

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

```
In [2]: pd.__version__
```

```
Out[2]: '2.2.2'
```

```
In [3]: emp = pd.read_excel(r"D:\Data Science & AI\Rawdata.xlsx")
```

```
In [4]: emp
```

```
Out[4]:
```

	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 [5]: emp.columns
```

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

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

```
Out[6]: (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	Hyderbad	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	Hyderbad	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

Out[10]:

	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 [12]: emp.isnull()

Out[12]:

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

```
Out[13]: Name      0
         Domain    0
         Age       2
         Location   2
         Salary     0
         Exp       1
         dtype: int64
```

DATA CLEANING OR DATA CLEANSING

```
In [14]: emp['Name']
```

```
Out[14]: 0      Mike
         1  Teddy^
         2  Uma#r
         3      Jane
         4  Uttam*
         5      Kim
         Name: Name, dtype: object
```

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

```
In [18]: emp['Name']
```

```
Out[18]: 0      Mike
         1  Teddy
         2  Umar
         3  Jane
         4  Uttam
         5  Kim
         Name: Name, dtype: object
```

```
In [19]: emp['Domain']
```

```
Out[19]: 0  Datascience#$
         1      Testing
         2  Dataanalyst^^#
         3  Ana^lytics
         4  Statistics
         5      NLP
         Name: Domain, dtype: object
```

```
In [22]: emp['Domain'] = emp['Domain'].str.replace(r'\W', '', regex = True)
```

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

```
Out[23]: 0  Datascience
         1      Testing
         2  Dataanalyst
         3  Analytics
         4  Statistics
         5      NLP
         Name: Domain, dtype: object
```

```
In [25]: emp['Location'] = emp['Location'].str.replace(r'\W', '', regex = True)
```

```
In [26]: emp['Location']
```

```
Out[26]: 0      Mumbai
1      Bangalore
2         NaN
3      Hyderabad
4         NaN
5         Delhi
Name: Location, dtype: object
```

```
In [27]: emp['Age'] = emp['Age'].str.extract('(\d+)')
```

```
<>:1: SyntaxWarning: invalid escape sequence '\d'
<>:1: SyntaxWarning: invalid escape sequence '\d'
C:\Users\HARIKAREDDY\AppData\Local\Temp\ipykernel_2204\1884116463.py:1: SyntaxWarning: invalid escape sequence '\d'
emp['Age'] = emp['Age'].str.extract('(\d+)')
```

```
In [28]: emp['Age']
```

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

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

```
In [30]: emp['Salary']
```

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

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

```
Out[31]:
```

	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	Hyderabad	20000	NaN
4	Uttam	Statistics	67	NaN	30000	5+ year
5	Kim	NLP	55	Delhi	60000	10+

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

```
Out[32]: 0      2+
        1      <3
        2      4> yrs
        3      NaN
        4      5+ year
        5      10+
        Name: Exp, dtype: object
```

```
In [33]: emp['Exp'] = emp['Exp'].str.extract('(\d+)')
```

```
<>:1: SyntaxWarning: invalid escape sequence '\d'
<>:1: SyntaxWarning: invalid escape sequence '\d'
C:\Users\HARIKAREDDY\AppData\Local\Temp\ipykernel_2204\4097307645.py:1: SyntaxWarning: invalid escape sequence '\d'
emp['Exp'] = emp['Exp'].str.extract('(\d+)')
```

```
In [34]: emp['Exp']
```

```
Out[34]: 0      2
        1      3
        2      4
        3      NaN
        4      5
        5      10
        Name: Exp, dtype: object
```

```
In [35]: emp
```

```
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

```
In [36]: clean_data = emp.copy()
```

- Missing value treatment for numerical data

```
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['Age']
```

```
Out[38]: 0      34
1      45
2      NaN
3      NaN
4      67
5      55
Name: Age, dtype: object
```

```
In [39]: import numpy as np
```

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

```
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 [44]: clean_data['Exp']
```

```
Out[44]: 0      2
1      3
2      4
3      NaN
4      5
5     10
Name: Exp, dtype: object
```

```
In [45]: clean_data['Exp'] = clean_data['Exp'].fillna(np.mean(pd.to_numeric(clean_data['E
```

```
In [46]: clean_data['Exp']
```

```
Out[46]: 0      2
1      3
2      4
3     4.8
4      5
5     10
Name: Exp, dtype: object
```

```
In [47]: clean_data
```

```
Out[47]:
```

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

```
Out[48]: 2
```

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

```
Out[49]: 0      Mumbai
1    Bangalore
2         NaN
3    Hyderbad
4         NaN
5        Delhi
Name: Location, dtype: object
```

```
In [50]: clean_data['Location'] = clean_data['Location'].fillna(clean_data['Location'].mode()[0])
```

```
In [51]: clean_data['Location']
```

```
Out[51]: 0      Mumbai
1    Bangalore
2    Bangalore
3    Hyderbad
4    Bangalore
5        Delhi
Name: Location, dtype: object
```

```
In [52]: clean_data
```

```
Out[52]:
```

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

```
In [55]: 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 [56]: clean_data['Salary'] = clean_data['Salary'].astype(int)
```

```
In [57]: clean_data['Exp'] = clean_data['Exp'].astype(int)
```

```
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     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 [59]: 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 [60]: 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 [61]: `clean_data`

Out[61]:

	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 [62]: `clean_data.to_csv('clean_data.csv')`

In [63]: `import os`
`os.getcwd()`

Out[63]: `'C:\\Users\\HARIKAREDDY'`

In [64]: `clean_data`

Out[64]:

	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

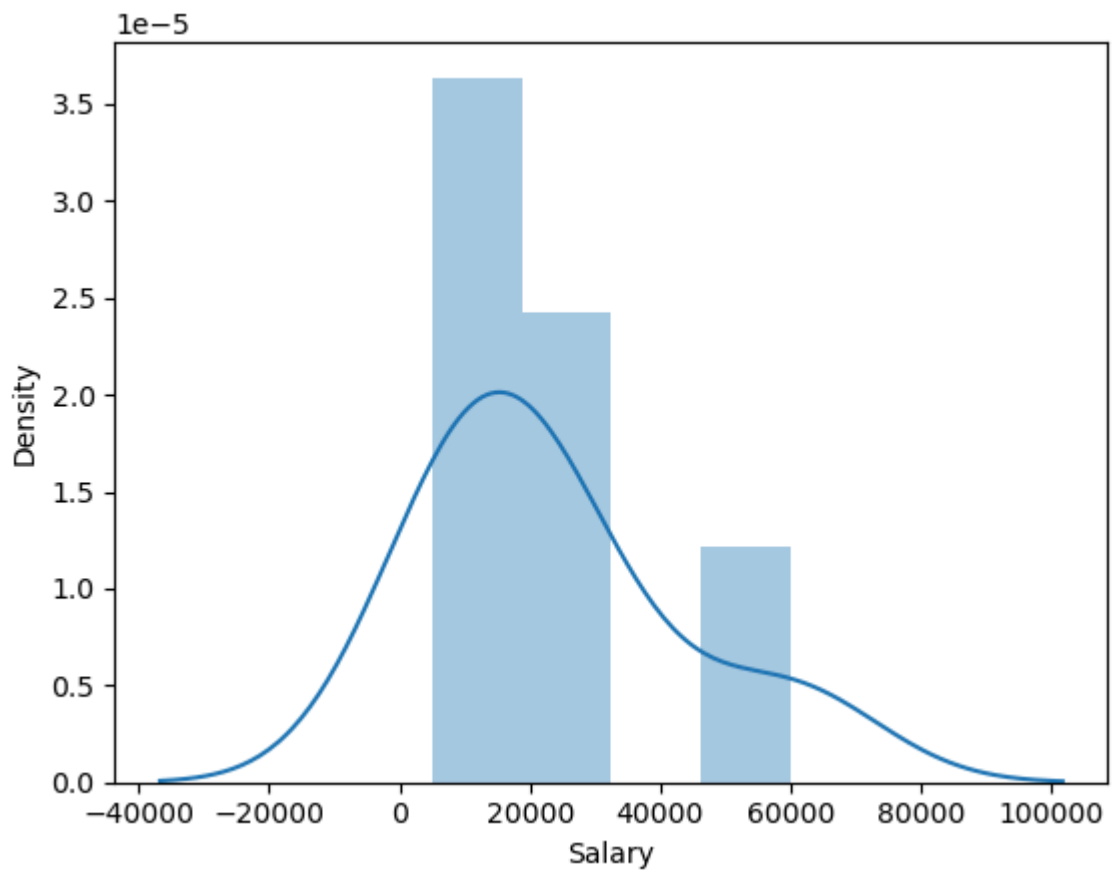
In [65]: `# EDA TECHNIQUE LETS APPLY`
`import matplotlib.pyplot as plt #visualization`
`import seaborn as sns`

In [66]: `import warnings`
`warnings.filterwarnings('ignore')`

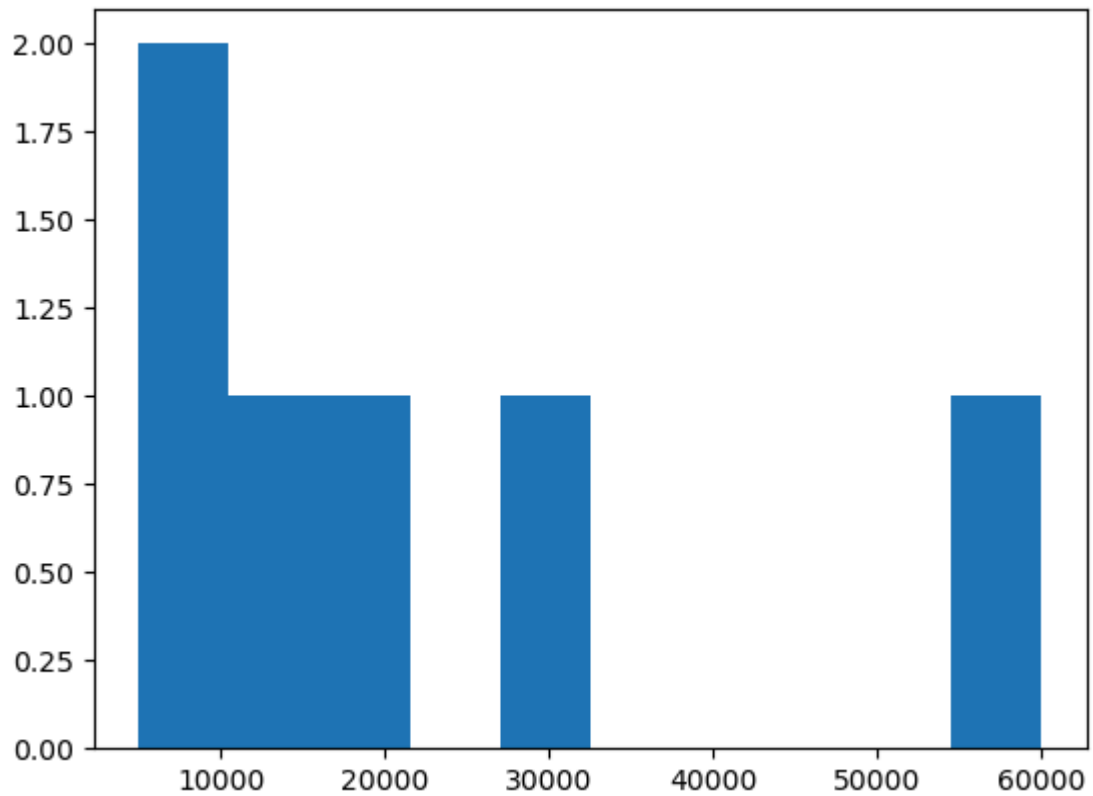
```
In [67]: clean_data['Salary']
```

```
Out[67]: 0      5000  
         1     10000  
         2     15000  
         3     20000  
         4     30000  
         5     60000  
         Name: Salary, dtype: int32
```

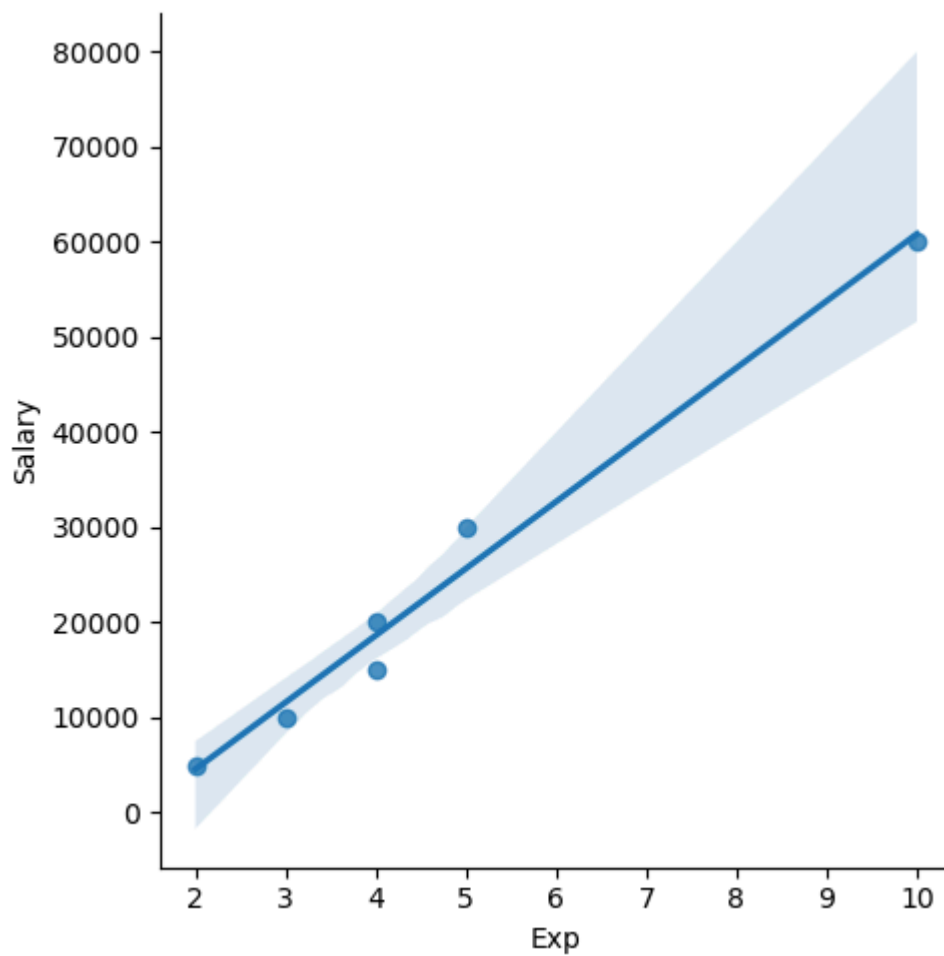
```
In [69]: vis1 = sns.distplot(clean_data['Salary'])
```



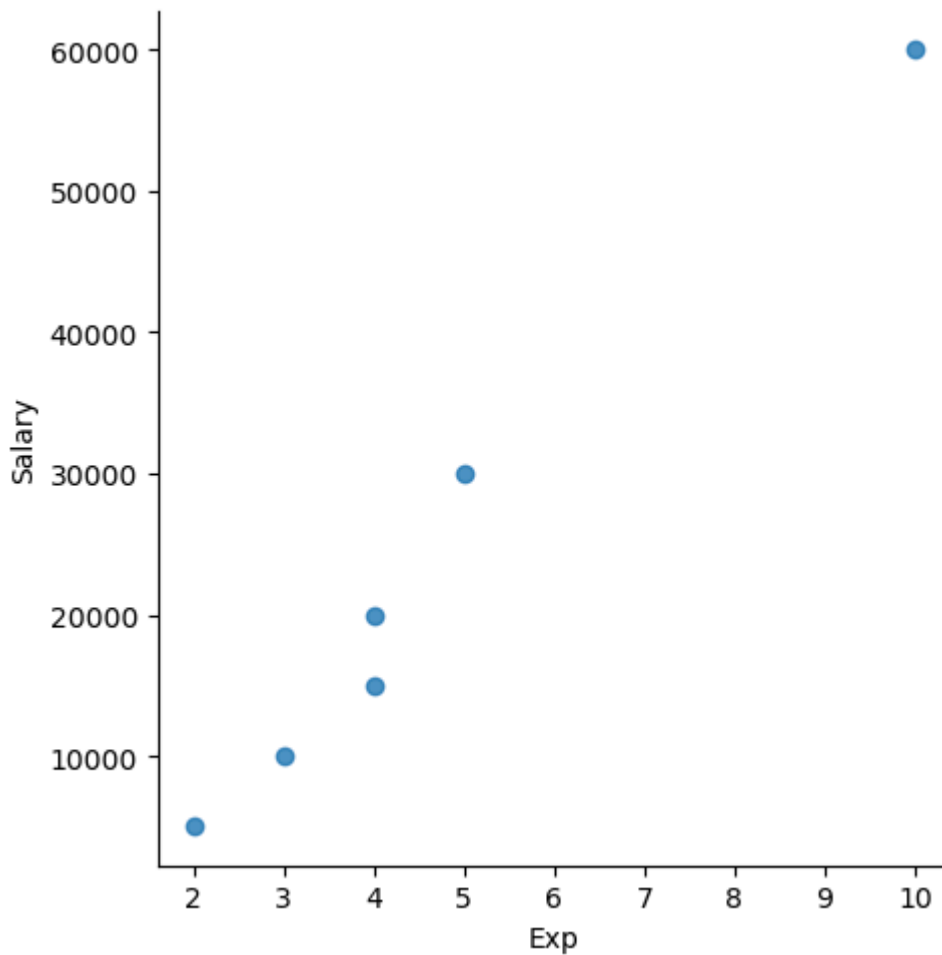
```
In [70]: vis2 = plt.hist(clean_data['Salary'])
```



```
In [71]: vis4 = sns.lmplot(data=clean_data,x = 'Exp',y = 'Salary')
```



```
In [72]: vis5 = sns.lmplot(data=clean_data,x = 'Exp',y = 'Salary',fit_reg = False)
```



```
In [73]: clean_data[:]
```

```
Out[73]:
```

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

```
Out[74]:
```

	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 [75]: clean_data[:, :-1]
```

Out[75]:

	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 [77]: `clean_data.columns`

Out[77]: `Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')`

In [78]: `x_iv = clean_data[['Name', 'Domain', 'Age', 'Location', 'Exp']]`

In [79]: `x_iv`

Out[79]:

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

In [83]: `y_dv`

Out[83]:

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

In [84]: `emp`

Out[84]:

	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 [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 [86]: `x_iv`

Out[86]:

	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 [87]: `y_dv`

Out[87]: **Salary**

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

In [88]: `clean_data`

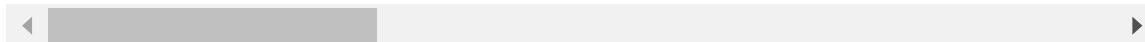
Out[88]:

	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 [96]: `imputation = pd.get_dummies(clean_data)`In [97]: `imputation`

Out[97]:

	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 [92]: `clean_data`

Out[92]:

	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 [93]: imputation

Out[93]:

	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 [99]: imputation.astype(int)

Out[99]:

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

In [100... imputation.columns

Out[100... 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 [101... len(imputation)

Out[101... 6

In [102... `imputation.shape`

Out[102... (6, 19)

In [103... `len(imputation.columns)`

Out[103... 19

In []: