

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

```
In [3]: emp=pd.read_excel(r'C:\Users\arsha_4tjdyqj\Downloads\Rawdata.xlsx')
emp
```

Out[3]:

	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 [5]: emp.shape
```

Out[5]: (6, 6)

```
In [7]: emp.head()
```

Out[7]:

	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

```
In [9]: emp.head(5)
```

Out[9]:

	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

```
In [11]: emp.tail()
```

Out[11]:

	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.isnull().sum()`

Out[13]:

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

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

DATA CLEANING OR DATA CLEANSING

In [17]:

```
emp['Name']=emp['Name'].str.replace(r'^[a-zA-Z0-9\s]+$', '', regex=True)
emp['Name']
```

Out[17]:

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

```
emp['Domain']=emp['Domain'].str.replace(r'^[a-zA-Z0-9\s]+$', '', regex=True)
emp['Domain']
```

Out[23]:

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

In [27]:

```
emp['Age']=emp['Age'].str.replace(r'^[0-9]+$', '', regex=True)
emp['Age']
```

Out[27]:

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

In [29]:

```
emp
```

Out[29]:

	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 [33]:

```
emp['Exp']=emp['Exp'].str.replace(r'^[0-9]+$', '', regex=True)
emp['Exp']
```

```
Out[33]: 0      2
         1      3
         2      4
         3    NaN
         4      5
         5     10
Name: Exp, dtype: object
```

```
In [35]: emp['Salary']=emp['Salary'].str.replace(r'^[a-zA-Z0-9\s]', '', regex=True)
emp['Salary']
```

```
Out[35]: 0      5000
         1     10000
         2    15000
         3   20000
         4   30000
         5   60000
Name: Salary, dtype: object
```

```
In [37]: 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 [39]: clean_data=emp.copy()
```

MISSING VALUES TREATMENT FOR NUMERICAL DATA

```
In [41]: 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 [43]: clean_data['Age']
```

```
Out[43]: 0    34
1    45
2    NaN
3    NaN
4    67
5    55
Name: Age, dtype: object
```

```
In [45]: import numpy as np
```

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

```
In [49]: clean_data['Age']
```

```
Out[49]: 0    34
1    45
2    50.25
3    50.25
4    67
5    55
Name: Age, dtype: object
```

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

```
Out[53]: 0    2
1    3
2    4
3    4.8
4    5
5    10
Name: Exp, dtype: object
```

```
In [55]: 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	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 [57]: clean_data['Location'].isnull().sum()
```

```
Out[57]: 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 [63]: clean_data['Location']=clean_data['Location'].fillna(clean_data['Location'].mode()[0])
clean_data['Location']
```

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

```
In [65]: 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	Hyderbad	20000	4.8
4	Uttam	Statistics	67	Bangalore	30000	5
5	Kim	NLP	55	Delhi	60000	10

```
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      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 [73]: clean_data['Age']=clean_data['Age'].astype(int)
clean_data['Salary']=clean_data['Salary'].astype(int)
clean_data['Exp']=clean_data['Exp'].astype(int)
```

```
In [77]: 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 [79]: 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 [81]: 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 [83]: 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	Hyderabad	20000	4
4	Uttam	Statistics	67	Bangalore	30000	5
5	Kim	NLP	55	Delhi	60000	10

```
In [85]: clean_data.to_csv('clean_data.csv')
```

```
In [87]: import os  
os.getcwd()
```

```
Out[87]: 'C:\\Users\\arsha_4tjdyqj'
```

EDA TECHNIQUES LETS APPLY

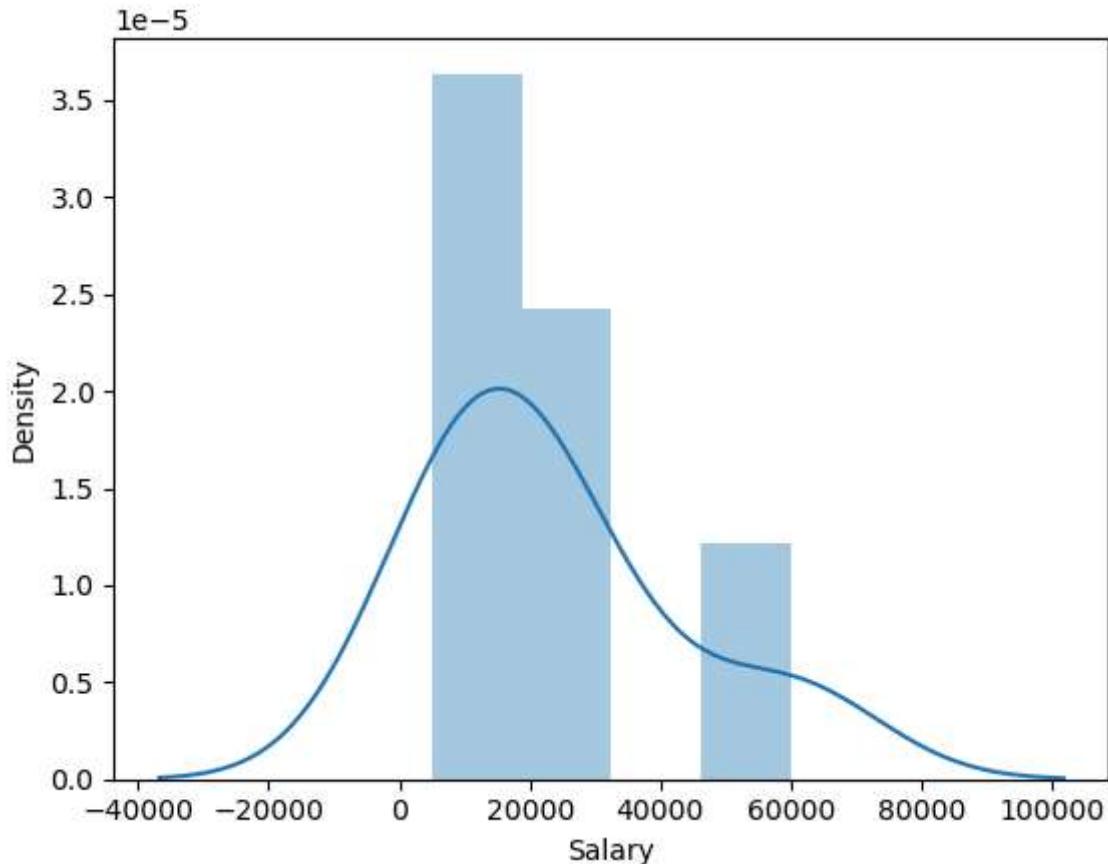
```
In [89]: import matplotlib.pyplot as plt  
import seaborn as sns
```

```
In [93]: import warnings  
warnings.filterwarnings('ignore')
```

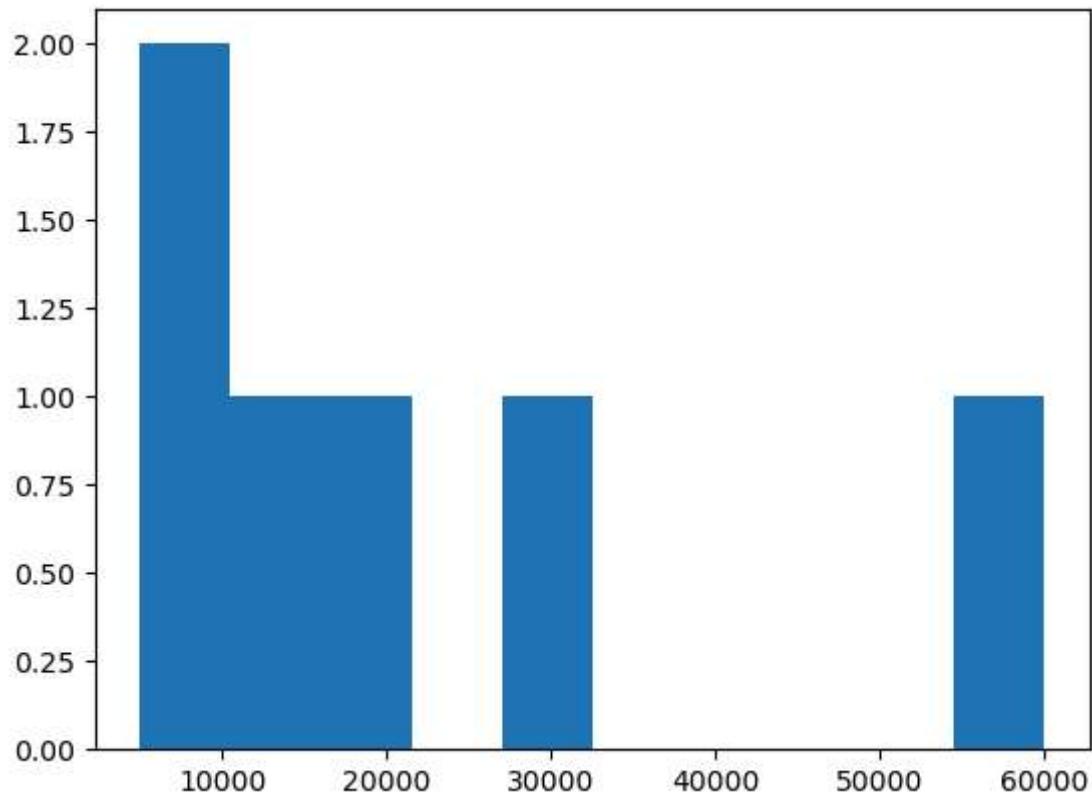
```
In [95]: clean_data['Salary']
```

```
Out[95]: 0      5000  
1     10000  
2     15000  
3     20000  
4     30000  
5     60000  
Name: Salary, dtype: int32
```

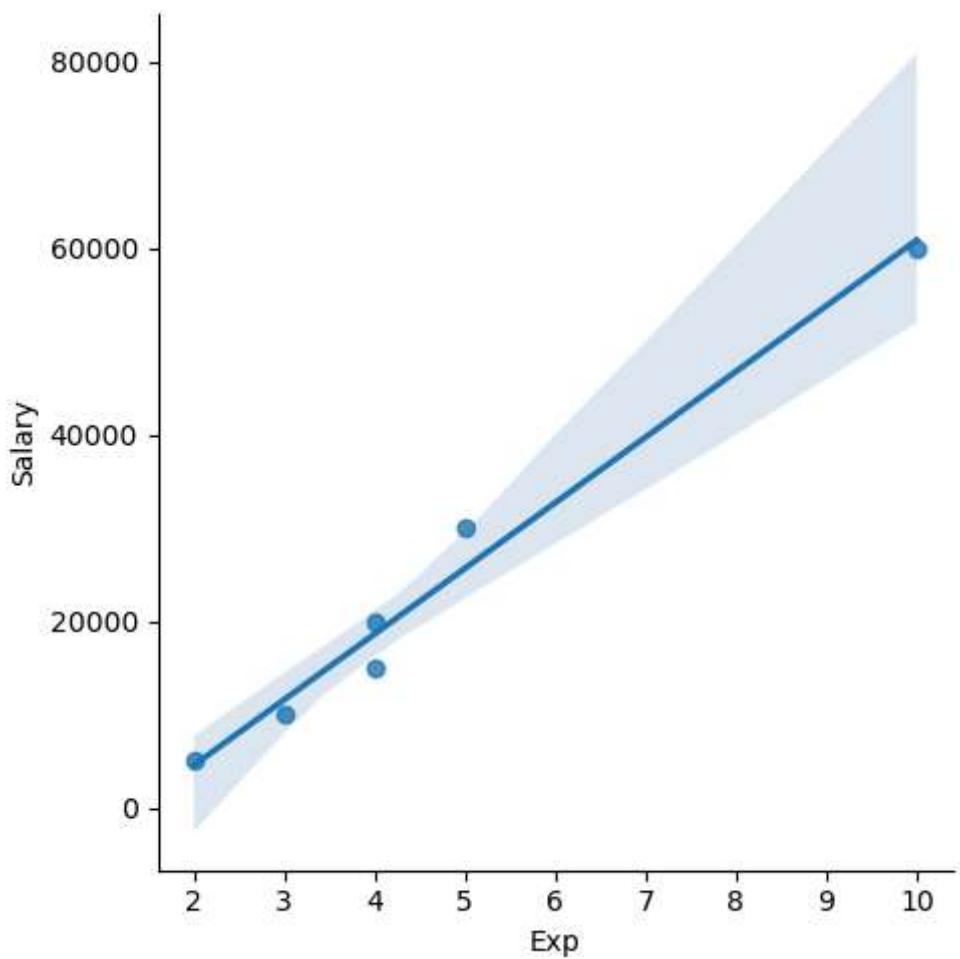
```
In [97]: vis1=sns.distplot(clean_data['Salary'])
```



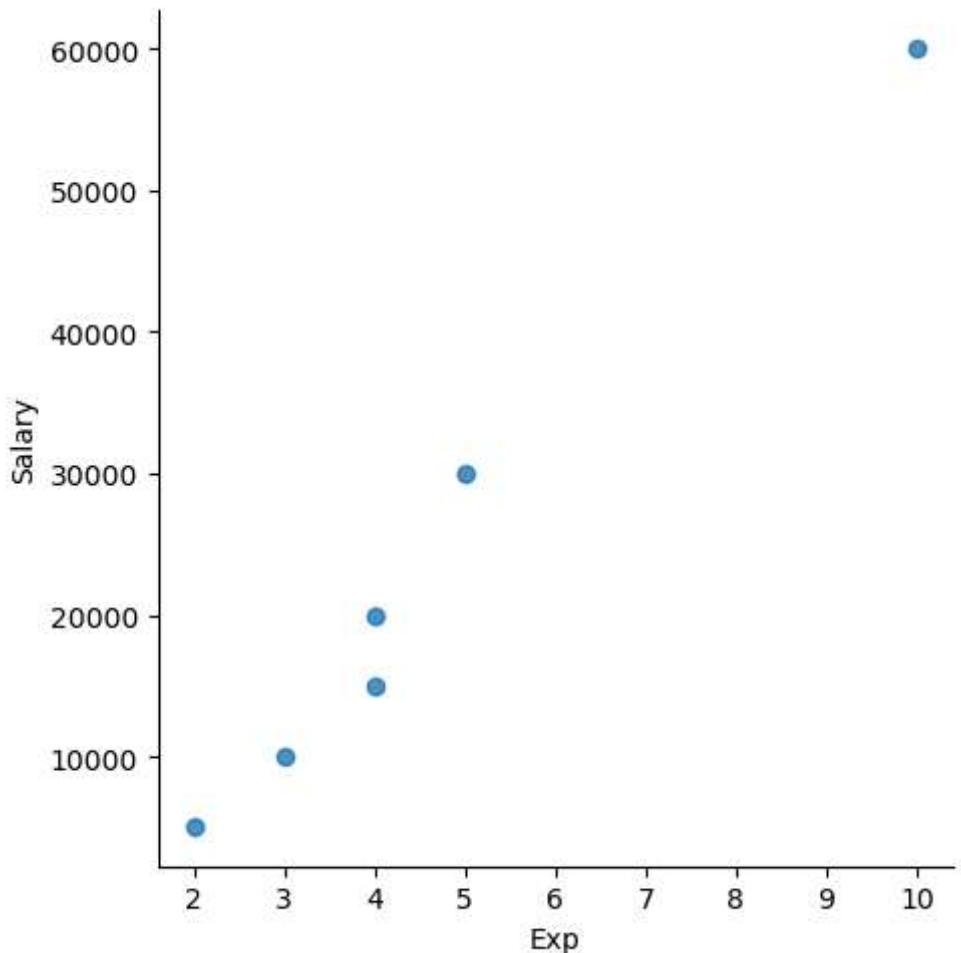
```
In [99]: vis2=plt.hist(clean_data['Salary'])#univariate
```



```
In [107...]: vis3=sns.lmplot(data=clean_data,x='Exp',y='Salary')#bivariate analysis
```



```
In [109]: vis4 = sns.lmplot(data=clean_data,x = 'Exp', y='Salary', fit_reg = False)
```



```
In [111]: clean_data[:]
```

```
Out[111]:
```

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

```
Out[113]:
```

	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 [115... clean_data[:::-1]

	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 [117... 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	Hyderabad	20000	4
4	Uttam	Statistics	67	Bangalore	30000	5
5	Kim	NLP	55	Delhi	60000	10

In [119... clean_data.columns

Out[119... Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')

In [121... x_indvar=clean_data[['Name', 'Domain', 'Age', 'Location', 'Exp']]

In [123... x_indvar

	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 [129...]  
y_depvar=clean_data[['Salary']]  
y_depvar
```

```
Out[129...]  
Salary
```

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

```
In [131...]  
y_depvar
```

```
Out[131...]  
Salary
```

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

```
In [134...]  
emp
```

```
Out[134...]  
Name    Domain    Age    Location    Salary    Exp
```

	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 [136...]  
clean_data
```

Out[136...]

	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 [138...]

x_indvar

Out[138...]

	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 [140...]

y_depvar

Out[140...]

Salary

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

In [142...]

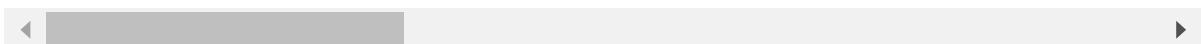
imputation=pd.get_dummies(clean_data)

In [144...]

imputation

Out[144...]

	Age	Salary	Exp	Name_Jane	Name_Kim	Name_Mike	Name_Teddy	Name_Umar	Nan
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 [146...]

clean_data

Out[146...]

	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 [148...]

imputation.shape

Out[148...]

(6, 19)