

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

```
In [2]: import warnings
warnings.filterwarnings('ignore')
```

```
In [3]: pd.__version__
```

```
Out[3]: '2.2.2'
```

```
In [4]: emp=pd.read_excel(r'C:\Users\HP\Downloads\Rawdata.xlsx')
```

```
In [5]: 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	Hyderabad	2000^0	NaN
4	Uttam*	Statistics	67-yr	NaN	30000-	5+ year
5	Kim	NLP	55yr	Delhi	6000^\$0	10+

```
In [6]: id(emp)
```

```
Out[6]: 2398234323056
```

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

```
Out[7]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')
```

```
In [8]: emp.shape
```

```
Out[8]: (6, 6)
```

```
In [9]: emp.head
```

	Name	Domain	Age	Location	Sal	ary	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	Hyderabad	2000^0		NaN
4	Uttam*	Statistics	67-yr	NaN	30000-		5+ year
5	Kim	NLP	55yr	Delhi	6000^\$0		10+>

```
In [10]: emp.tail
```

```
Out[10]: <bound method NDFrame.tail of      Name          Domain        Age  Location   Sal
           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 [11]: `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 [12]: `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 [13]: `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 [14]: `emp`

Out[14]:

	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	Hyderabad	2000^0	NaN
4	Uttam*	Statistics	67-yr	NaN	30000-	5+ year
5	Kim	NLP	55yr	Delhi	6000^\$0	10+

In [15]:

`emp.isnull()`

Out[15]:

	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 [16]:

`emp.isna()`

Out[16]:

	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 [17]:

`emp.isnull().sum()`

Out[17]:

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

```
In [18]: emp.columns
```

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

```
In [19]: emp
```

```
Out[19]:
```

	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	Hyderabad	2000^0	NaN
4	Uttam*	Statistics	67-yr	NaN	30000-	5+ year
5	Kim	NLP	55yr	Delhi	6000^\$0	10+

## Data Cleaning or Data Cleansing

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

```
Out[21]: 0      Mike
1      Teddy^
2      Uma#r
3      Jane
4      Uttam*
5      Kim
Name: Name, dtype: object
```

```
In [22]: emp['Name']=emp['Name'].str.replace(r'\W',' ',regex=True) # W-Non word Character
```

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

```
Out[23]: 0      Mike
1      Teddy
2      Umar
3      Jane
4      Uttam
5      Kim
Name: Name, dtype: object
```

```
In [24]: emp['Domain']=emp['Domain'].str.replace(r'\W',' ',regex=True)
```

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

```
Out[25]: 0    Datascienc
          1    Testing
          2  Dataanalyst
          3   Analytics
          4  Statistics
          5      NLP
Name: Domain, dtype: object
```

```
In [26]: 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 [27]: emp['Age']=emp['Age'].str.replace(r'\W',' ',regex=True)
emp['Age']
```

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

```
In [28]: emp['Age']=emp['Age'].str.extract('(\d+)')  # r(r'(\d+')
emp['Age']
```

```
Out[28]: 0    34
          1    45
          2    NaN
          3    NaN
          4    67
          5    55
Name: Age, dtype: object
```

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

```
Out[29]: 0      Mumbai
          1    Bangalore
          2      NaN
          3   Hyderabad
          4      NaN
          5      Delhi
Name: Location, dtype: object
```

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

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

```
In [31]: emp['Exp']=emp['Exp'].str.replace(r'\W',' ',regex=True)
emp['Exp']
```

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

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

```
Out[32]: 0      2
         1      3
         2      4
         3     NaN
         4      5
         5     10
Name: Exp, dtype: object
```

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

	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 [34]: clean_data=emp.copy()
```

```
In [35]: clean_data
```

Out[35]:

	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

## EDA TECHNIQUE LET'S APPLY

- Missing value treatment

In [37]: clean\_data

Out[37]:

	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 [38]: clean\_data.isnull().sum()

Out[38]:

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

In [39]: clean\_data['Age']

Out[39]:

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

```
In [40]: import numpy as np
```

```
In [41]: clean_data['Age']=clean_data['Age'].fillna(np.mean(pd.to_numeric(clean_data['Age'])))
```

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

```
Out[42]: 0      34
         1      45
         2    50.25
         3    50.25
         4      67
         5      55
Name: Age, dtype: object
```

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

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

```
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
```

	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	NaN	15000	4
<b>3</b>	Jane	Analytics	50.25	Hyderabad	20000	4.8
<b>4</b>	Uttam	Statistics	67	NaN	30000	5
<b>5</b>	Kim	NLP	55	Delhi	60000	10

```
In [47]: clean_data['Location'].isnull().sum()
```

Out[47]: 2

In [48]: `clean_data['Location']=clean_data['Location'].fillna(clean_data['Location'].mode()[0])`In [49]: `clean_data['Location']`

Out[49]:

0	Mumbai
1	Bangalore
2	Bangalore
3	Hyderabad
4	Bangalore
5	Delhi

Name: Location, dtype: object

In [50]: `clean_data`

Out[50]:

	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	Hyderabad	20000	4.8
4	Uttam	Statistics	67	Bangalore	30000	5
5	Kim	NLP	55	Delhi	60000	10

In [51]: `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 [52]: `clean_data['Age']=clean_data['Age'].astype(int)`  
`clean_data['Age']`

```
Out[52]: 0    34
         1    45
         2    50
         3    50
         4    67
         5    55
Name: Age, dtype: int32
```

```
In [53]: 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 [54]: clean_data['Salary']=clean_data['Salary'].astype(int)
clean_data['Salary']
```

```
Out[54]: 0    5000
         1   10000
         2   15000
         3   20000
         4   30000
         5   60000
Name: Salary, dtype: int32
```

```
In [55]: clean_data['Exp']=clean_data['Exp'].astype(int)
clean_data['Exp']
```

```
Out[55]: 0    2
         1    3
         2    4
         3    4
         4    5
         5   10
Name: Exp, dtype: int32
```

```
In [56]: 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 [57]: 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 [58]: 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 [59]: clean_data
```

```
Out[59]:
```

	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	Bangalore	15000	4
<b>3</b>	Jane	Analytics	50	Hyderabad	20000	4
<b>4</b>	Uttam	Statistics	67	Bangalore	30000	5
<b>5</b>	Kim	NLP	55	Delhi	60000	10

```
In [60]: clean_data.to_csv('clean_data.csv')
```

```
In [61]: import os
os.getcwd()
```

```
Out[61]: 'C:\\Users\\HP'
```

```
In [62]: clean_data
```

```
Out[62]:
```

	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

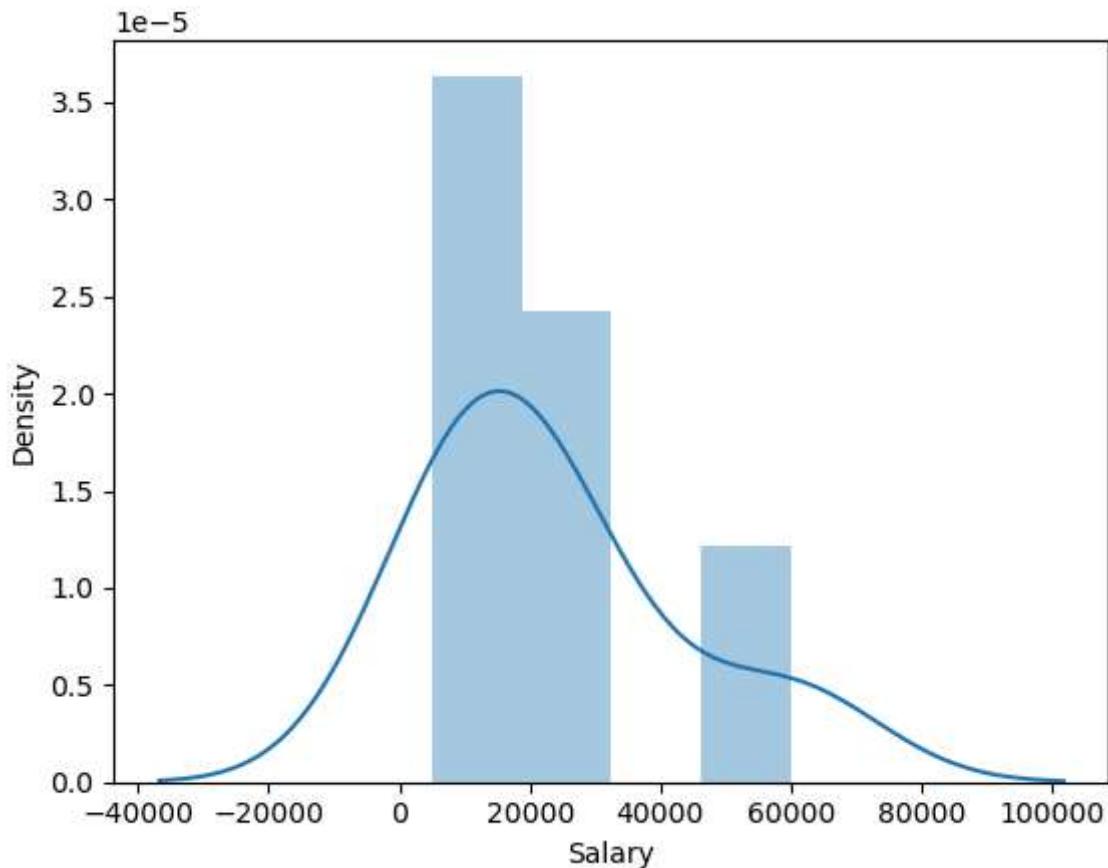
```
In [63]: import matplotlib.pyplot as plt  
import seaborn as sns
```

```
In [64]: import warnings  
warnings.filterwarnings('ignore')
```

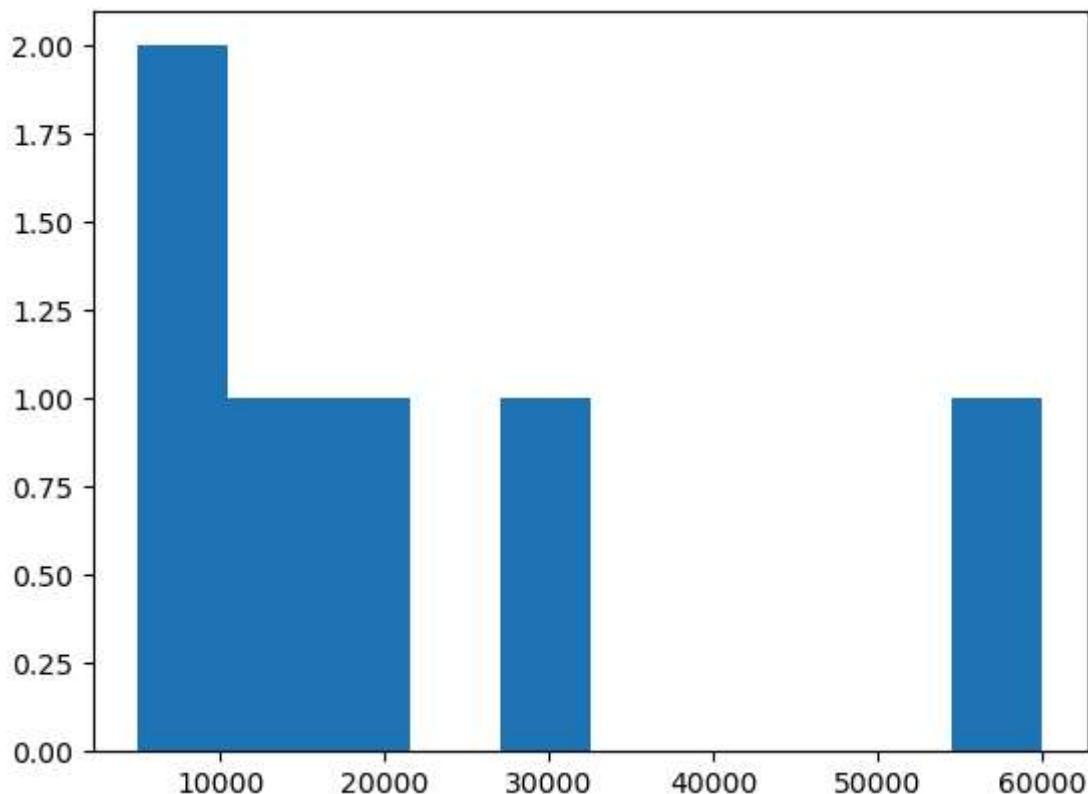
```
In [65]: clean_data['Salary']
```

```
Out[65]: 0      5000  
1     10000  
2     15000  
3     20000  
4     30000  
5     60000  
Name: Salary, dtype: int32
```

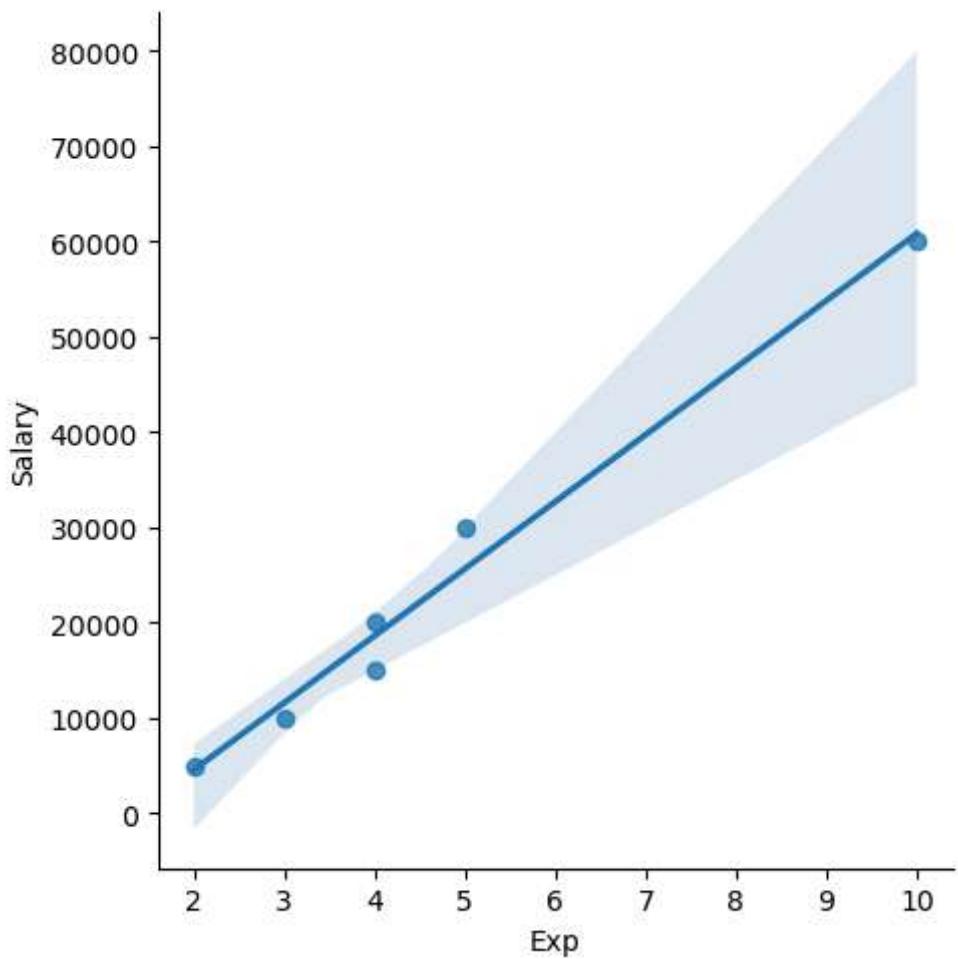
```
In [66]: vis1 = sns.distplot(clean_data['Salary'])
```



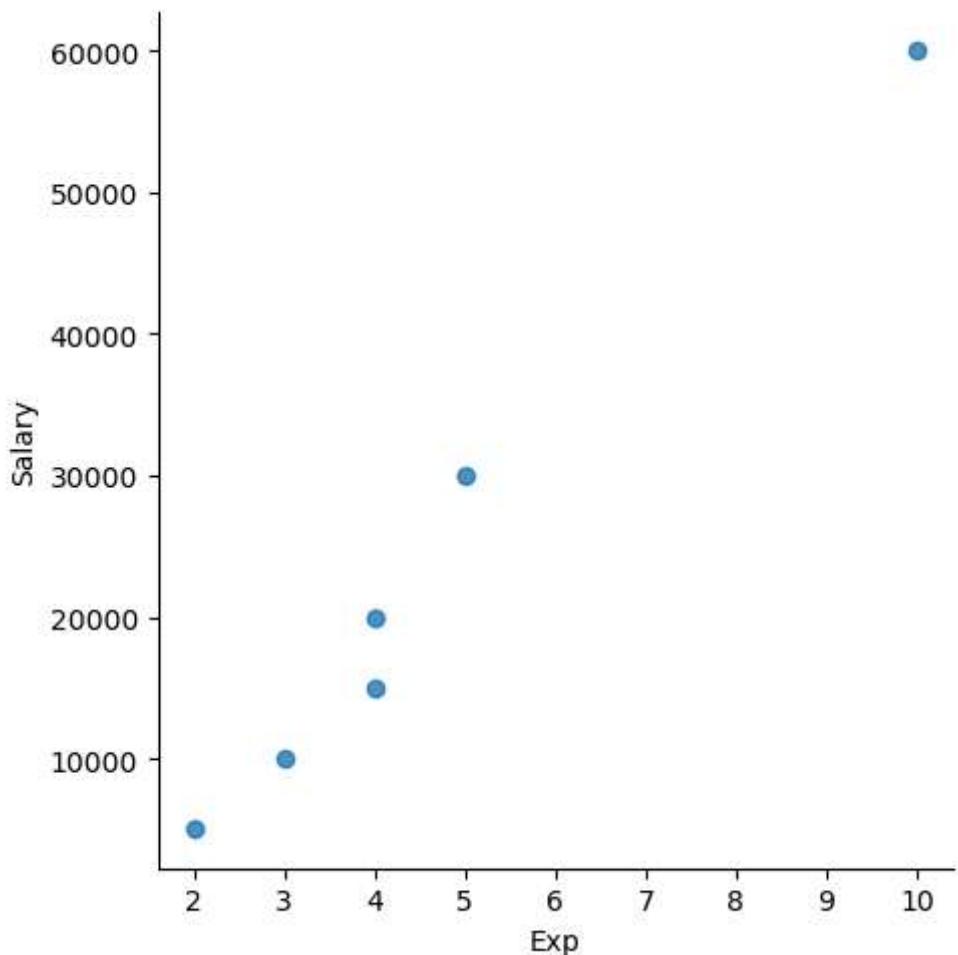
```
In [67]: vis2 = plt.hist(clean_data['Salary'])
```



```
In [68]: vis4 = sns.lmplot(data=clean_data,x = 'Exp', y='Salary')
```



```
In [69]: vis5 = sns.lmplot(data=clean_data,x = 'Exp', y='Salary', fit_reg = False)
```



```
In [70]: clean_data[:]
```

```
Out[70]:
```

	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 [71]: clean_data[0:6:2]
```

```
Out[71]:
```

	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 [72]: `clean_data[:::-1]`

Out[72]:

	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 [73]: `clean_data.columns`

Out[73]: `Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')`

In [74]: `X_iv = clean_data[['Name', 'Domain', 'Age', 'Location', 'Exp']]`

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

In [76]: `y_dv = clean_data[['Salary']]  
y_dv`

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 [79]: 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 [80]: y\_dv

Out[80]:

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

In [81]:

clean\_data

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	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 [82]:

```
imputation = pd.get_dummies(clean_data)
imputation
```

Out[82]:

	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 [83]:

clean\_data

Out[83]:

	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 [84]:

imputation

Out[84]:

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



In [ ]:

In [ ]:

In [ ]: