

Raw data to clean data conversion using python EDA

```
In [7]: import pandas as pd
```

```
In [8]: emp = pd.read_excel(r"E:\naresh_it\20feb2025\19th - EDA Practicle\19th - EDA Practi
```

```
In [9]: 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 [10]: emp.columns
```

```
Out[10]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')
```

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

```
Out[16]: (6, 6)
```

```
In [18]: 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 [20]: emp.tail()
```

Out[20]:

	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 [22]:

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 [24]:

emp

Out[24]:

	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 [26]:

emp['Domain']

Out[26]:

```

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

```

In [28]:

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 [30]: `emp.isnull().sum()`

Out[30]:

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

In [32]: `emp['Name']`

Out[32]:

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

In [50]: `emp['Name'] = emp['Name'].str.replace(r'\W', ' ', regex=True)`

In [52]: `emp['Name']`

Out[52]:

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

In [54]: `emp['Domain']`

Out[54]:

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

```
In [56]: emp['Domain'] = emp['Domain'].str.replace(r'\W', ' ', regex=True)
```

```
In [58]: emp['Domain']
```

```
Out[58]: 0    Datascienc
         1        Testing
         2   Dataanalyst
         3     Analytics
         4   Statistics
         5        NLP
Name: Domain, dtype: object
```

```
In [60]: 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 [62]: emp['Age']
```

```
Out[62]: 0    34 years
         1    45' yr
         2      NaN
         3      NaN
         4    67-yr
         5    55yr
Name: Age, dtype: object
```

```
In [64]: emp['Age'] = emp['Age'].str.replace(r'\W', ' ', regex=True)
```

```
In [66]: emp['Age']
```

```
Out[66]: 0    34years
         1    45yr
         2      NaN
         3      NaN
         4    67yr
         5    55yr
Name: Age, dtype: object
```

```
In [72]: emp['Age'] = emp['Age'].astype(str).str.extract('(\d+)')
```

```
In [74]: emp['Age']
```

```
Out[74]: 0    34
         1    45
         2    NaN
         3    NaN
         4    67
         5    55
Name: Age, dtype: object
```

```
In [76]: emp
```

	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%0000	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 [78]: emp['Location'] = emp['Location'].str.replace(r'\W', '')
```

```
In [80]: emp['Location']
```

```
Out[80]: 0      Mumbai
         1    Bangalore
         2      NaN
         3    Hyderbad
         4      NaN
         5      Delhi
Name: Location, dtype: object
```

```
In [82]: emp['Salary']
```

```
Out[82]: 0    5^00#0
         1  10%0000
         2  1$5%0000
         3   2000^0
         4   30000-
         5  6000^$0
Name: Salary, dtype: object
```

```
In [89]: emp['Salary'] = emp['Salary'].str.replace(r'\W', '', regex=True)
```

```
In [91]: emp['Salary']
```

```
Out[91]: 0    5000
         1   10000
         2   15000
         3   20000
         4   30000
         5   60000
Name: Salary, dtype: object
```

In [93]: 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> yrs
3	Jane	Analytics	NaN	Hyderbad	20000	NaN
4	Uttam	Statistics	67	NaN	30000	5+ year
5	Kim	NLP	55	Delhi	60000	10+

In [95]: emp['Exp']

```
Out[95]: 0      2+
         1      <3
         2      4> yrs
         3      NaN
         4      5+ year
         5      10+
Name: Exp, dtype: object
```

In [99]: emp['Exp'] = emp['Exp'].astype(str).str.extract('(\d+)')

In [101...]: emp['Exp']

```
Out[101...]: 0      2
             1      3
             2      4
             3      NaN
             4      5
             5     10
Name: Exp, dtype: object
```

In [103...]: emp

Out[103...]

	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 [105...]

clean_data = emp.copy()

- missing values treatment for numerical data

In [108...]

clean_data

Out[108...]

	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 [110...]

clean_data['Age']

Out[110...]

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

In [112...]

import numpy as np

In [114...]

clean_data['Age'] = clean_data['Age'].fillna(np.mean(pd.to_numeric(clean_data['Age'])))

In [116...]

clean_data['Age']

```
Out[116... 0      34
       1      45
       2    50.25
       3    50.25
       4      67
       5      55
Name: Age, dtype: object
```

```
In [118... clean_data['Exp']
```

```
Out[118... 0      2
       1      3
       2      4
       3    NaN
       4      5
       5     10
Name: Exp, dtype: object
```

```
In [120... clean_data['Exp'] = clean_data['Exp'].fillna(np.mean(pd.to_numeric(clean_data['Exp'])))
```

```
In [122... clean_data['Exp']
```

```
Out[122... 0      2
       1      3
       2      4
       3    4.8
       4      5
       5     10
Name: Exp, dtype: object
```

```
In [124... 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		15000	4
3	Jane	Analytics	50.25	Hyderbad	20000	4.8
4	Uttam	Statistics	67		30000	5
5	Kim	NLP	55	Delhi	60000	10

```
In [126... clean_data['Location'].isnull().sum()
```

```
Out[126... np.int64(2)
```

```
In [128... clean_data['Location']
```

```
Out[128... 0      Mumbai
           1      Bangalore
           2      NaN
           3      Hyderabad
           4      NaN
           5      Delhi
Name: Location, dtype: object
```

```
In [130... clean_data['Location'] = clean_data['Location'].fillna(clean_data['Location'].mode()
```

```
In [133... clean_data['Location']
```

```
Out[133... 0      Mumbai
           1      Bangalore
           2      Bangalore
           3      Hyderabad
           4      Bangalore
           5      Delhi
Name: Location, dtype: object
```

```
In [135... 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 [137... 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 [139... clean_data['Age'] = clean_data['Age'].astype(int)
```

```
In [141... 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 [143...]: clean_data['Salary'] = clean_data['Salary'].astype(int)
clean_data['Exp'] = clean_data['Exp'].astype(int)
```

```
In [145...]: 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 [147...]: 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 [149...]: 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 [151...]: clean_data
```

Out[151...]

	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 [153...]

clean_data.to_csv('clean_data.csv')

In [155...]

```
import os
os.getcwd()
```

Out[155...]

'C:\\Users\\krish'

In [157...]

clean_data

Out[157...]

	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 [159...]

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

In [160...]

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

In [163...]

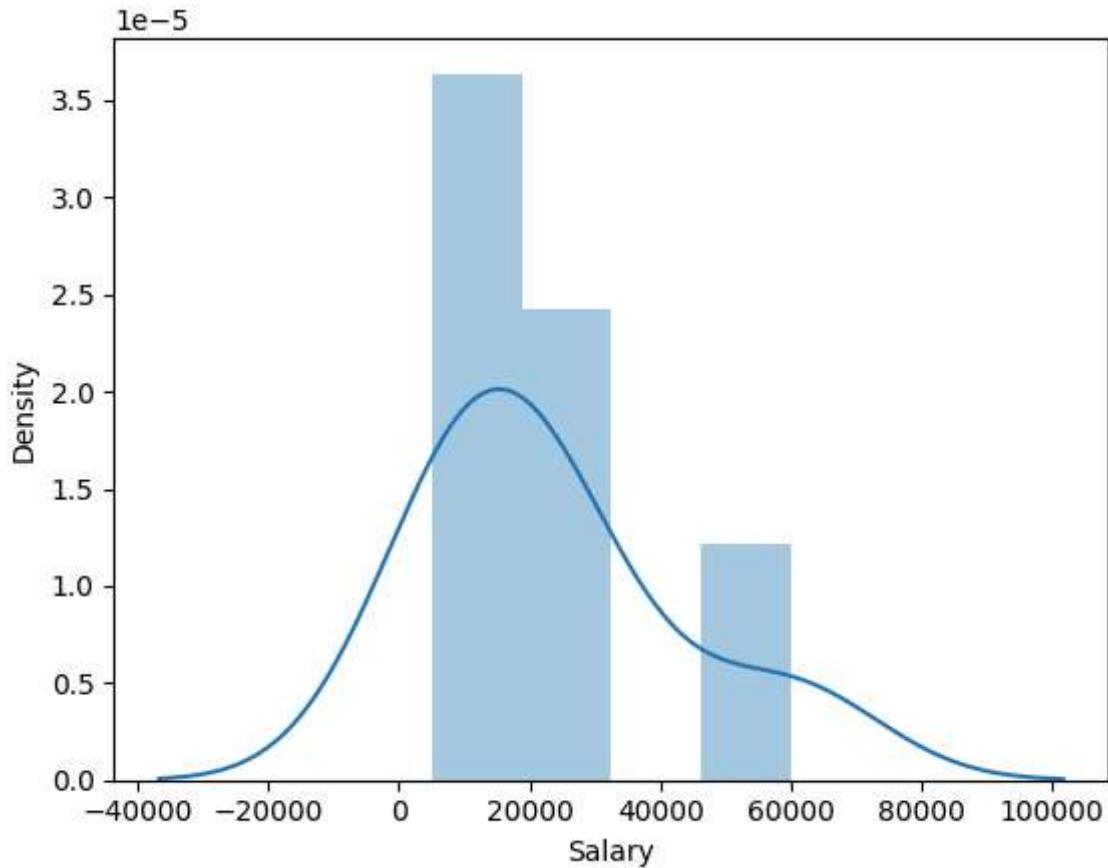
clean_data['Salary']

Out[163...]

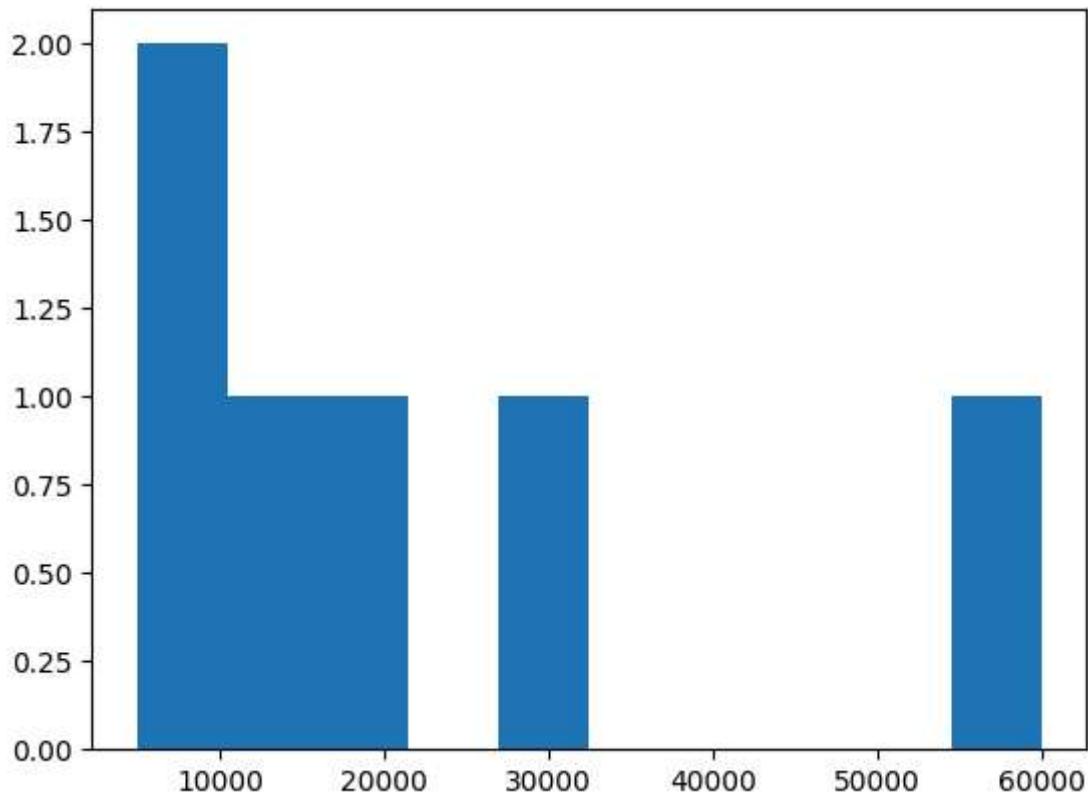
```
0      5000
1     10000
2     15000
3     20000
4     30000
5     60000
Name: Salary, dtype: int64
```

In [165...]

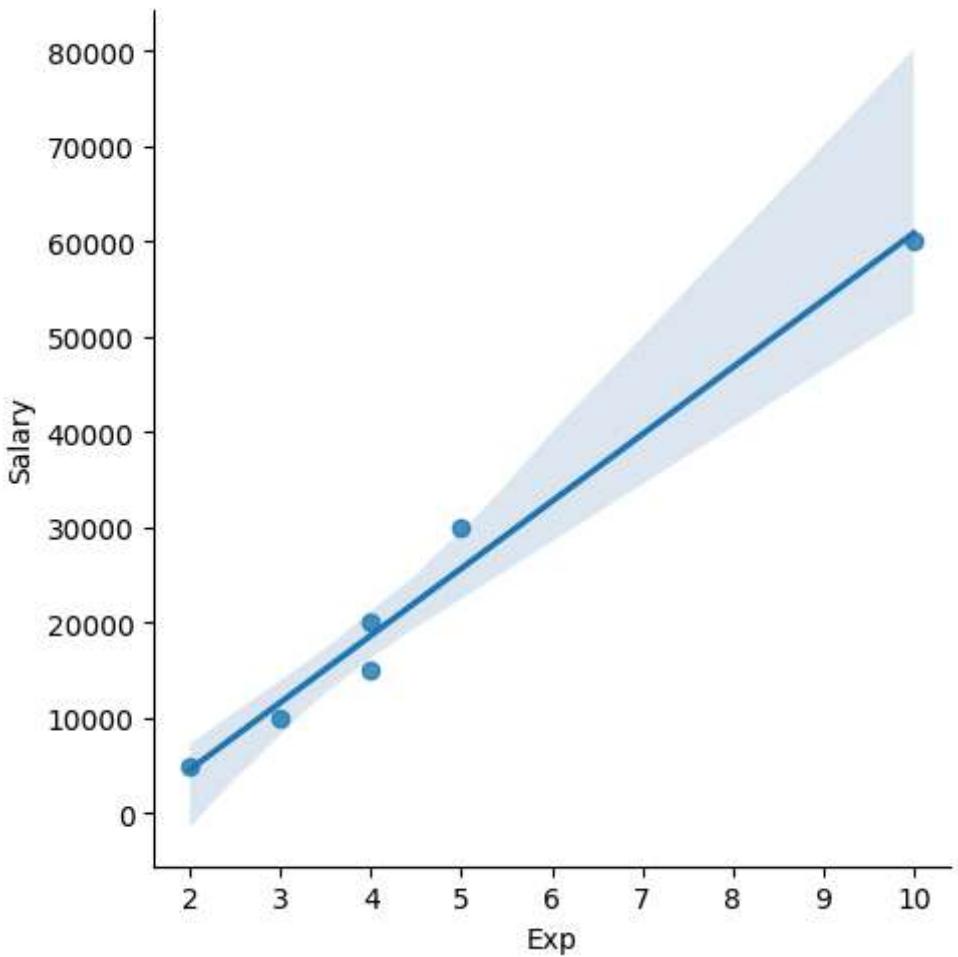
vis1 = sns.distplot(clean_data['Salary'])



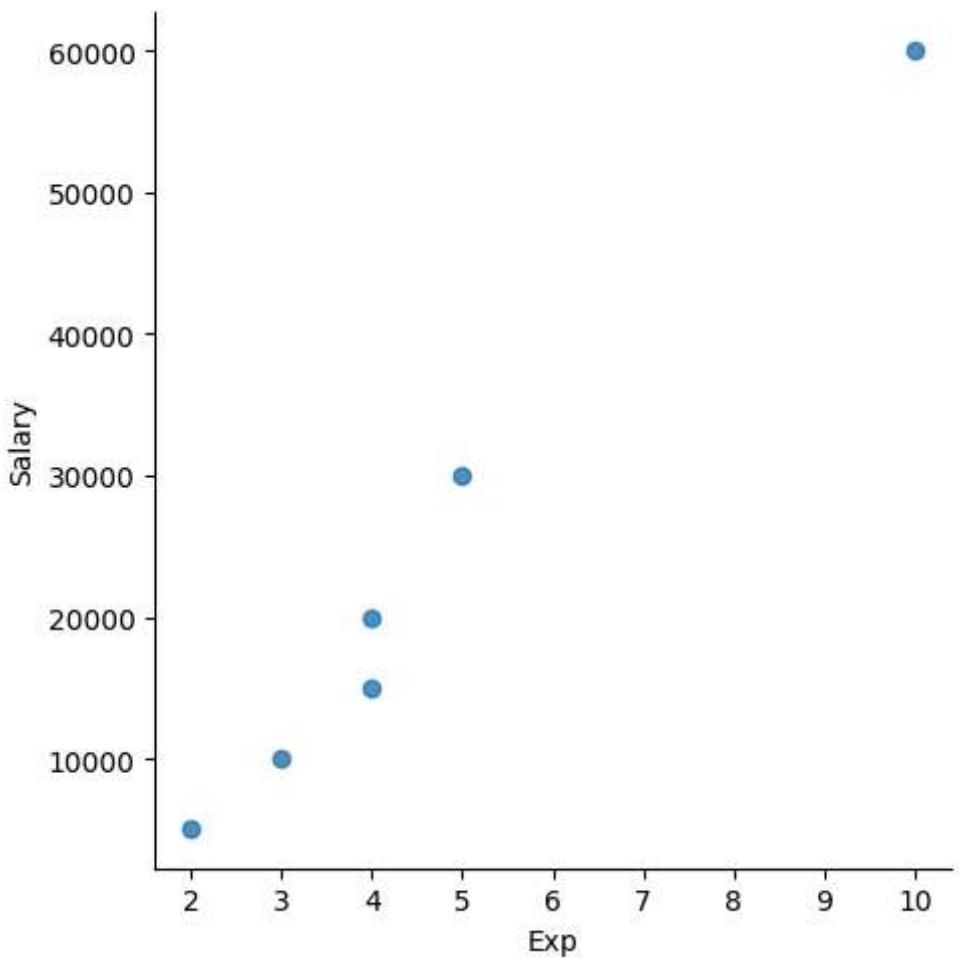
```
In [167]: vis2 = plt.hist(clean_data['Salary'])
```



```
In [169]: vis4 = sns.lmplot(data=clean_data,x = 'Exp', y='Salary')
```



```
In [171]: vis5 = sns.lmplot(data=clean_data,x = 'Exp', y='Salary', fit_reg = False)
```



```
In [173...]: clean_data[:]
```

```
Out[173...]:
```

	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 [175...]: clean_data[0:6:2]
```

```
Out[175...]:
```

	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 [177... `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 [179... `clean_data.columns`

Out[179... `Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')`

In [181... `X_iv = clean_data[['Name', 'Domain', 'Age', 'Location', 'Exp']]`

In [183... `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	Hyderabad	4
4	Uttam	Statistics	67	Bangalore	5
5	Kim	NLP	55	Delhi	10

In [185... `y_dv = clean_data[['Salary']]`

In [187... `y_dv`

Out[187... `Salary`

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

In [189... 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 [191... 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 [193... 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 [195... y_dv

Out[195...]

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

In [197...]

clean_data

Out[197...]

	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 [199...]

imputation = pd.get_dummies(clean_data)

In [201...]

imputation

Out[201...]

	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 [203...]

clean_data

Out[203...]

	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 []: