

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

```
In [2]: emp = pd.read_excel(r"C:\Users\Asus.LAPTOP-EMBE8J7O\Downloads\Rawdata.xlsx")
```

```
In [3]: 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	Hyderbad	2000^0	NaN
4	Uttam*	Statistics	67-yr	NaN	30000-	5+ year
5	Kim	NLP	55yr	Delhi	6000^\$0	10+

```
In [4]: pd.__version__
```

```
Out[4]: '1.4.2'
```

```
In [5]: id(emp)
```

```
Out[5]: 1505493363824
```

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

```
Out[6]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')
```

```
In [7]: emp.shape
```

```
Out[7]: (6, 6)
```

```
In [8]: emp.head()
```

```
Out[8]:
```

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

Out[9]:

	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 [10]:

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: 416.0+ bytes
```

In [11]:

emp

Out[11]:

	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.isna()`

Out[13]:

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

Out[14]:

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

In [15]: `emp.columns`

Out[15]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')

DATA CLEANING OR DATA CLEANSING

In [16]: `emp['Name']`

Out[16]:

```
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)`
`emp['Name']`

Out[17]:

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

In [18]: `emp`

Out[18]:

	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 [19]: emp['Domain']

Out[19]:

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

In [20]: emp['Domain'] = emp['Domain'].str.replace(r'\W','',regex=True)

In [21]: emp['Domain']

Out[21]:

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

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

In [23]: emp['Age']

Out[23]:

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

In [24]: emp['Age'] = emp['Age'].str.extract('(\d+)')

In [25]: emp['Age']

Out[25]:

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

In [26]: emp

Out[26]:

	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 [27]: emp['Location']

Out[27]:

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

In [28]: emp['Location'] = emp['Location'].str.replace(r'\W', '', regex=True)

In [29]: emp['Location']

Out[29]:

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

In [30]: emp

Out[30]:

	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 [31]: emp['Salary']

```
Out[31]: 0      5^00#0
1      10%%000
2      1$5%000
3      2000^0
4      30000-
5      6000^$0
Name: Salary, dtype: object
```

```
In [32]: emp['Salary'] = emp['Salary'].str.replace(r'\W', '', regex=True)
```

```
In [33]: emp['Salary']
```

```
Out[33]: 0      5000
1     10000
2     15000
3     20000
4     30000
5     60000
Name: Salary, dtype: object
```

```
In [34]: emp
```

```
Out[34]:
```

	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 [35]: emp['Exp']
```

```
Out[35]: 0      2+
1      <3
2      4> yrs
3      NaN
4      5+ year
5      10+
Name: Exp, dtype: object
```

```
In [36]: emp['Exp'] = emp['Exp'].str.extract('(\d+)')
```

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

```
Out[37]: 0      2
1      3
2      4
3      NaN
4      5
5     10
Name: Exp, dtype: object
```

```
In [38]: emp
```

Out[38]:

	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()`

In [40]: `clean_data`

Out[40]:

	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 [41]: `clean_data.isnull().sum()`

Out[41]:

```

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

```

In [42]: `import numpy as np`

In [43]: `clean_data['Age'] = clean_data['Age'].fillna(np.mean(pd.to_numeric(clean_data['Age'])))`

In [44]: `clean_data['Age']`

Out[44]:

```

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

```

In [45]: `clean_data`

Out[45]:

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

In [46]: `clean_data['Exp'] = clean_data['Exp'].fillna(np.mean(pd.to_numeric(clean_data['Exp'])))`

In [47]: `clean_data['Exp']`

Out[47]:

0	2
1	3
2	4
3	4.8
4	5
5	10

Name: Exp, dtype: object

In [48]: `clean_data`

Out[48]:

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

Out[49]: 2

In [50]: `clean_data['Location']`

Out[50]:

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

Name: Location, dtype: object

In [51]: `clean_data['Location'] = clean_data['Location'].fillna(clean_data['Location'].mode()[0])`

In [52]: `clean_data['Location']`


```
Out[52]: 0      Mumbai
1      Bangalore
2      Bangalore
3      Hyderabad
4      Bangalore
5      Delhi
Name: Location, dtype: object
```

```
In [53]: clean_data
```

```
Out[53]:
```

	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 [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      object
3   Location      6 non-null      object
4   Salary        6 non-null      object
5   Exp           6 non-null      object
dtypes: object(6)
memory usage: 416.0+ bytes
```

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

```
In [57]: clean_data
```

```
Out[57]:
```

	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 [59]: 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: 392.0+ bytes
```

```
In [61]: clean_data['Exp'] = clean_data['Exp'].astype(int)
```

```
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	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 [63]: 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      int32
dtypes: int32(2), object(4)
memory usage: 368.0+ bytes
```

```
In [64]: clean_data['Salary'] = clean_data['Salary'].astype(int)
```

```
In [65]: clean_data
```

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	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 [66]: `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: 344.0+ bytes
```

In [67]: `clean_data`

Out[67]:

	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 [70]: `clean_data['Name'] = clean_data['Name'].astype('category')`

In [71]: `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     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: category(1), int32(3), object(2)
memory usage: 522.0+ bytes
```

```
In [72]: clean_data['Domain'] = clean_data['Domain'].astype('category')
```

```
In [73]: 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     object
4   Salary      6 non-null     int32
5   Exp         6 non-null     int32
dtypes: category(2), int32(3), object(1)
memory usage: 700.0+ bytes
```

```
In [74]: clean_data['Location'] = clean_data['Location'].astype('category')
```

```
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     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: 862.0 bytes
```

```
In [79]: clean_data
```

Out[79]:

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

In [81]: `import os`
`os.getcwd()`

Out[81]: 'C:\\Users\\Asus.LAPTOP-EMBE8J70'

In [82]: `clean_data`

Out[82]:

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

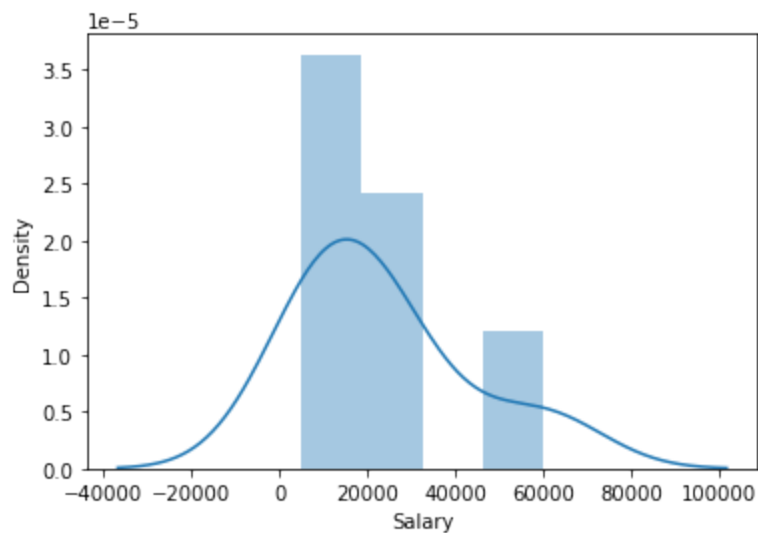
In [87]: `import warnings`
`warnings.filterwarnings('ignore')`

In [88]: `clean_data['Salary']`

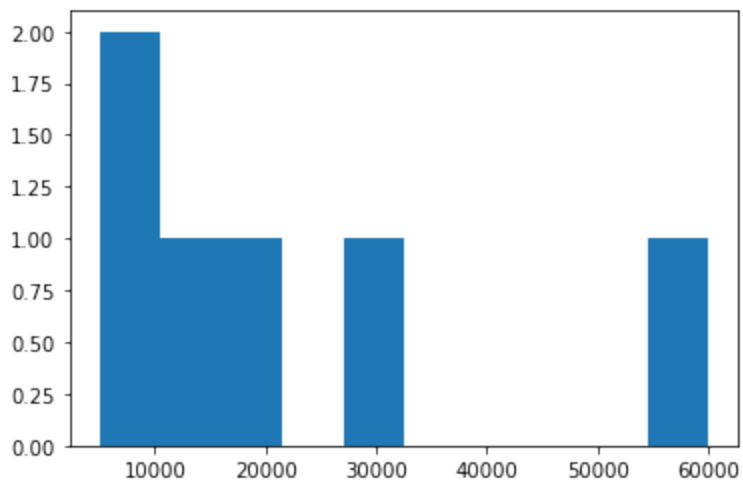
Out[88]:

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

In [91]: `vis1 = sns.distplot(clean_data['Salary'])`



```
In [93]: vis2 = plt.hist(clean_data['Salary'])
```

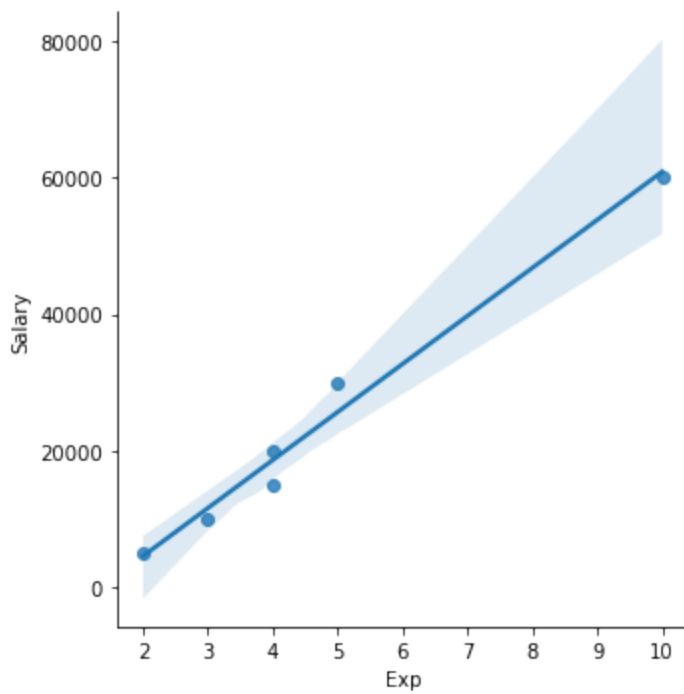


```
In [94]: clean_data
```

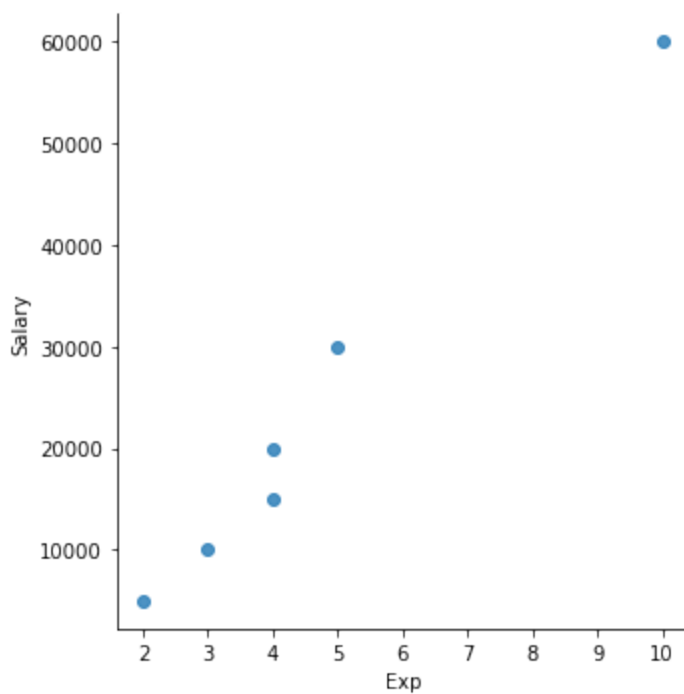
```
Out[94]:
```

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



```
In [100...] vis5 = sns.lmplot(data = clean_data, x = 'Exp', y = 'Salary', fit_reg=False)
```



```
In [101...] clean_data[0:6:2]
```

```
Out[101]:
```

	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 [102...] clean_data[:, :-1]
```

Out[102]:

	Name	Domain	Age	Location	Salary	Exp
5	Kim	NLP	55	Delhi	60000	10
4	Uttam	Statistics	67	Bangalore	30000	5
3	Jane	Analytics	50	Hyderbad	20000	4
2	Umar	Dataanalyst	50	Bangalore	15000	4
1	Teddy	Testing	45	Bangalore	10000	3
0	Mike	Datascience	34	Mumbai	5000	2

In [103... `clean_data.columns`

Out[103]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')

In [105... `clean_data`

Out[105]:

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

In [108... `x_iv`

Out[108]:

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

In [112... `y_dv`

Out[112]: **Salary**

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

In [113... emp

Out[113]:

	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 [114... clean_data

Out[114]:

	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 [115... imputation = pd.get_dummies(clean_data, dtype=int)

In [116... imputation

Out[116]:

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

In [117... clean_data

Out[117]:

	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 [118... len(clean_data)

Out[118]: 6

In [119... imputation.columns

Out[119]: 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 [120... len(imputation.columns)

Out[120]: 19

In []: