

```
In [1]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
```

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

```
Out[2]: '2.2.3'
```

```
In [3]: emp = pd.read_excel(r"C:\Users\rajendra damahe\Downloads\Rawdata.xlsx")
```

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

	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]: id(emp)
```

```
Out[5]: 2424344099008
```

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

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

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

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

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

	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 [9]: `emp.tail()`

	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()`

	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()`

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

# DATA CLEANNING

```
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', ' ', regex=True)
```

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

```
Out[15]: 0      Mike
          1    Teddy
          2    Umar
          3    Jane
          4    Uttam
          5      Kim
Name: Name, dtype: object
```

```
In [16]: emp.columns
```

```
Out[16]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')
```

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

```
Out[17]: 0      Datascience$$
          1        Testing
          2  Dataanalyst^^#
          3    Ana^^lytics
          4      Statistics
          5        NLP
Name: Domain, dtype: object
```

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

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

```
Out[21]: 0      Mumbai
         1    Bangalore
         2        NaN
         3   Hyderbad
         4        NaN
         5      Delhi
Name: Location, dtype: object
```

```
In [22]: emp['Age'] = emp['Age'].str.extract('(\d+)') # r(r'(\d+)')
```

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

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

```
In [24]: emp
```

	Name	Domain	Age	Location	Salary	Exp
<b>0</b>	Mike	Datascience	34	Mumbai	5^00#0	2+
<b>1</b>	Teddy	Testing	45	Bangalore	10%0000	<3
<b>2</b>	Umar	Dataanalyst	NaN	NaN	1\$5%0000	4> yrs
<b>3</b>	Jane	Analytics	NaN	Hyderbad	2000^0	NaN
<b>4</b>	Uttam	Statistics	67	NaN	30000-	5+ year
<b>5</b>	Kim	NLP	55	Delhi	6000^\$0	10+

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

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

```
Out[26]: 0    5000
         1   10000
         2   15000
         3   20000
         4   30000
         5   60000
Name: Salary, dtype: object
```

```
In [27]: emp['Exp'] = emp['Exp'].str.extract('(\d+)')
```

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

```
Out[28]: 0      2
         1      3
         2      4
         3    NaN
         4      5
         5     10
Name: Exp, dtype: object
```

```
In [29]: clean_data=emp.copy()
```

```
In [30]: clean_data
```

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

```
In [31]: clean_data.isnull().sum()
```

```
Out[31]: Name      0
          Domain    0
          Age       2
          Location  2
          Salary    0
          Exp       1
          dtype: int64
```

```
In [32]: clean_data['Age']
```

```
Out[32]: 0      34
         1      45
         2    NaN
         3    NaN
         4      67
         5      55
Name: Age, dtype: object
```

```
In [39]: clean_data['Age'] = clean_data['Age'].fillna(np.mean(pd.to_numeric(clean_data['Age'])))
```

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

```
Out[40]: 0      34
         1      45
         2    50.25
         3    50.25
         4      67
         5      55
Name: Age, dtype: object
```

```
In [41]: clean_data['Exp']
```

```
Out[41]: 0      2
         1      3
         2      4
         3    4.8
         4      5
         5     10
Name: Exp, dtype: object
```

```
In [42]: clean_data['Exp'] = clean_data['Exp'].fillna(np.mean(pd.to_numeric(clean_data['Exp'])))
```

```
In [43]: clean_data['Exp']
```

```
Out[43]: 0      2
         1      3
         2      4
         3    4.8
         4      5
         5     10
Name: Exp, dtype: object
```

```
In [44]: clean_data
```

```
Out[44]:   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  Hyderabad  20000    4.8
4  Uttam      Statistics  67      NaN  30000     5
5    Kim        NLP  55    Delhi  60000    10
```

```
In [45]: clean_data['Location'].isnull().sum()
```

```
Out[45]: np.int64(2)
```

```
In [46]: clean_data['Location']
```

```
Out[46]: 0      Mumbai
         1      Bangalore
         2      NaN
         3      Hyderabad
         4      NaN
         5      Delhi
Name: Location, dtype: category
Categories (4, object): ['Bangalore', 'Delhi', 'Hyderabad', 'Mumbai']
```

```
In [47]: clean_data['Location'] = clean_data['Location'].fillna(clean_data['Location'].mode())
```

