: 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 [72]: 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      int64   
 3   Location    6 non-null      category
 4   Salary      6 non-null      int64   
 5   Exp         6 non-null      int64   
dtypes: category(3), int64(3)
memory usage: 938.0 bytes
```

```
In [73]: clean_data
```

```
Out[73]:
```

| | 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 |

```
In [75]: clean_data.to_csv('clean_data.csv')
```

```
In [76]: import os
os.getcwd()
```

```
Out[76]: 'C:\\\\Users\\\\ANITHA'
```

```
In [77]: clean_data
```

Out[77]:

| | Name | Domain | Age | Location | Salary | Exp |
|---|-------|-------------|-----|-----------|--------|-----|
| 0 | Mike | Datascienc | 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 |

EDA TECHNIQUE LETS APPLY

In [78]:

```
import matplotlib.pyplot as plt # visualization
import seaborn as sns
```

In [80]:

```
import warnings
warnings.filterwarnings('ignore')
```

In [81]:

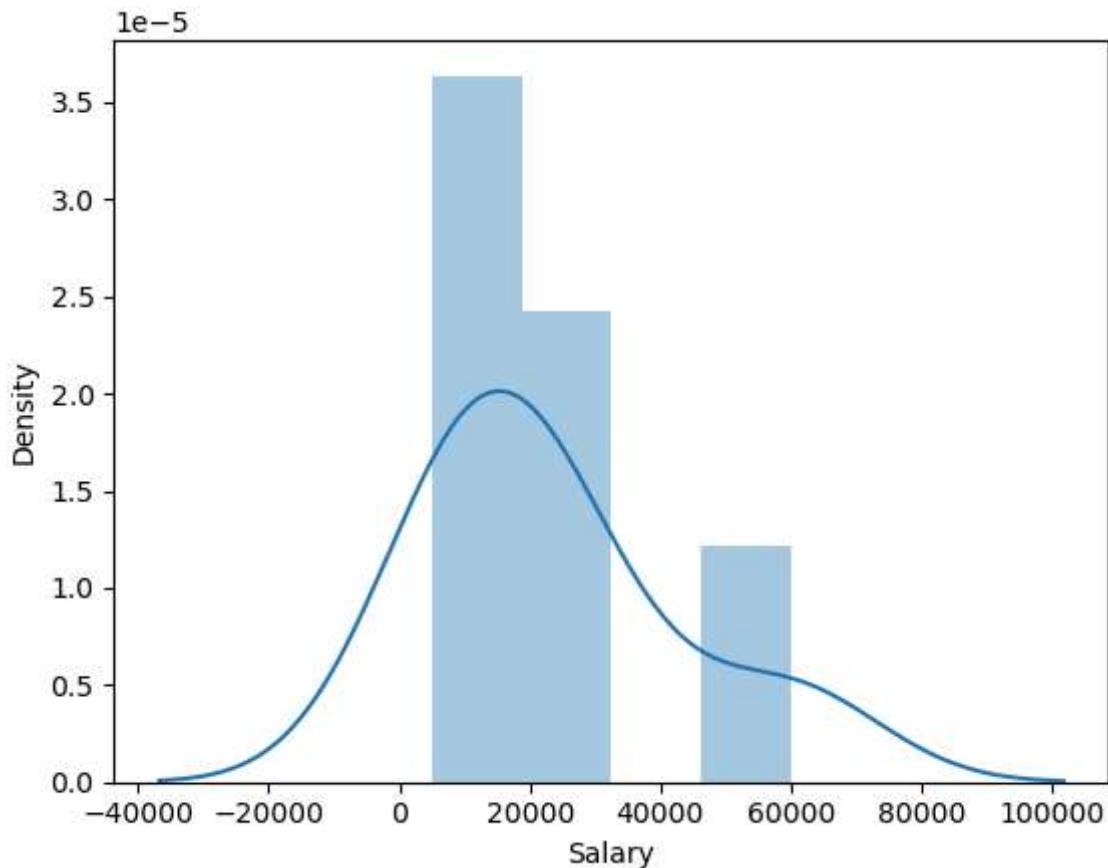
```
clean_data['Salary']
```

Out[81]:

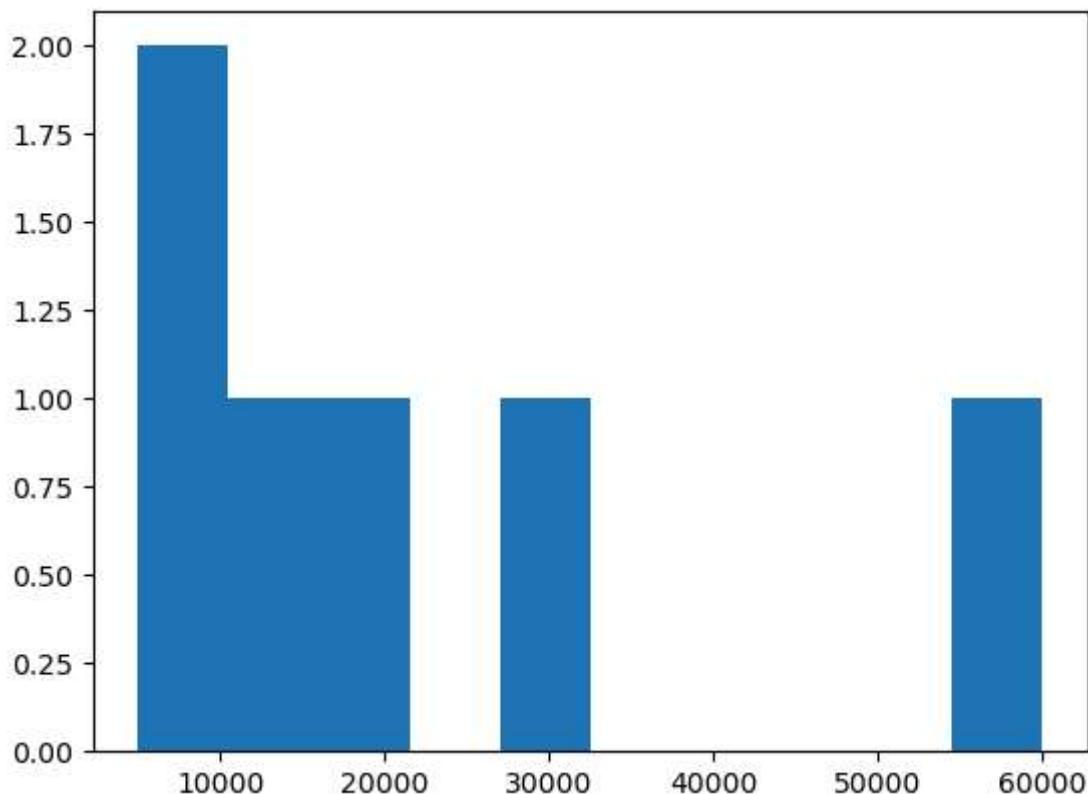
```
0    5000
1    10000
2    15000
3    20000
4    30000
5    60000
Name: Salary, dtype: int64
```

In [82]:

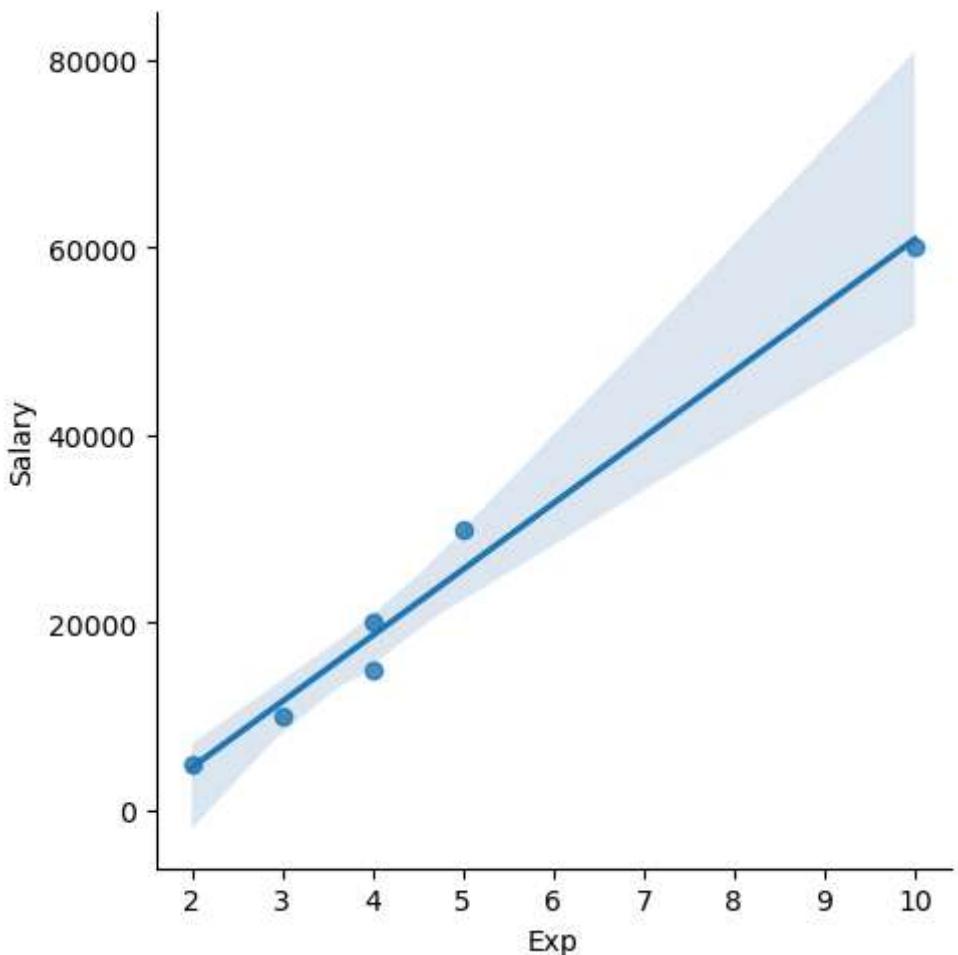
```
is1 = sns.distplot(clean_data['Salary'])
```



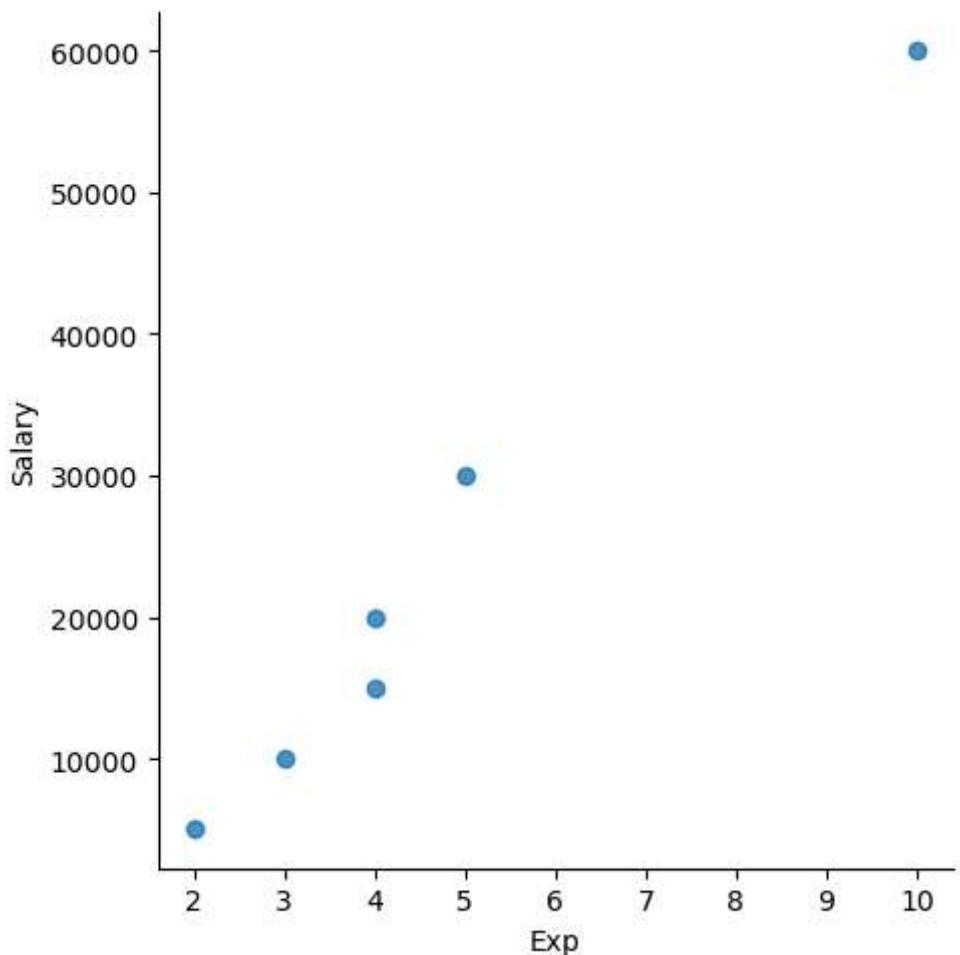
```
In [83]: vis2 = plt.hist(clean_data['Salary'])
```



```
In [84]: vis4 = sns.lmplot(data=clean_data,x = 'Exp', y='Salary')
```



```
In [85]: vis5 = sns.lmplot(data=clean_data,x = 'Exp', y='Salary', fit_reg = False)
```



```
In [86]: clean_data[:]
```

```
Out[86]:
```

| | 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 |

```
In [87]: clean_data[0:6:2]
```

```
Out[87]:
```

| | 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 [89]: `clean_data[:::-1]`

Out[89]:

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

In [90]: `clean_data.columns`

Out[90]: `Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')`

In [91]: `X_iv = clean_data[['Name', 'Domain', 'Age', 'Location', 'Exp']]`

In [92]: `X_iv`

Out[92]:

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

In [93]: `y_dv = clean_data[['Salary']]`

In [94]: `y_dv`

Out[94]:

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

In [95]: 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 [96]: 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 [97]: 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 |

In [98]: y_dv

Out[98]:

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

In [99]:

clean_data

Out[99]:

| | 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 |

In [100...]

imputation = pd.get_dummies(clean_data)

In [101...]

imputation

Out[101...]

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



In [102...]

clean_data

Out[102...]

| | Name | Domain | Age | Location | Salary | Exp |
|---|-------|-------------|-----|-----------|--------|-----|
| 0 | Mike | Datascienc | 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 |

In [103...]

imputation

Out[103...]

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

raw data with lot of regex, missing, uncleandata**regex, clean****fill missing numerical & cateigroica****clean_dataset (data cleaning) 3 month - 5mont****outlier treatement, univati, bivariate, corelation****split the data into x_i.v & y_dv****impute cateogrica data to numerical**

In []: