

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

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

```
Out[2]: '2.0.3'
```

```
In [3]: emp = pd.read_excel("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]: id(emp)
```

```
Out[5]: 1677872730384
```

```
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()
```

	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+

```
>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 [11]: emp.isnull()
```

Out[11]:

	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]: emp.isna()
```

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
```

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

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

```
In [15]: emp['Domain']
```

```
Out[15]: 0    Datascience#$
         1      Testing
         2  Dataanalyst^^#
         3    Ana^^lytics
         4    Statistics
         5           NLP
         Name: Domain, dtype: object
```

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

```
Out[16]: 0    34 years
         1    45' yr
         2      NaN
         3      NaN
         4    67-yr
         5    55yr
         Name: Age, dtype: object
```

```
In [17]: emp['Location']
```

```
Out[17]: 0      Mumbai
         1    Bangalore
         2      NaN
         3    Hyderbad
         4      NaN
         5      Delhi
         Name: Location, dtype: object
```

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

```
Out[18]: 0      5^00#0
         1    10%%000
         2    1$5%000
         3    2000^0
         4    30000-
         5    6000^$0
         Name: Salary, dtype: object
```

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

Out[19]:

	Name	Domain
0	Mike	Datascience#\$
1	Teddy^	Testing
2	Uma#r	Dataanalyst^^#
3	Jane	Ana^^lytics
4	Uttam*	Statistics
5	Kim	NLP

```
In [20]: emp[['Name', 'Domain', 'Age']]
```

Out[20]:

	Name	Domain	Age
0	Mike	Datascience#\$	34 years
1	Teddy^	Testing	45' yr
2	Uma#r	Dataanalyst^^#	NaN
3	Jane	Ana^^lytics	NaN
4	Uttam*	Statistics	67-yr
5	Kim	NLP	55yr

```
In [21]: emp[['Name', 'Domain', 'Age', 'Location', 'Exp']]
```

Out[21]:

	Name	Domain	Age	Location	Exp
0	Mike	Datascience#\$	34 years	Mumbai	2+
1	Teddy^	Testing	45' yr	Bangalore	<3
2	Uma#r	Dataanalyst^^#	NaN	NaN	4> yrs
3	Jane	Ana^^lytics	NaN	Hyderbad	NaN
4	Uttam*	Statistics	67-yr	NaN	5+ year
5	Kim	NLP	55yr	Delhi	10+

## Data Cleaning Or Data Cleansing

```
In [22]: emp['Name'] #regex = #$$^
```

```
Out[22]: 0      Mike
          1      Teddy^
          2      Uma#r
          3      Jane
          4      Uttam*
          5      Kim
          Name: Name, dtype: object
```

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

```
In [24]: emp['Name']
```

```
Out[24]: 0      Mike
          1      Teddy
          2      Umar
          3      Jane
          4      Uttam
          5      Kim
          Name: Name, dtype: object
```

```
In [25]: emp
```

```
Out[25]:
```

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

```
In [27]: emp['Domain']
```

```
Out[27]: 0      Datascience
          1      Testing
          2      Dataanalyst
          3      Analytics
          4      Statistics
          5      NLP
          Name: Domain, dtype: object
```

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

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

```
Out[29]: 0      34years
          1      45yr
          2      NaN
          3      NaN
          4      67yr
          5      55yr
          Name: Age, dtype: object
```

```
In [30]: emp['Age']=emp['Age'].str.extract('(\d+)')    #r'(\d+)'
```

```
In [31]: emp['Age']
```

```
Out[31]: 0      34
          1      45
          2      NaN
          3      NaN
          4      67
          5      55
          Name: Age, dtype: object
```

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

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

```
Out[33]: 0      Mumbai
          1  Bangalore
          2      NaN
          3  Hyderabad
          4      NaN
          5      Delhi
          Name: Location, dtype: object
```

```
In [34]: emp['Salary']=emp['Salary'].str.replace(r'\W', '', regex=True)
```

```
In [35]: emp['Salary']
```

```
Out[35]: 0      5000
          1     10000
          2     15000
          3     20000
          4     30000
          5     60000
          Name: Salary, dtype: object
```

```
In [36]: emp['Exp']=emp['Exp'].str.replace(r'\W', '', regex=True)
```

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

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

```
In [38]: emp['Exp']=emp['Exp'].str.extract('(\d+)')
```

```
In [39]: emp['Exp']
```

```
Out[39]: 0      2
         1      3
         2      4
         3    NaN
         4      5
         5     10
         Name: Exp, dtype: object
```

```
In [40]: clean_data = emp.copy()
```

```
In [41]: clean_data
```

```
Out[41]:
```

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

## Clean missing value treatment

```
In [42]: emp
```

```
Out[42]:
```

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

```
In [43]: clean_data
```

```
Out[43]:
```

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

```
In [44]: clean_data.isnull().sum()
```

```
Out[44]: Name      0
Domain    0
Age       2
Location  2
Salary    0
Exp       1
dtype: int64
```

```
In [45]: 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         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 [46]: import numpy as np
```

```
In [47]: clean_data['Age'] = clean_data['Age'].fillna(np.mean(pd.to_numeric(clean_data['Age']
```

```
In [48]: clean_data['Age']
```

```
Out[48]: 0      34
1      45
2    50.25
3    50.25
4      67
5      55
Name: Age, dtype: object
```

```
In [49]: clean_data['Exp'] = clean_data['Exp'].fillna(np.mean(pd.to_numeric(clean_data['Ex
```

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



```
Out[50]: 0      2
         1      3
         2      4
         3    4.8
         4      5
         5     10
         Name: Exp, dtype: object
```

```
In [51]: clean_data
```

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

```
Out[52]: 2
```

```
In [53]: clean_data['Location'] = clean_data['Location'].fillna(clean_data['Location'].mode()[0])
```

```
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	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 [55]: 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 [56]: 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 [57]: clean_data['Age'] = clean_data['Age'].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     object
5   Exp         6 non-null     object
dtypes: int32(1), object(5)
memory usage: 396.0+ bytes

```

```
In [59]: clean_data['Salary'] = clean_data['Salary'].astype(int)
```

```
In [60]: clean_data['Exp'] = clean_data['Exp'].astype(int)
```

```
In [61]: 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 [62]: clean_data['Name'] = clean_data['Name'].astype('category')
```

```
In [63]: clean_data['Domain'] = clean_data['Domain'].astype('category')
```

```
In [64]: clean_data['Location'] = clean_data['Location'].astype('category')
```

```
In [65]: 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 [66]: clean_data
```

```
Out[66]:
```

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

```
In [68]: import os
os.getcwd()
```

```
Out[68]: 'C:\\Users\\Prachi\\FSDS SENAPATI SIR\\Projects'
```

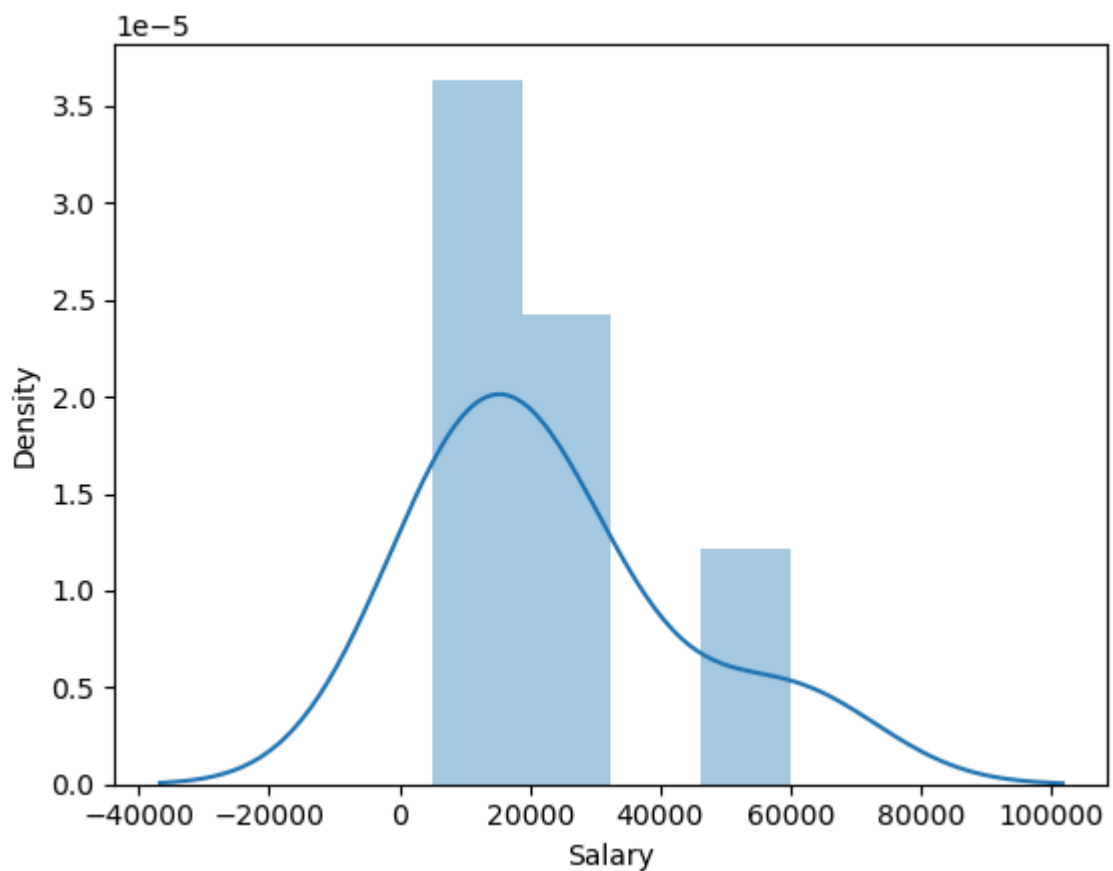
```
In [69]: import matplotlib.pyplot as plt          #visualization
import seaborn as sns
```

```
In [70]: import warnings
warnings.filterwarnings('ignore')
```

```
In [71]: clean_data['Salary']
```

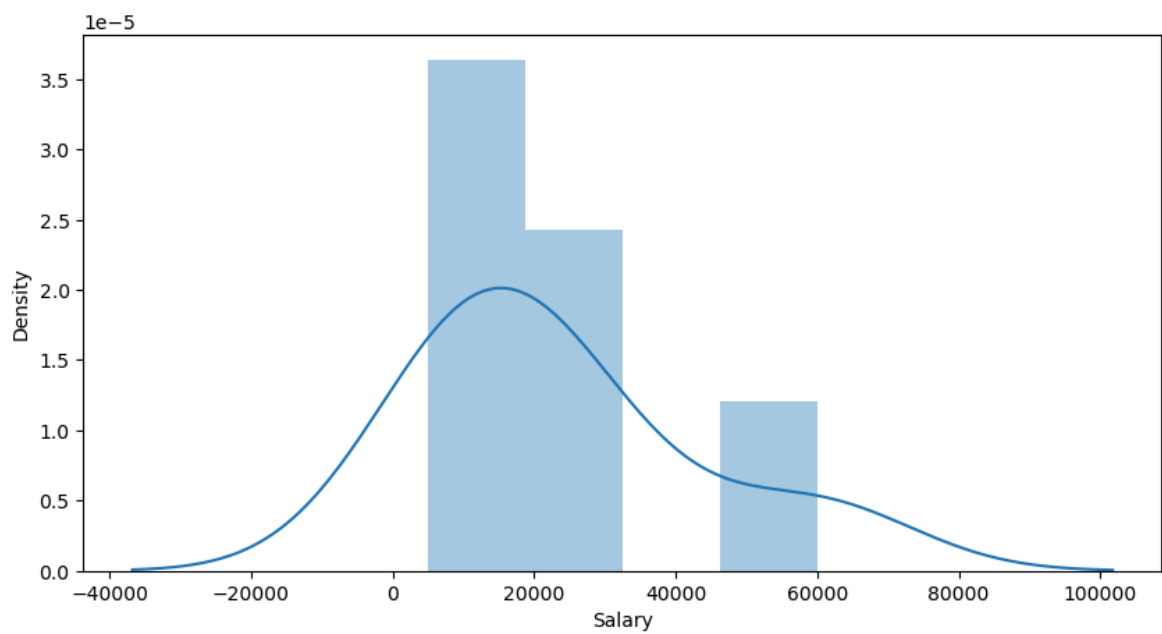
```
Out[71]: 0      5000
1     10000
2     15000
3     20000
4     30000
5     60000
Name: Salary, dtype: int32
```

```
In [72]: vis1 = sns.distplot(clean_data['Salary'])
```

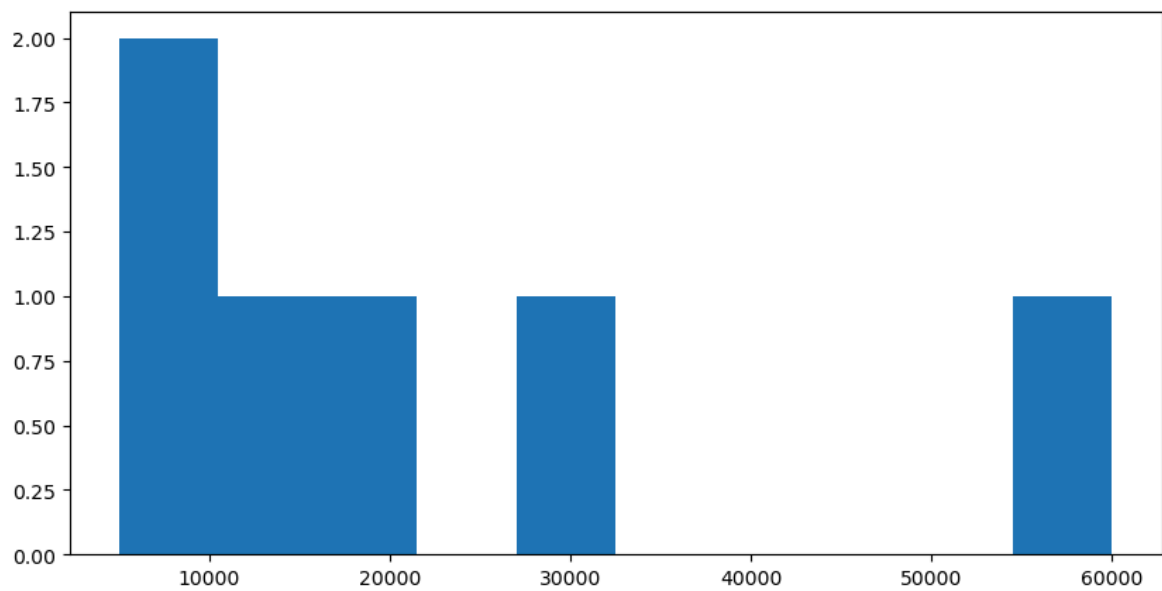


```
In [73]: plt.rcParams['figure.figsize']=10,5
```

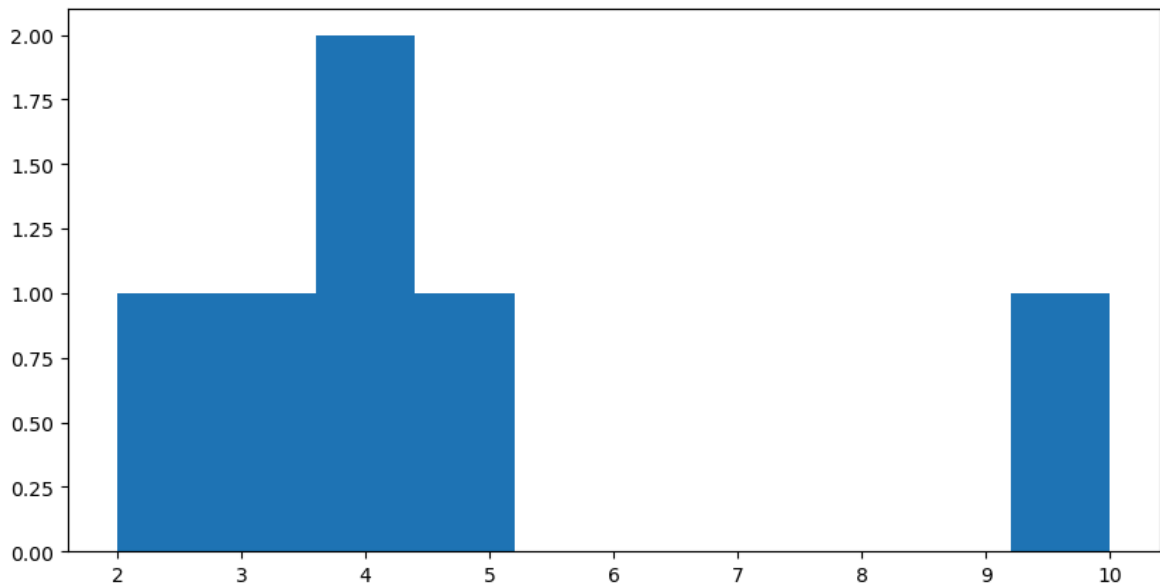
```
In [74]: vis1 = sns.distplot(clean_data['Salary'])
```



```
In [75]: vis2 = plt.hist(clean_data['Salary'])
```



```
In [76]: vis3 = plt.hist(clean_data['Exp'])
```

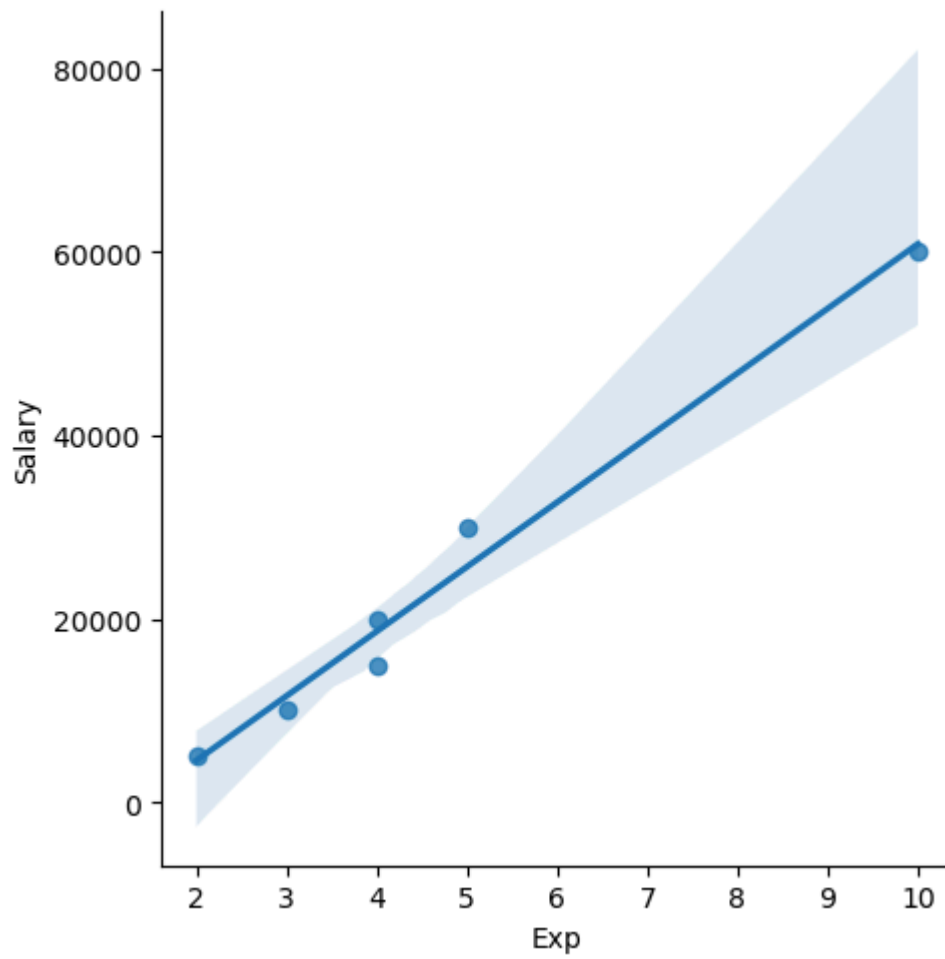


```
In [77]: clean_data
```

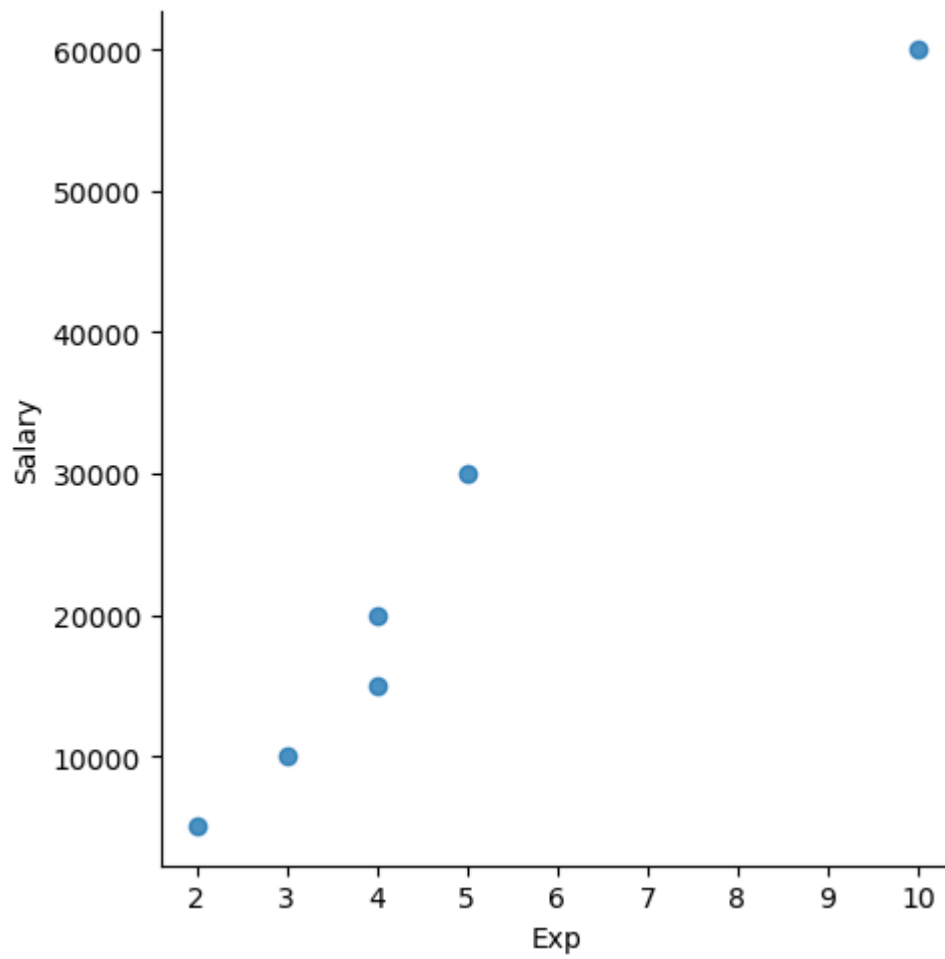
```
Out[77]:
```

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

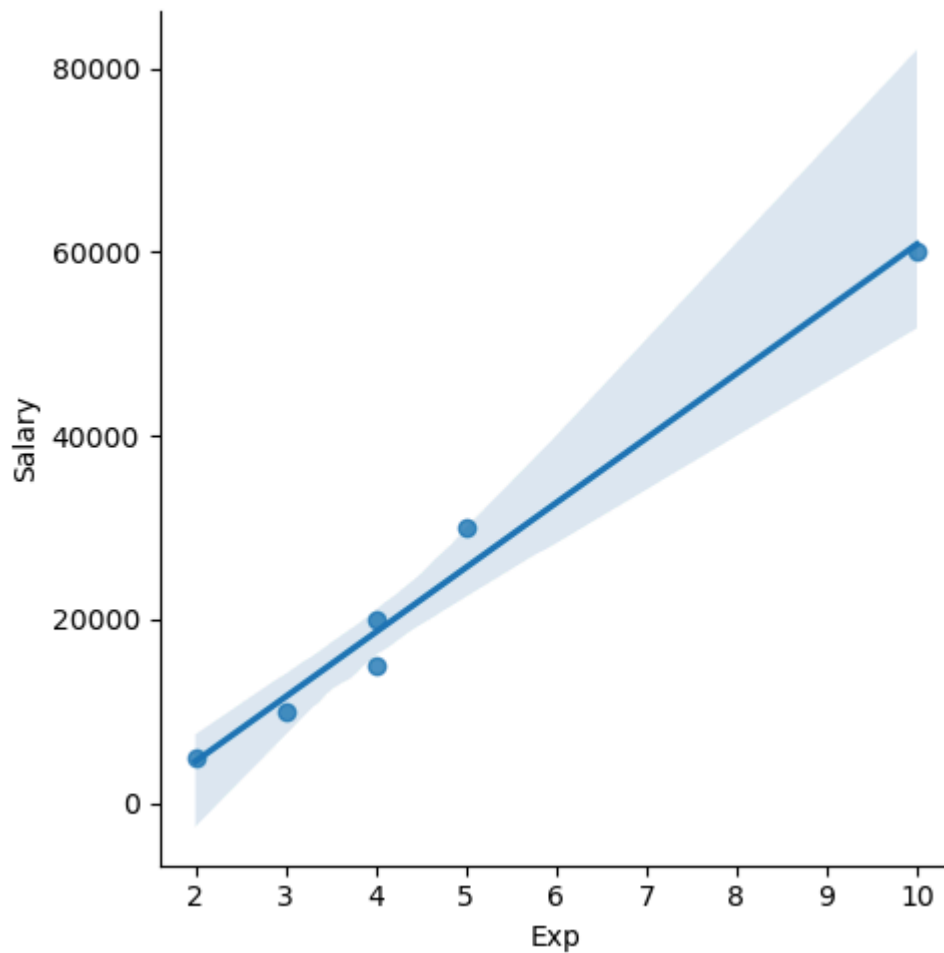


```
In [83]: vis5 = sns.lmplot(data=clean_data, x='Exp', y='Salary', fit_reg = False)
```



```
In [84]: vis6 = sns.lmplot(data=clean_data, x='Exp', y='Salary', fit_reg=True)
```





```
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	Hyderabad	20000	4
4	Uttam	Statistics	67	Bangalore	30000	5
5	Kim	NLP	55	Delhi	60000	10

```
In [86]: clean_data[:,]
```

Out[86]:

	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 [87]: `clean_data[:2]`

Out[87]:

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascience	34	Mumbai	5000	2
1	Teddy	Testing	45	Bangalore	10000	3

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 [89]: `clean_data[0:1]`

Out[89]:

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascience	34	Mumbai	5000	2

In [90]: `clean_data`

```
Out[90]:
```

	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 [91]: x_iv = clean_data.drop(['Salary'],axis=1)
```

```
In [92]: x_iv
```

```
Out[92]:
```

	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 [93]: x_iv.columns
```

```
Out[93]: Index(['Name', 'Domain', 'Age', 'Location', 'Exp'], dtype='object')
```

```
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	Hyderbad	20000	4
4	Uttam	Statistics	67	Bangalore	30000	5
5	Kim	NLP	55	Delhi	60000	10

```
In [95]: y_dv = clean_data.drop(['Name', 'Domain', 'Age', 'Location', 'Exp'],axis=1)
```

```
In [96]: y_dv
```

Out[96]: **Salary**

<b>0</b>	5000
<b>1</b>	10000
<b>2</b>	15000
<b>3</b>	20000
<b>4</b>	30000
<b>5</b>	60000

In [97]: `clean_data`

Out[97]: **Name      Domain   Age   Location   Salary   Exp**

<b>0</b>	Mike	Datascience	34	Mumbai	5000	2
<b>1</b>	Teddy	Testing	45	Bangalore	10000	3
<b>2</b>	Umar	Dataanalyst	50	Bangalore	15000	4
<b>3</b>	Jane	Analytics	50	Hyderbad	20000	4
<b>4</b>	Uttam	Statistics	67	Bangalore	30000	5
<b>5</b>	Kim	NLP	55	Delhi	60000	10

In [98]: `x_iv`

Out[98]: **Name      Domain   Age   Location   Exp**

<b>0</b>	Mike	Datascience	34	Mumbai	2
<b>1</b>	Teddy	Testing	45	Bangalore	3
<b>2</b>	Umar	Dataanalyst	50	Bangalore	4
<b>3</b>	Jane	Analytics	50	Hyderbad	4
<b>4</b>	Uttam	Statistics	67	Bangalore	5
<b>5</b>	Kim	NLP	55	Delhi	10

In [100... `y_dv`

Out[100...

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

In [101...

clean\_data

Out[101...

**Name Domain Age Location Salary Exp****0** Mike Data science 34 Mumbai 5000 2**1** Teddy Testing 45 Bangalore 10000 3**2** Umar Data analyst 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 [103...

imputation = pd.get\_dummies(clean\_data)

In [104...

imputation

Out[104...

**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