```
In [48]: clean_data['Location']
```

```
Out[48]: 0      Mumbai
         1      Bangalore
         2      Bangalore
         3      Hyderabad
         4      Bangalore
         5      Delhi
Name: Location, dtype: category
Categories (4, object): ['Bangalore', 'Delhi', 'Hyderabad', 'Mumbai']
```

```
In [49]: 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	50.25	Bangalore	15000	4
3	Jane	Analytics	50.25	Hyderabad	20000	4.8
4	Uttam	Statistics	67	Bangalore	30000	5
5	Kim	NLP	55	Delhi	60000	10

```
In [50]: 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    category
 2   Age        6 non-null    object  
 3   Location   6 non-null    category
 4   Salary     6 non-null    object  
 5   Exp        6 non-null    object  
dtypes: category(2), object(4)
memory usage: 760.0+ bytes
```

```
In [51]: # 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 [52]: `clean_data.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6 entries, 0 to 5
Data columns (total 6 columns):
 #   Column      Non-Null Count  Dtype  
---  --          -----          object 
 0   Name        6 non-null     object  
 1   Domain      6 non-null     category
 2   Age         6 non-null     object  
 3   Location    6 non-null     category
 4   Salary      6 non-null     object  
 5   Exp         6 non-null     object  
dtypes: category(2), object(4)
memory usage: 760.0+ bytes
```

In [53]: `clean_data.to_csv('clean_data.csv')`

In [54]: `import os  
os.getcwd()`

Out[54]: 'C:\\\\Users\\\\rajendra\_damahe\\\\DATA SCIENCE'

In [55]: `clean_data`

Out[55]:

	Name	Domain	Age	Location	Salary	Exp
<b>0</b>	Mike	Datascience	34	Mumbai	5000	2
<b>1</b>	Teddy	Testing	45	Bangalore	10000	3
<b>2</b>	Umar	Dataanalyst	50.25	Bangalore	15000	4
<b>3</b>	Jane	Analytics	50.25	Hyderabad	20000	4.8
<b>4</b>	Uttam	Statistics	67	Bangalore	30000	5
<b>5</b>	Kim	NLP	55	Delhi	60000	10

## EDA TECHNIQUE LETS APPLY

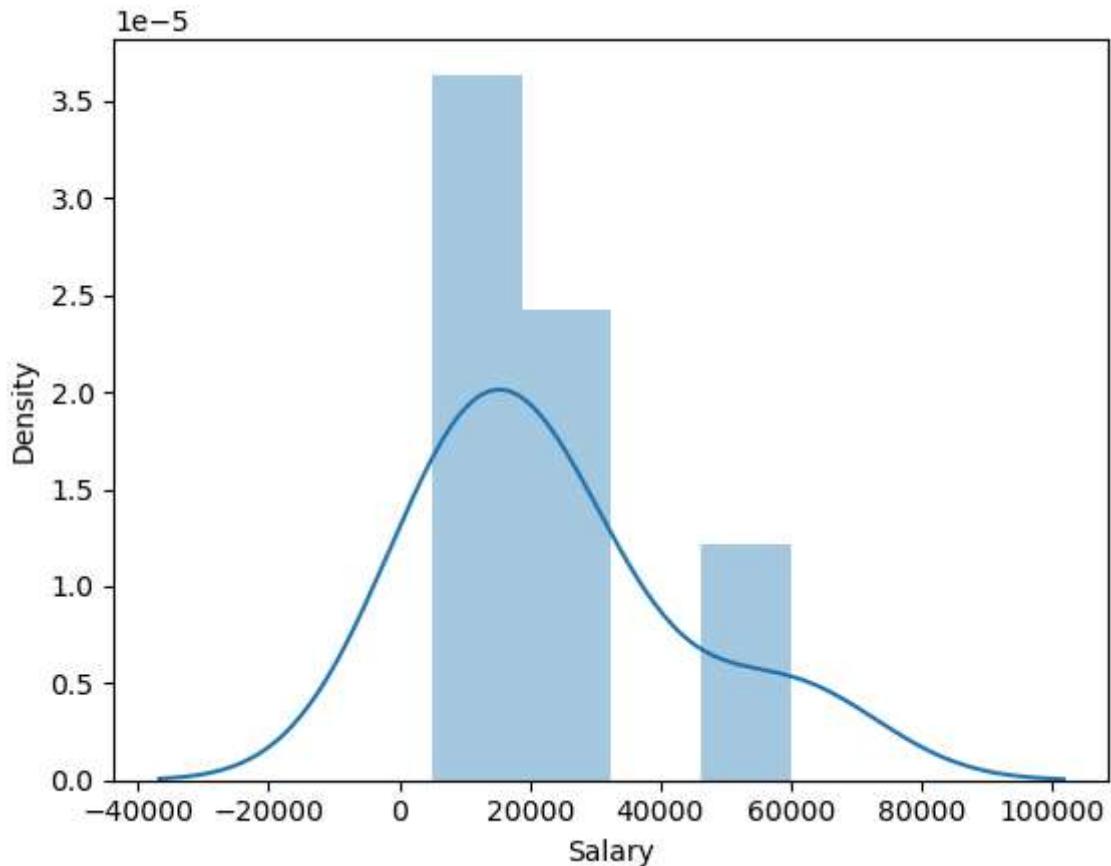
In [71]: `import matplotlib.pyplot as plt # visualization  
import seaborn as sns`

In [72]: `import warnings  
warnings.filterwarnings('ignore')`

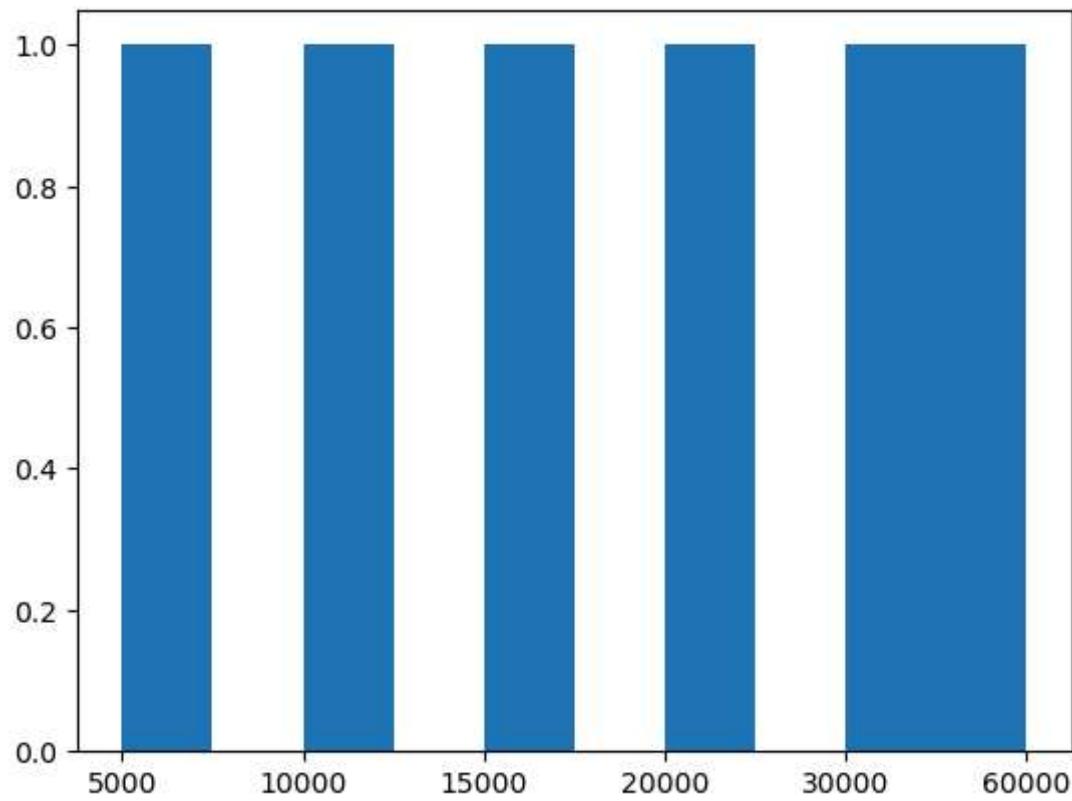
In [73]: `clean_data['Salary']`

```
Out[73]: 0    5000
         1   10000
         2   15000
         3   20000
         4   30000
         5   60000
Name: Salary, dtype: object
```

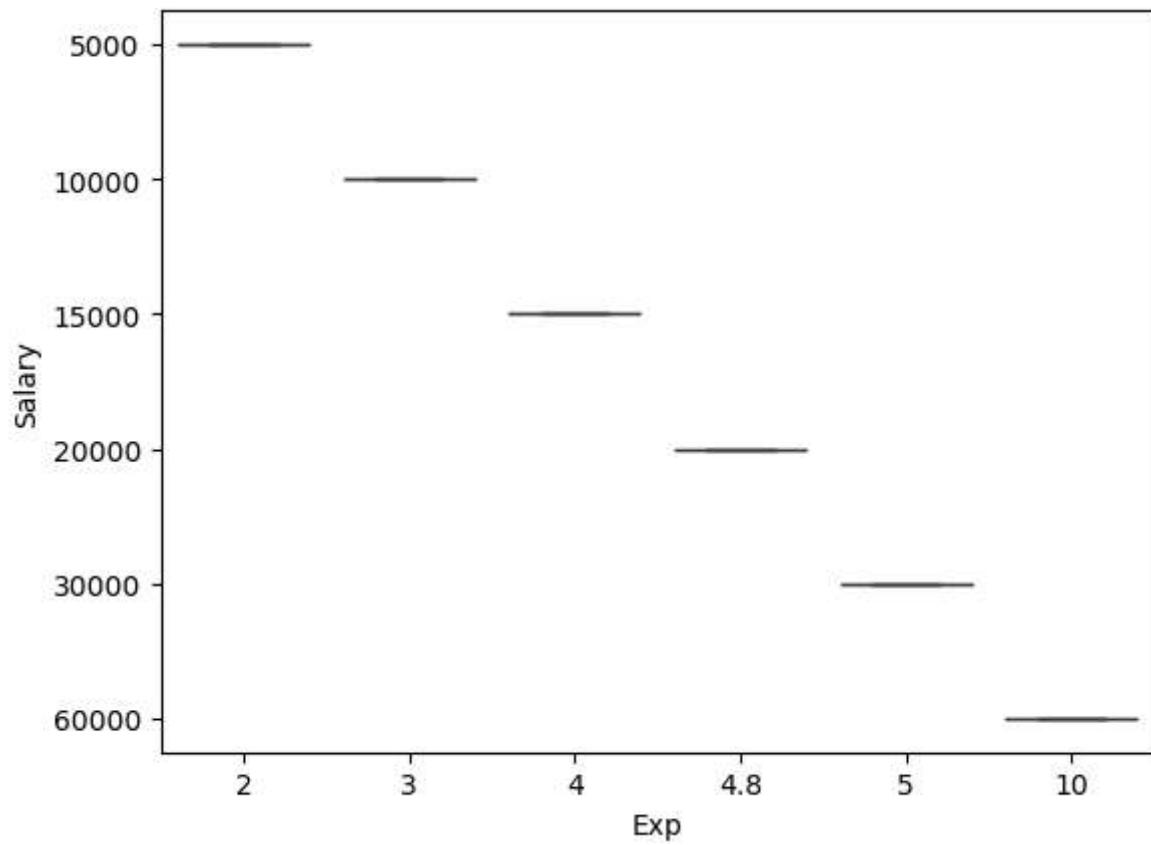
```
In [74]: vis1 = sns.distplot(clean_data['Salary'])
```



```
In [75]: vis2 = plt.hist(clean_data['Salary'])
```



```
In [56]: vis4 = sns.boxplot(data=clean_data,x = 'Exp', y='Salary')
```



```
In [57]: clean_data[0:6:2]
```

Out[57]:

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascience	34	Mumbai	5000	2
2	Umar	Dataanalyst	50.25	Bangalore	15000	4
4	Uttam	Statistics	67	Bangalore	30000	5

In [58]: `clean_data[0:6:3]`

Out[58]:

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascience	34	Mumbai	5000	2
3	Jane	Analytics	50.25	Hyderbad	20000	4.8

In [59]: `clean_data[::-1]`

Out[59]:

	Name	Domain	Age	Location	Salary	Exp
5	Kim	NLP	55	Delhi	60000	10
4	Uttam	Statistics	67	Bangalore	30000	5
3	Jane	Analytics	50.25	Hyderbad	20000	4.8
2	Umar	Dataanalyst	50.25	Bangalore	15000	4
1	Teddy	Testing	45	Bangalore	10000	3
0	Mike	Datascience	34	Mumbai	5000	2

In [61]: `clean_data[::-1]`

Out[61]:

	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 [63]: `clean_data[0:4:2]`

Out[63]:

	Name	Domain	Age	Location	Salary	Exp
<b>0</b>	Mike	Datascience	34	Mumbai	5000	2
<b>2</b>	Umar	Dataanalyst	50.25	Bangalore	15000	4

In [84]: `clean_data.columns`

Out[84]: `Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp', 'Age'].fillna(np.mean(pd.))', dtype='object')`

In [85]: `X_iv = clean_data[['Name', 'Domain', 'Age', 'Location', 'Exp']]`

In [86]: `X_iv`

Out[86]:

	Name	Domain	Age	Location	Exp
<b>0</b>	Mike	Datascience	34	Mumbai	2
<b>1</b>	Teddy	Testing	45	Bangalore	3
<b>2</b>	Umar	Dataanalyst	NaN	Bangalore	4
<b>3</b>	Jane	Analytics	NaN	Hyderbad	4.8
<b>4</b>	Uttam	Statistics	67	Bangalore	5
<b>5</b>	Kim	NLP	55	Delhi	10

In [87]: `y_dv = clean_data[['Salary']] = clean_data[['Salary']]`

In [88]: `y_dv`

Out[88]:

	Salary
<b>0</b>	5000
<b>1</b>	10000
<b>2</b>	15000
<b>3</b>	20000
<b>4</b>	30000
<b>5</b>	60000

In [89]: `emp`

Out[89]:

	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 [90]: `clean_data`

Out[90]:

	Name	Domain	Age	Location	Salary	Exp	Age].fillna(np.mean(pd.))
0	Mike	Datascience	34	Mumbai	5000	2	34
1	Teddy	Testing	45	Bangalore	10000	3	45
2	Umar	Dataanalyst	NaN	Bangalore	15000	4	50.25
3	Jane	Analytics	NaN	Hyderbad	20000	4.8	50.25
4	Uttam	Statistics	67	Bangalore	30000	5	67
5	Kim	NLP	55	Delhi	60000	10	55

In [91]: `x_iv`

Out[91]:

	Name	Domain	Age	Location	Exp
0	Mike	Datascience	34	Mumbai	2
1	Teddy	Testing	45	Bangalore	3
2	Umar	Dataanalyst	NaN	Bangalore	4
3	Jane	Analytics	NaN	Hyderbad	4.8
4	Uttam	Statistics	67	Bangalore	5
5	Kim	NLP	55	Delhi	10

In [92]: `y_dv`

Out[92]:

Salary	
<b>0</b>	5000
<b>1</b>	10000
<b>2</b>	15000
<b>3</b>	20000
<b>4</b>	30000
<b>5</b>	60000

In [93]:

clean\_data

Out[93]:

	Name	Domain	Age	Location	Salary	Exp	Age].fillna(np.mean(pd.))
<b>0</b>	Mike	Datascience	34	Mumbai	5000	2	34
<b>1</b>	Teddy	Testing	45	Bangalore	10000	3	45
<b>2</b>	Umar	Dataanalyst	NaN	Bangalore	15000	4	50.25
<b>3</b>	Jane	Analytics	NaN	Hyderbad	20000	4.8	50.25
<b>4</b>	Uttam	Statistics	67	Bangalore	30000	5	67
<b>5</b>	Kim	NLP	55	Delhi	60000	10	55

In [94]:

imputation = pd.get\_dummies(clean\_data)

In [95]:

imputation

Out[95]:

	Name_Jane	Name_Kim	Name_Mike	Name_Teddy	Name_Umar	Name_Uttam	Domain_
<b>0</b>	False	False	True	False	False	False	False
<b>1</b>	False	False	False	True	False	False	False
<b>2</b>	False	False	False	False	True	False	False
<b>3</b>	True	False	False	False	False	False	False
<b>4</b>	False	False	False	False	False	True	False
<b>5</b>	False	True	False	False	False	False	False

6 rows × 37 columns



In [96]:

clean\_data

Out[96]:

	Name	Domain	Age	Location	Salary	Exp	Age].fillna(np.mean(pd.))
0	Mike	Datascienc	34	Mumbai	5000	2	34
1	Teddy	Testing	45	Bangalore	10000	3	45
2	Umar	Dataanalyst	NaN	Bangalore	15000	4	50.25
3	Jane	Analytics	NaN	Hyderbad	20000	4.8	50.25
4	Uttam	Statistics	67	Bangalore	30000	5	67
5	Kim	NLP	55	Delhi	60000	10	55

In [97]: imputation

Out[97]:

	Name_Jane	Name_Kim	Name_Mike	Name_Teddy	Name_Umar	Name_Uttam	Domain_
0	False	False	True	False	False	False	False
1	False	False	False	True	False	False	False
2	False	False	False	False	True	False	False
3	True	False	False	False	False	False	False
4	False	False	False	False	False	True	False
5	False	True	False	False	False	False	False

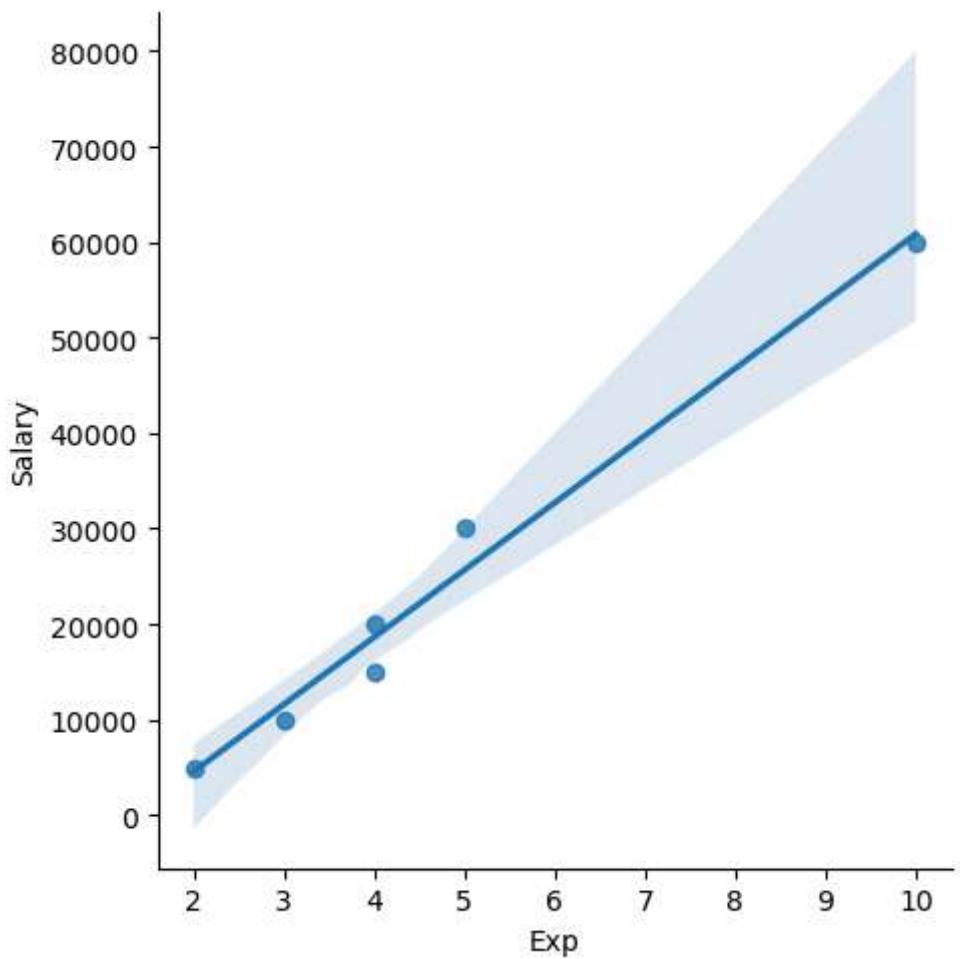
6 rows × 37 columns



In [64]:

```
clean_data['Salary'] = clean_data['Salary'].round().astype(int)
clean_data['Exp'] = clean_data['Exp'].round().astype(int)
vis4 = sns.lmplot(data = clean_data, x='Exp', y ='Salary')
vis4
```

Out[64]: <seaborn.axisgrid.FacetGrid at 0x23476456510>



In [ ]: