

```
In [127... import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [4]: emp = pd.read_excel(r"C:\Users\mohap\Rawdata.xlsx")
```

```
In [5]: emp
```

```
Out[5]:
```

	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 [6]: id(emp)
```

```
Out[6]: 2737312324096
```

```
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 [8]: `emp.tail()`

Out[8]:

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

In [9]: `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 [10]: emp.isnull()
```

```
Out[10]:
```

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

```
Out[21]: Name      0
Domain    0
Age       2
Location  2
Salary    0
Exp       1
dtype: int64
```

```
In [23]: emp
```

```
Out[23]:
```

	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 [25]: emp['Name']
```

```
Out[25]: 0    Mike
          1    Teddy^
          2    Uma#r
          3    Jane
          4    Uttam*
          5    Kim
          Name: Name, dtype: object
```

```
In [27]: emp['Name']=emp['Name'].str.replace(r'\W','',regex=True) # non word char
```

```
In [29]: emp['Name']
```

```
Out[29]: 0    Mike
          1    Teddy
          2    Umar
          3    Jane
          4    Uttam
          5    Kim
          Name: Name, dtype: object
```

```
In [31]: emp
```

```
Out[31]:
```

	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 [35]: emp['Domain']=emp['Domain'].str.replace(r'\W','',regex=True) # non word char
```

```
In [37]: emp['Domain']
```

```
Out[37]: 0    Datascience
         1      Testing
         2    Dataanalyst
         3      Analytics
         4      Statistics
         5          NLP
         Name: Domain, dtype: object
```

```
In [39]: emp['Age']=emp['Age'].str.replace(r'\W','',regex=True) # non word char
```

```
In [41]: emp['Age']
```

```
Out[41]: 0    34years
         1     45yr
         2      NaN
         3      NaN
         4     67yr
         5     55yr
         Name: Age, dtype: object
```

```
In [45]: emp['Age'] = emp['Age'].str.extract("(\\d+)")
```

```
In [47]: emp['Age']
```

```
Out[47]: 0     34
         1     45
         2    NaN
         3    NaN
         4     67
         5     55
         Name: Age, dtype: object
```

```
In [49]: emp
```

Out[49]:

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

```
In [53]: emp['Location']
```

Out[53]:

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

Name: Location, dtype: object

```
In [55]: emp['Salary'] = emp['Salary'].str.replace(r'\W', '', regex=True)
```

```
In [57]: emp['Salary']
```

Out[57]:

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

Name: Salary, dtype: object

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

```
Out[59]: 0      2+
          1      <3
          2      4> yrs
          3      NaN
          4      5+ year
          5      10+
          Name: Exp, dtype: object
```

```
In [61]: emp['Exp'] = emp['Exp'].str.extract('(\d+)')
```

```
In [63]: emp['Exp']
```

```
Out[63]: 0      2
          1      3
          2      4
          3      NaN
          4      5
          5      10
          Name: Exp, dtype: object
```

```
In [65]: emp
```

```
Out[65]:
```

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

```
In [69]: clean_data
```

```
Out[69]:
```

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

```
In [75]: clean_data['Age']
```

```
Out[75]: 0      34
1      45
2    50.25
3    50.25
4      67
5      55
Name: Age, dtype: object
```

```
In [77]: clean_data['Exp']
```

```
Out[77]: 0      2
1      3
2      4
3    NaN
4      5
5     10
Name: Exp, dtype: object
```

```
In [79]: clean_data['Exp'].isnull().sum()
```

```
Out[79]: 1
```



```
In [83]: # Convert Exp column from string to numeric
clean_data['Exp'] = pd.to_numeric(clean_data['Exp'], errors='coerce')
# Now fill missing value with median and convert to int
clean_data['Exp'] = clean_data['Exp'].fillna(clean_data['Exp'].median()).astype(int)
```

```
In [85]: clean_data['Exp']
```

```
Out[85]: 0      2
         1      3
         2      4
         3      4
         4      5
         5     10
         Name: Exp, dtype: int32
```

```
In [89]: clean_data['Location'].isnull().sum()
```

```
Out[89]: 2
```

```
In [93]: clean_data['Location']=clean_data['Location'].fillna(clean_data['Location'].mode()[0])
```

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

```
Out[95]: 0      Mumbai
         1    Bangalore
         2    Bangalore
         3     Hyderabad
         4    Bangalore
         5         Delhi
         Name: Location, dtype: object
```

```
In [97]: 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 [99]: `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      int32
dtypes: int32(1), object(5)
memory usage: 396.0+ bytes

```

In [101... `clean_data['Age'] = clean_data['Age'].astype(int)`

In [103... `clean_data['Salary'] = clean_data['Salary'].astype(int)`

In [105... `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 [107... 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 [109... 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 [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	Hyderbad	20000	4
4	Uttam	Statistics	67	Bangalore	30000	5
5	Kim	NLP	55	Delhi	60000	10

In [113...

```
clean_data.to_csv("clean_data.csv")
```

In [115...

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

Out[115...

```
'C:\\Users\\mohap'
```

In [117...

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

In [119...

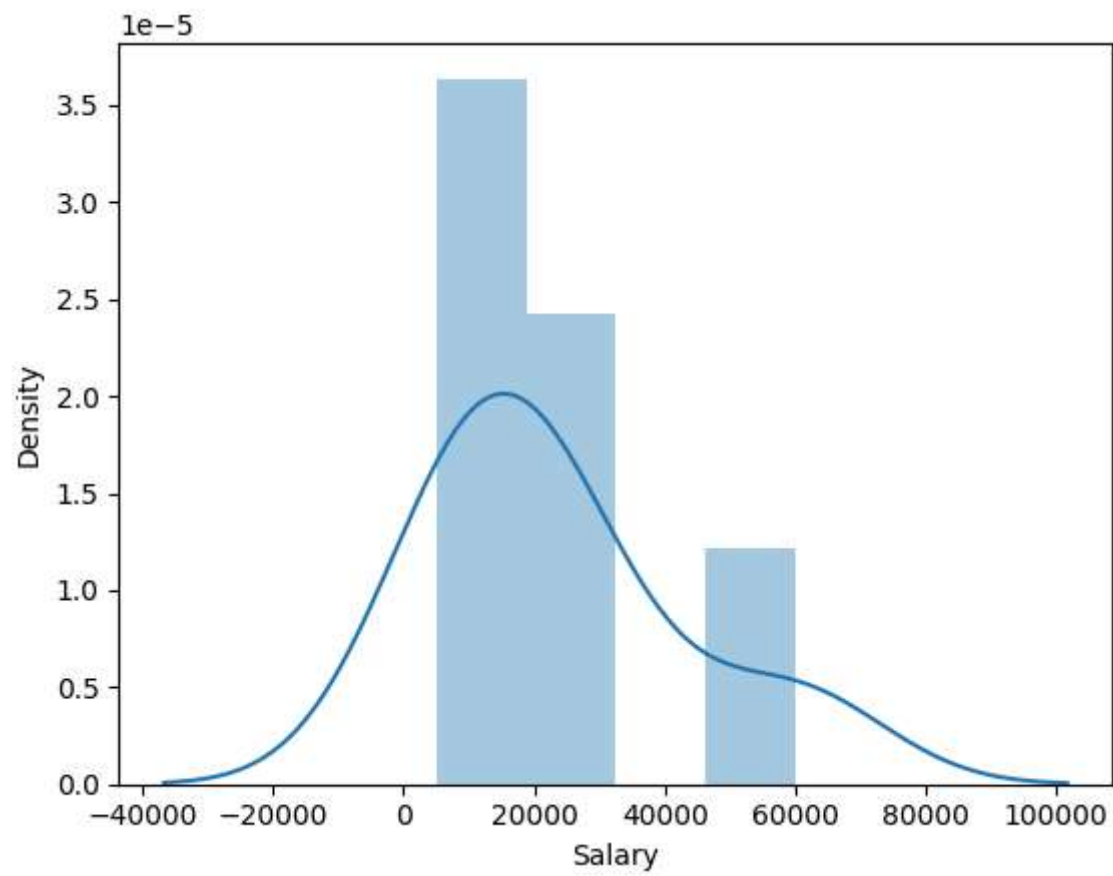
```
clean_data['Salary']
```

Out[119...

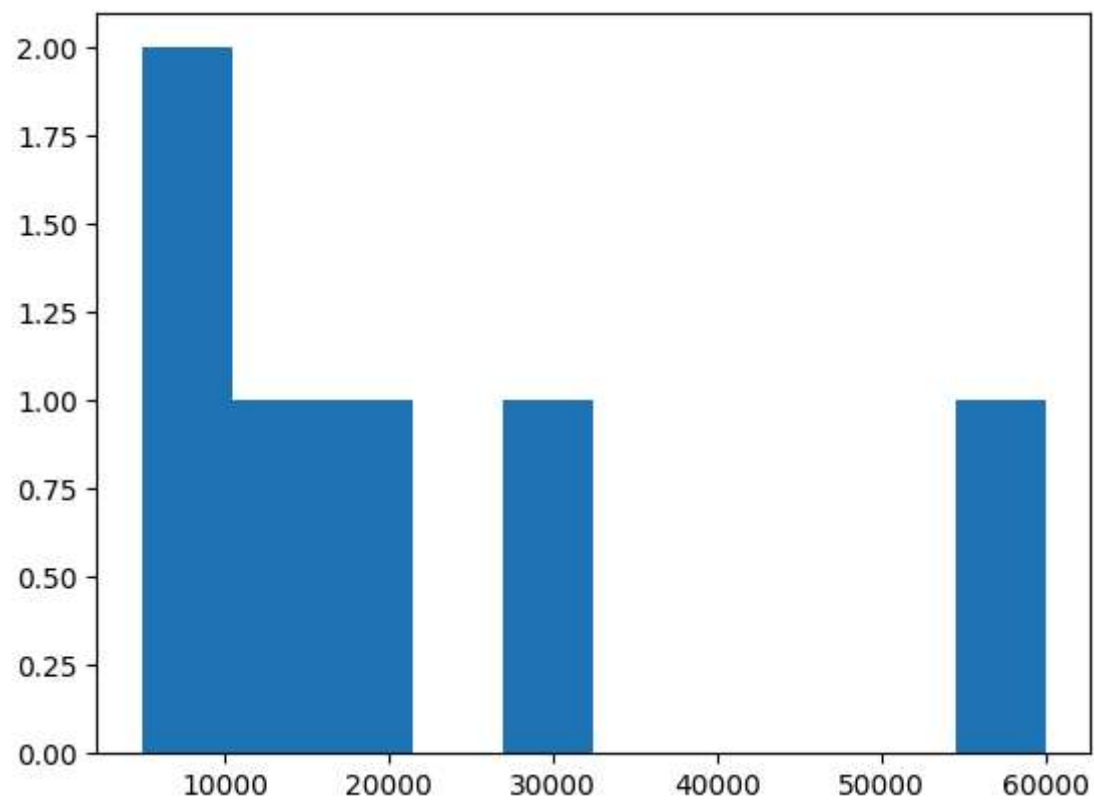
```
0    5000  
1   10000  
2   15000  
3   20000  
4   30000  
5   60000  
Name: Salary, dtype: int32
```

In [121...

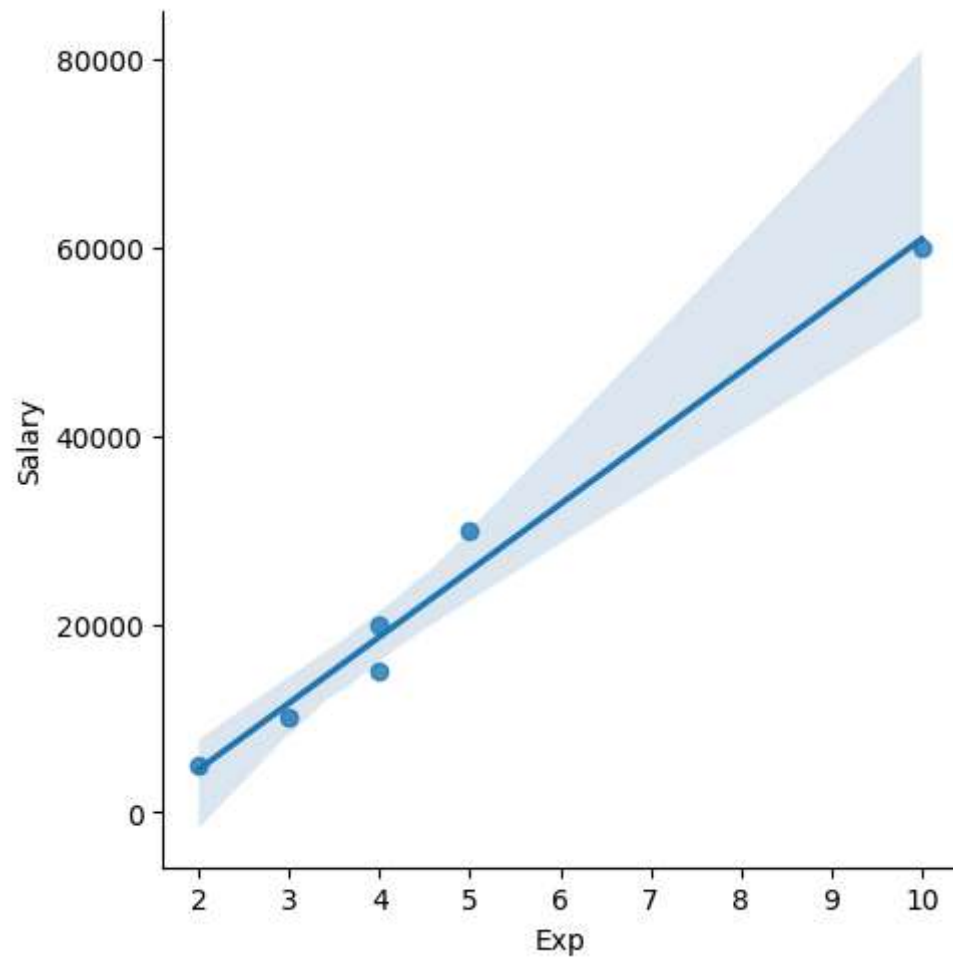
```
vis1 = sns.distplot(clean_data['Salary'])
```



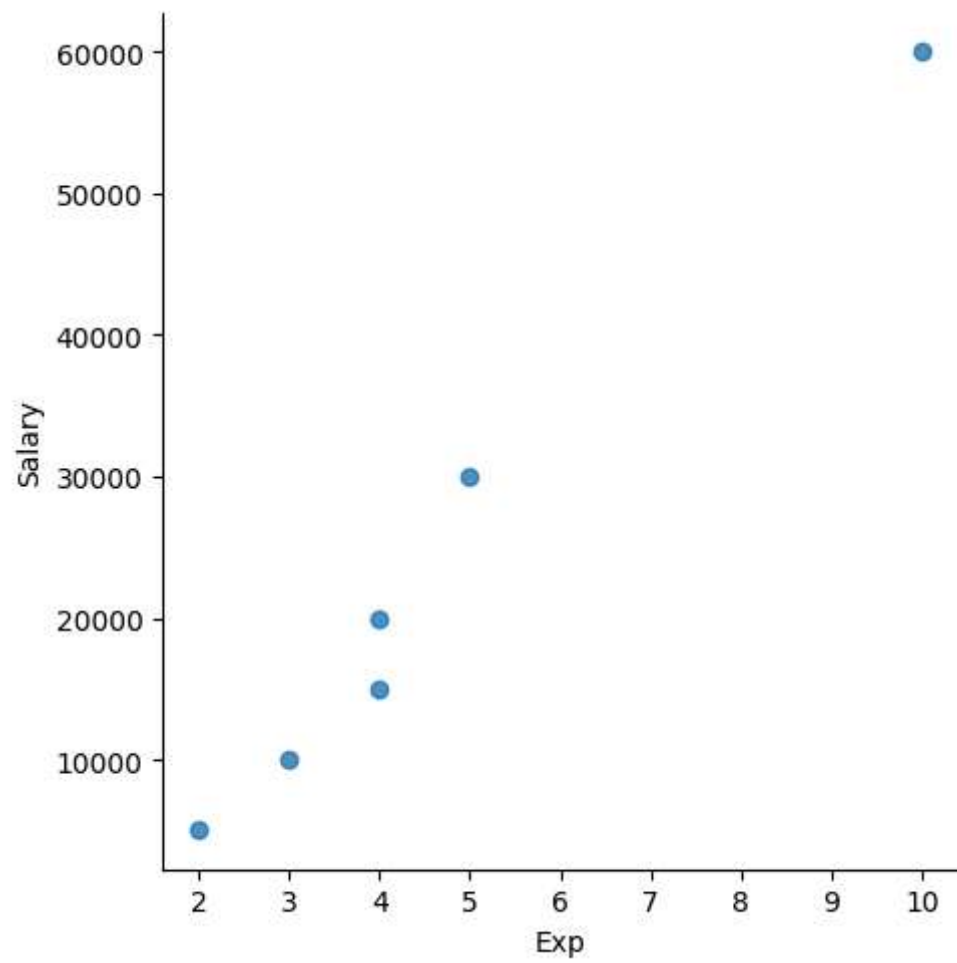
```
In [129... vis2 =plt.hist(clean_data['Salary'])  
plt.show()
```



```
In [131... vis4 = sns.lmplot(data=clean_data,x='Exp',y='Salary')
```



```
In [135... vis5 = sns.lmplot(data = clean_data,x = 'Exp',y = 'Salary',fit_reg=False)
```



```
In [137... x_iv = clean_data[['Name','Domain','Age','Location','Exp']]
```

```
In [139... x_iv
```


Out[139...

	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 [141...

```
y_dv = clean_data[['Salary']]
```

In [143...

```
y_dv
```

Out[143...

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

In [145...

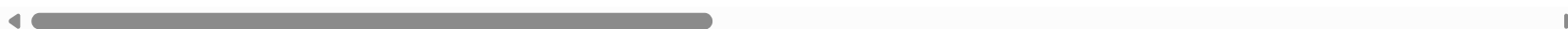
```
imputation = pd.get_dummies(clean_data, dtype=int)
```

In [147...

```
imputation
```

```
Out[147...
```

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



```
In [151...] len(clean_data)
```

```
Out[151...] 6
```

```
In [153...] imputation.columns
```

```
Out[153...] Index(['Age', 'Salary', 'Exp', 'Name_Jane', 'Name_Kim', 'Name_Mike',
      'Name_Teddy', 'Name_Umar', 'Name_Uttam', 'Domain_Analytics',
      'Domain_Dataanalyst', 'Domain_Datascience', 'Domain_NLP',
      'Domain_Statistics', 'Domain_Testing', 'Location_Bangalore',
      'Location_Delhi', 'Location_Hyderabad', 'Location_Mumbai'],
      dtype='object')
```

```
In [155...] len(imputation.columns)
```

```
Out[155...] 19
```

```
In [ ]:
```