

```
In [1]: import numpy as np
```

```
In [3]: import pandas as pd
```

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

	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]: Raw.columns
```

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

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

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

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

	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 [8]: Raw.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 [9]: `Raw.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 [10]: `Raw['Domain']`

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

In [11]: `Raw.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]: `Raw.isnull().sum()`

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

In [13]: `Raw['Name']`

```
Out[13]: 0      Mike
          1      Teddy^
          2      Uma#r
          3      Jane
          4      Uttam*
          5      Kim
Name: Name, dtype: object
```

DATA CLEANING OR CLEANSING

```
In [15]: Raw['Name'] = Raw['Name'].str.replace(r'\W', ' ', regex= True)
Raw['Name']
```

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

```
In [16]: Raw['Domain']
```

```
Out[16]: 0      Datascience#$#
          1      Testing
          2      Dataanalyst^^#
          3      Ana^^lytics
          4      Statistics
          5      NLP
Name: Domain, dtype: object
```

```
In [17]: Raw['Domain'] = Raw['Domain'].str.replace(r'\W', ' ', regex= True)
Raw['Domain']
```

```
Out[17]: 0      Datascience
          1      Testing
          2      Dataanalyst
          3      Analytics
          4      Statistics
          5      NLP
Name: Domain, dtype: object
```

```
In [18]: Raw['Location']
```

```
Out[18]: 0      Mumbai
          1      Bangalore
          2      NaN
          3      Hyderabad
          4      NaN
          5      Delhi
Name: Location, dtype: object
```

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

```
Out[19]: 0      Mumbai
          1      Bangalore
          2      NaN
          3      Hyderabad
          4      NaN
          5      Delhi
Name: Location, dtype: object
```

```
In [20]: Raw['Age']=Raw['Age'].str.extract('(\d+)')
```

```
<>:1: SyntaxWarning: invalid escape sequence '\d'
<>:1: SyntaxWarning: invalid escape sequence '\d'
C:\Users\mdtan\AppData\Local\Temp\ipykernel_11920\2585282648.py:1: SyntaxWarning:
invalid escape sequence '\d'
Raw['Age']=Raw['Age'].str.extract('(\d+)')
```

```
In [21]: Raw['Age']
```

```
Out[21]: 0      34
          1      45
          2      NaN
          3      NaN
          4      67
          5      55
Name: Age, dtype: object
```

```
In [22]: Raw['Salary']
```

```
Out[22]: 0      5^00#0
          1      10%%000
          2      1$5%000
          3      2000^0
          4      30000-
          5      6000^$0
Name: Salary, dtype: object
```

```
In [23]: Raw['Salary']=Raw['Salary'].str.replace(r'\W',' ',regex = True)
Raw['Salary']
```

```
Out[23]: 0      5000
          1      10000
          2      15000
          3      20000
          4      30000
          5      60000
Name: Salary, dtype: object
```

```
In [24]: Raw
```

Out[24]:

	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 [46]: `Raw['Exp']`

Out[46]:

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

In [48]: `Raw['Exp'] = Raw['Exp'].str.extract('(\d+)')`

```
<>:1: SyntaxWarning: invalid escape sequence '\d'
<>:1: SyntaxWarning: invalid escape sequence '\d'
C:\Users\mdtan\AppData\Local\Temp\ipykernel_11920\3954229274.py:1: SyntaxWarning:
invalid escape sequence '\d'
Raw['Exp'] = Raw['Exp'].str.extract('(\d+)')
```

In [49]: `Raw['Exp']`

Out[49]:

```
0      2
1      3
2      4
3      NaN
4      5
5     10
Name: Exp, dtype: object
```

In [51]: `Raw`

Out[51]:

	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 [52]: `clean_data=Raw.copy()`

In [54]: `clean_data`

Out[54]:

	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 [55]: `clean_data['Age']`

Out[55]:

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

In [56]: `import numpy as np`

In [58]: `clean_data['Age']= clean_data['Age'].fillna(np.mean(np.mean(pd.to_numeric(clean_data['Age']))))`

Out[58]:

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

In [59]: `clean_data['Exp']`

Out[59]:

```
0    2
1    3
2    4
3    NaN
4    5
5    10
Name: Exp, dtype: object
```

In [61]: `clean_data['Exp']=clean_data['Exp'].fillna(np.mean(pd.to_numeric(clean_data['Exp'])))`

Out[61]:

```
0    2
1    3
2    4
3    4.8
4    5
5    10
Name: Exp, dtype: object
```

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	50.25	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 [63]: `clean_data['Location'].isnull().sum()`

Out[63]: 2

In [72]: `clean_data['Location']`

Out[72]:

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

Name: Location, dtype: object

In [74]: `clean_data['Location'] = clean_data['Location'].fillna(clean_data['Location'].mode()[0])`
clean_data

Out[74]:

	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 [75]: `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 [76]: `clean_data['Age'] = clean_data['Age'].astype(int)`

In [78]: `clean_data`

Out[78]:

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

In [80]: `clean_data['Salary'] = clean_data['Salary'].astype(int)`
`clean_data['Exp'] = clean_data['Exp'].astype(int)`

In [82]: `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 [83]: `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 [84]: `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 [85]: clean_data

Out[85]:

	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 [99]:

```
clean_data.to_csv('clean_data.csv')
import os
os.getcwd()
```

Out[99]:

```
'C:\\\\Users\\\\mdtan'
```

In [101...]:

clean_data

Out[101...]:

	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

EDA TECHNIQUE

In [104...]:

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

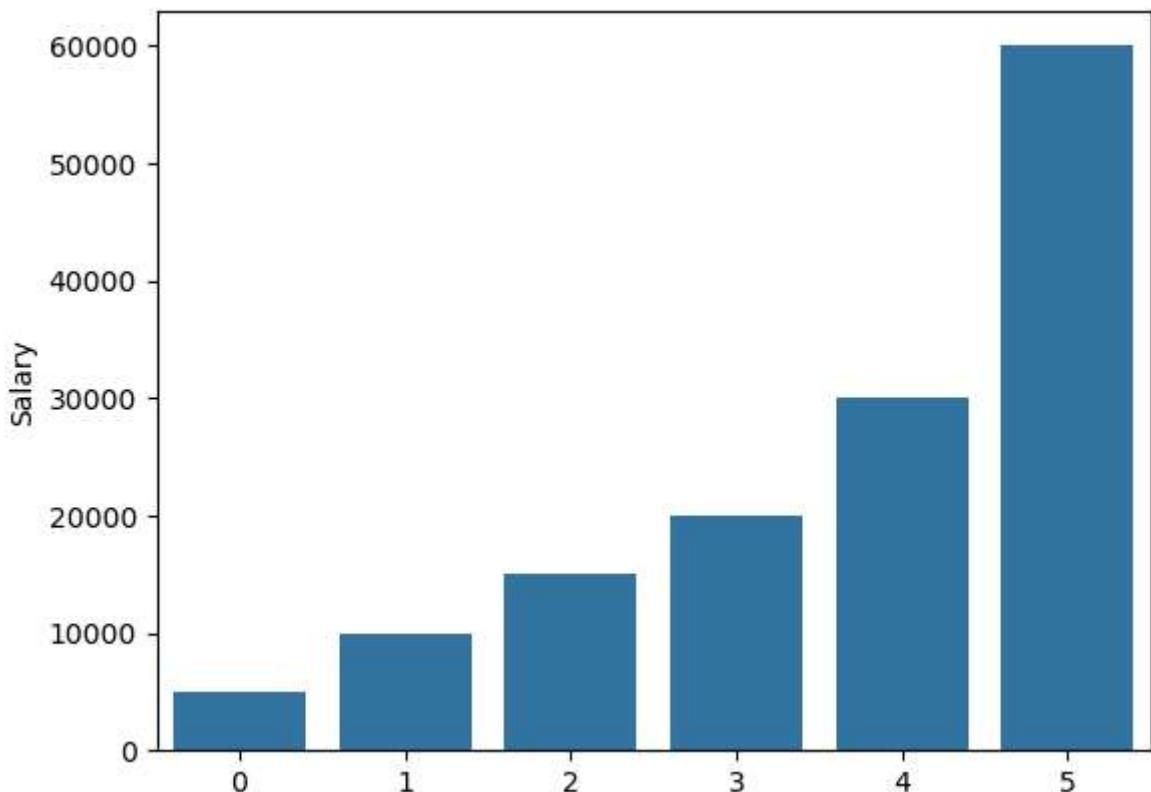
In [106...]:

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

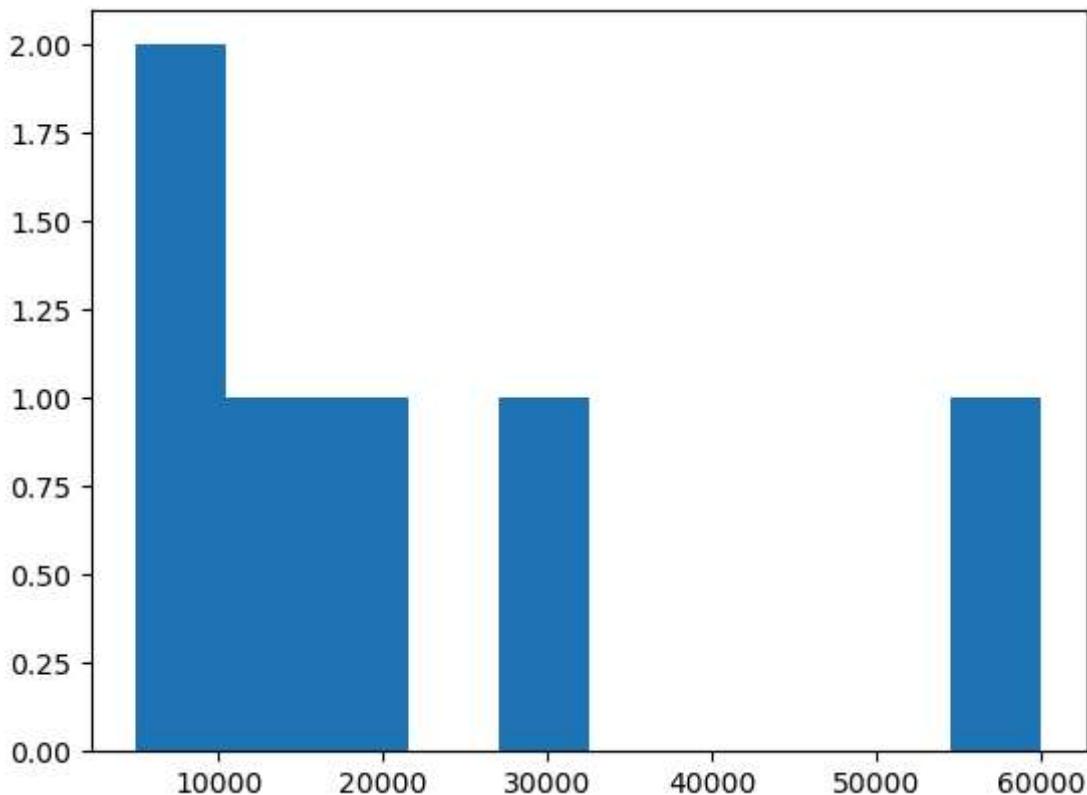
```
In [108...]: clean_data['Salary']
```

```
Out[108...]: 0      5000
  1     10000
  2    15000
  3   20000
  4   30000
  5   60000
Name: Salary, dtype: int32
```

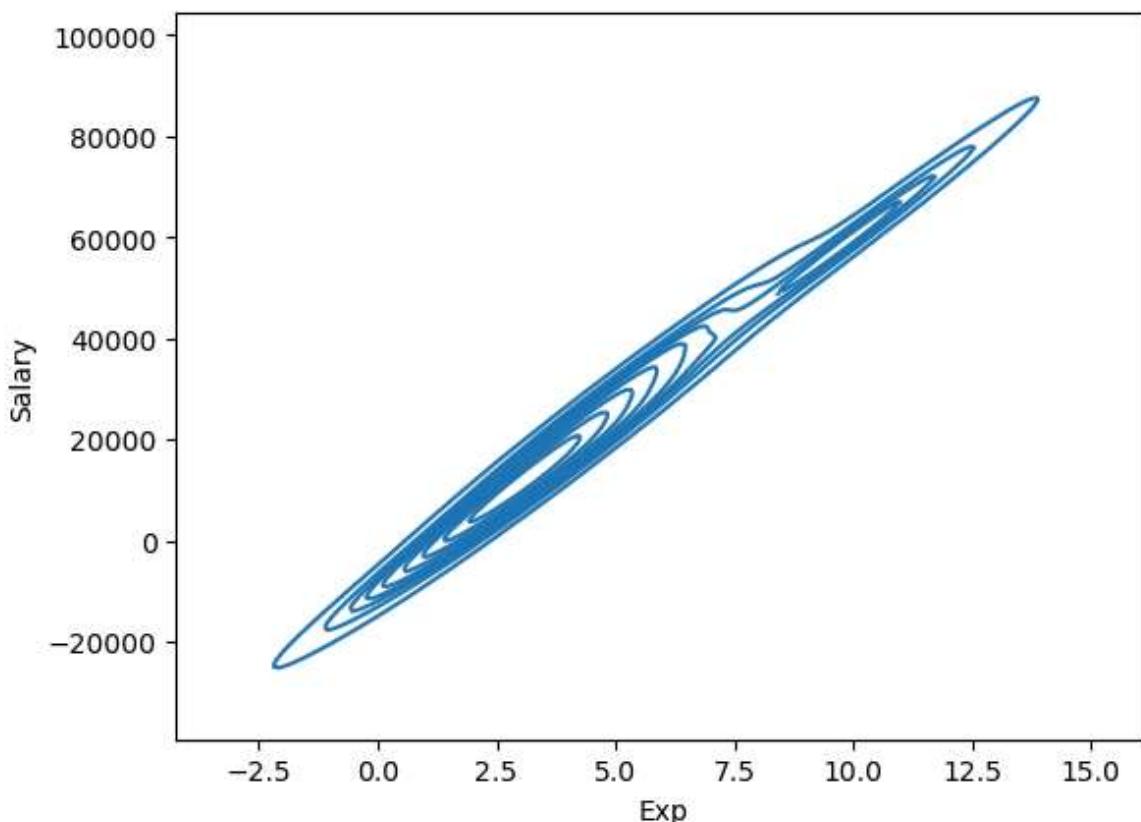
```
In [132...]: dell=sns.barplot(clean_data['Salary'])
```



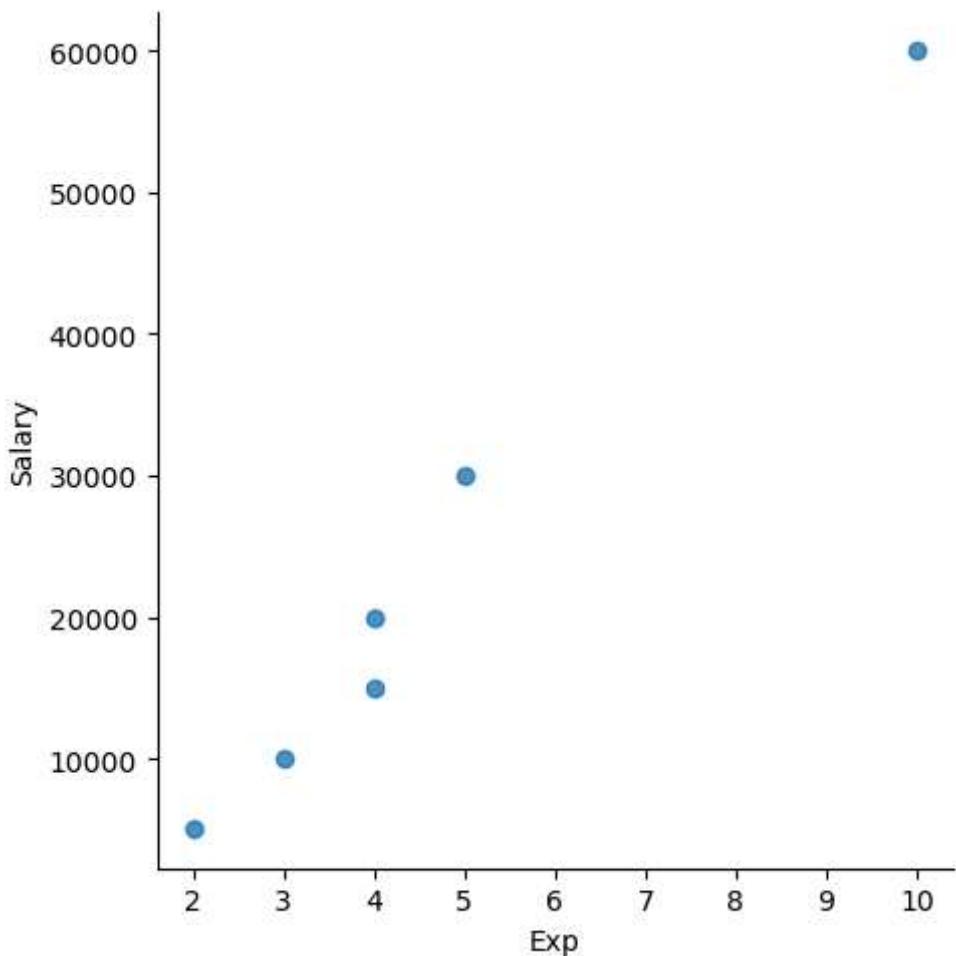
```
In [140...]: dell= plt.hist(clean_data['Salary'])
```



```
In [146...]: dell = sns.kdeplot(data=clean_data, x = 'Exp', y='Salary')
```



```
In [148...]: dell = sns.lmplot(data=clean_data, x = 'Exp', y='Salary', fit_reg = False)
```



In [152...]

clean_data[:]

Out[152...]

	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 [156...]

clean_data[0:6]

Out[156...]

	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 [164...]

clean_data[0:4:]

Out[164...]

	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

In [166...]

clean_data[::-1]

Out[166...]

	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 [168...]

clean_data.columns

Out[168...]

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

In [170...]

hp = clean_data[['Name', 'Domain', 'Age', 'Location', 'Exp']]

In [172...]

hp

Out[172...]

	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 [174...]

BMW = clean_data[['Salary']]

In [176...]

BMW

Out[176...]

Salary

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

In [178...]

Raw

Out[178...]

	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 [180...]

clean_data

Out[180...]

	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 [182...]

hp

Out[182...]

	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 [184...]

BMW

Out[184...]

Salary

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

In [186...]

clean_data

Out[186...]

	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 [188...]

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

Out[188...]

	Age	Salary	Exp	Name_Jane	Name_Kim	Name_Mike	Name_Teddy	Name_Umar
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 [190...]

```
clean_data
```

Out[190...]

	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 [194...]

```
imputation
```

Out[194...]

	Age	Salary	Exp	Name_Jane	Name_Kim	Name_Mike	Name_Teddy	Name_Umar
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 []: