

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

```
In [2]: emp = pd.read_excel(r'C:\Users\kripal\EDA\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]: emp.shape
```

```
Out[4]: (6, 6)
```

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

```
Out[5]: 6
```

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

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

```
In [7]: 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 [8]: emp.isnull().sum()
```

```
Out[8]: Name        0
Domain      0
Age         2
Location    2
Salary      0
Exp         1
dtype: int64
```

In [9]: emp

Out[9]:

	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['Name']

Out[10]:

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

In [11]: emp['Name'] = emp['Name'].str.replace(r'\W', '', regex=True)
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	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 [12]: emp['Domain'] = emp['Domain'].str.replace(r'\W', '', regex=True)

In [13]: emp['Age'] = emp['Age'].str.replace(r'\W', '', regex=True)
emp['Location'] = emp['Location'].str.replace(r'\W', '', regex=True)
emp['Salary'] = emp['Salary'].str.replace(r'\W', '', regex=True)

In [14]: emp

Out[14]:

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

In [15]: `emp['Age'] = emp['Age'].str.extract('(\d+)') # give the output only numbers rest emp`

```
<>:1: SyntaxWarning: invalid escape sequence '\d'
<>:1: SyntaxWarning: invalid escape sequence '\d'
C:\Users\kripal\AppData\Local\Temp\ipykernel_1896\3836655718.py:1: SyntaxWarning:
invalid escape sequence '\d'
emp['Age'] = emp['Age'].str.extract('(\d+)') # give the output only numbers rest
t of item gone #extract the digits
```

Out[15]:

	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 [16]: `emp['Exp'] = emp['Exp'].str.replace(r'\W','',regex=True)
emp['Exp'] = emp['Exp'].str.extract('(\d+)')`

```
<>:2: SyntaxWarning: invalid escape sequence '\d'
<>:2: SyntaxWarning: invalid escape sequence '\d'
C:\Users\kripal\AppData\Local\Temp\ipykernel_1896\3744867905.py:2: SyntaxWarning:
invalid escape sequence '\d'
emp['Exp'] = emp['Exp'].str.extract('(\d+)')
```

In [17]: `emp`

Out[17]:

	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 [18]: `emp[['Name', 'Domain', 'Age', 'Exp']]`

Out[18]:

	Name	Domain	Age	Exp
0	Mike	Datascience	34	2
1	Teddy	Testing	45	3
2	Umar	Dataanalyst	NaN	4
3	Jane	Analytics	NaN	NaN
4	Uttam	Statistics	67	5
5	Kim	NLP	55	10

Missing Value Treatment

In [20]: `clean_data = emp.copy()`

In [21]: `clean_data`

Out[21]:

	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 [22]: `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 [23]: import numpy as np
```

```
In [24]: clean_data.head(1)
```

```
Out[24]:
```

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

```
In [25]: clean_data['Age']
```

```
Out[25]: 0      34
         1      45
         2     NaN
         3     NaN
         4      67
         5      55
         Name: Age, dtype: object
```

```
In [26]: clean_data['Age'] = clean_data['Age'].fillna(np.mean(pd.to_numeric(clean_data['A
```

```
In [27]: clean_data['Age']
```

```
Out[27]: 0      34
         1      45
         2    50.25
         3    50.25
         4      67
         5      55
         Name: Age, dtype: object
```

```
In [28]: emp
```

Out[28]:

	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 [29]: clean_data

Out[29]:

	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 [30]: clean_data['Exp'] = clean_data['Exp'].fillna(np.mean(pd.to_numeric(clean_data['E

In [31]: clean_data

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	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 [32]: clean_data['Location'] = clean_data['Location'].fillna(clean_data['Location'].mo

In [33]: clean_data

Out[33]:

	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 [34]: `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 [35]:

```
clean_data['Age'] = clean_data['Age'].astype(int)
clean_data['Salary'] = clean_data['Salary'].astype(int)
clean_data['Exp'] = clean_data['Exp'].astype(int)

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 [36]: `clean_data[['Name', 'Domain']] = clean_data[['Name', 'Domain']].astype('category')`

In [37]: `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: 704.0+ bytes
```

```
In [38]: clean_data['Location'] = clean_data['Location'].astype('category') #capital mist
```

```
In [39]: 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 [40]: clean_data['Location'] = clean_data['Location'].astype('category')
```

```
In [41]: clean_data = clean_data.drop('location',axis=1)
```



```

-----
KeyError                                Traceback (most recent call last)
Cell In[41], line 1
----> 1 clean_data = clean_data.drop('location',axis=1)

File ~\anaconda3\Lib\site-packages\pandas\core\frame.py:5581, in DataFrame.drop(self, labels, axis, index, columns, level, inplace, errors)
    5433 def drop(
    5434     self,
    5435     labels: IndexLabel | None = None,
    (... )
    5442     errors: IgnoreRaise = "raise",
    5443 ) -> DataFrame | None:
    5444     """
    5445     Drop specified labels from rows or columns.
    5446     (...)
    5579         weight  1.0      0.8
    5580     """
-> 5581     return super().drop(
    5582         labels=labels,
    5583         axis=axis,
    5584         index=index,
    5585         columns=columns,
    5586         level=level,
    5587         inplace=inplace,
    5588         errors=errors,
    5589     )

File ~\anaconda3\Lib\site-packages\pandas\core\generic.py:4788, in NDFrame.drop(self, labels, axis, index, columns, level, inplace, errors)
    4786 for axis, labels in axes.items():
    4787     if labels is not None:
-> 4788         obj = obj._drop_axis(labels, axis, level=level, errors=errors)
    4790 if inplace:
    4791     self._update_inplace(obj)

File ~\anaconda3\Lib\site-packages\pandas\core\generic.py:4830, in NDFrame._drop_axis(self, labels, axis, level, errors, only_slice)
    4828     new_axis = axis.drop(labels, level=level, errors=errors)
    4829     else:
-> 4830     new_axis = axis.drop(labels, errors=errors)
    4831     indexer = axis.get_indexer(new_axis)
    4833 # Case for non-unique axis
    4834 else:

File ~\anaconda3\Lib\site-packages\pandas\core\indexes\base.py:7070, in Index.drop(self, labels, errors)
    7068 if mask.any():
    7069     if errors != "ignore":
-> 7070         raise KeyError(f"{labels[mask].tolist()} not found in axis")
    7071     indexer = indexer[~mask]
    7072     return self.delete(indexer)

KeyError: "[ 'location' ] not found in axis"

```

In [101... 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 [103... `clean_data`

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	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 [105... `clean_data.to_csv('clean_data.csv')`

In [107... `import os`
`os.getcwd()`

Out[107... `'C:\\Users\\kripal'`

In [108... `clean_data.columns`

Out[108... `Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')`

In [109... `import matplotlib.pyplot as plt`
`%matplotlib inline`
`import seaborn as sns`

In [111... `import warnings`
`warnings.filterwarnings('ignore')`

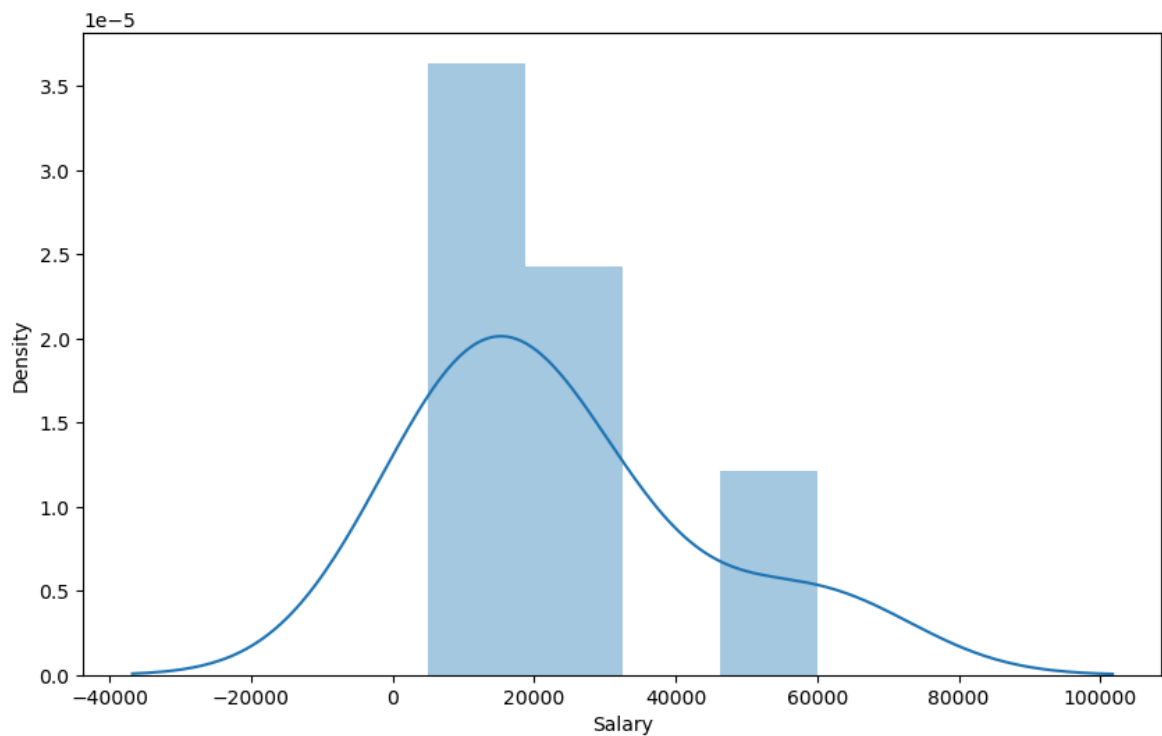
In [113... `clean_data['Salary']`

Out[113...

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

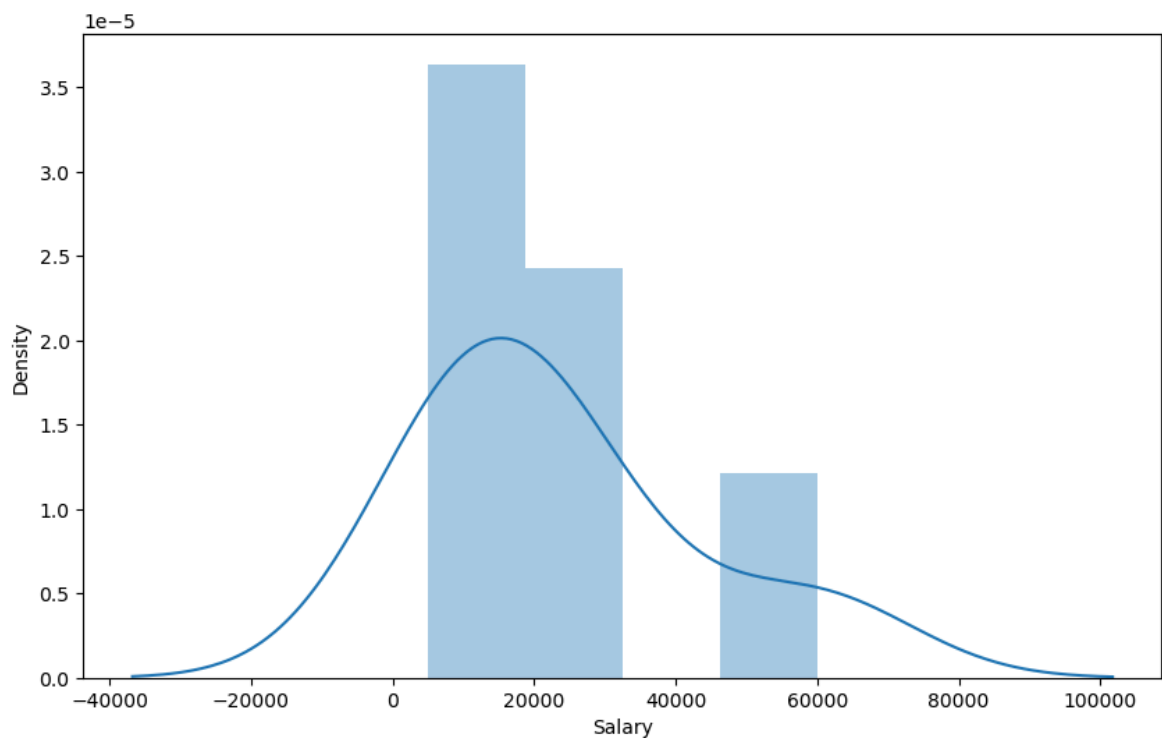
Name: Salary, dtype: int32

```
In [115... vis1 = sns.distplot(clean_data['Salary'])
```

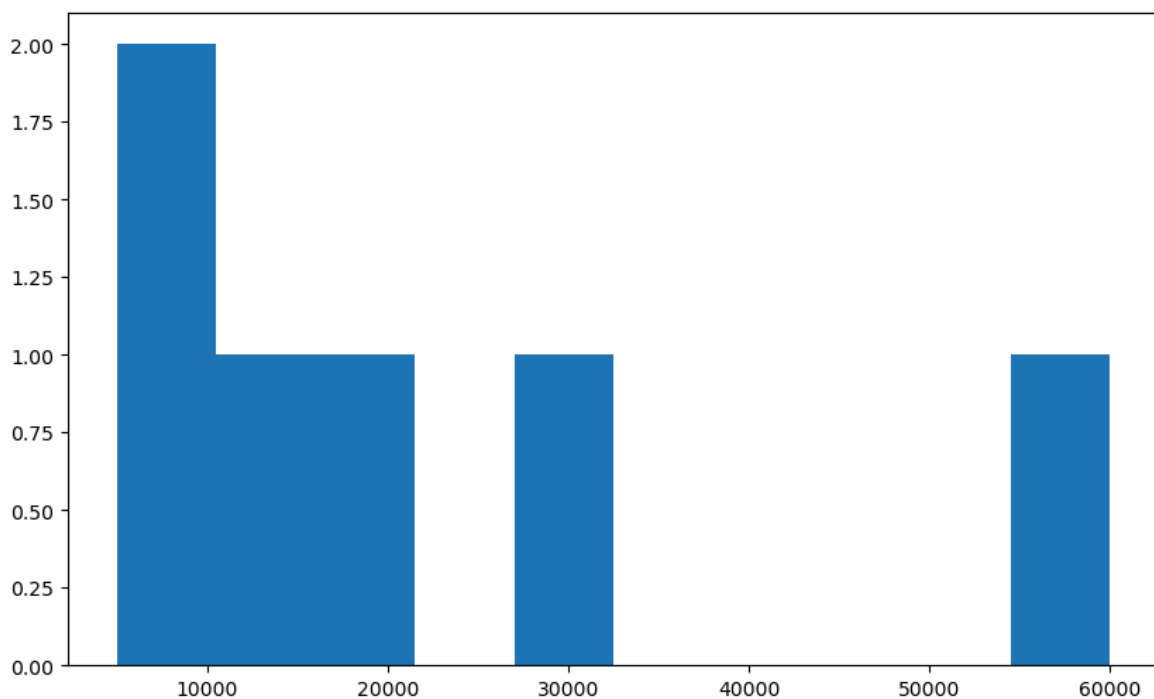


```
In [116... plt.rcParams['figure.figsize']=10,6
```

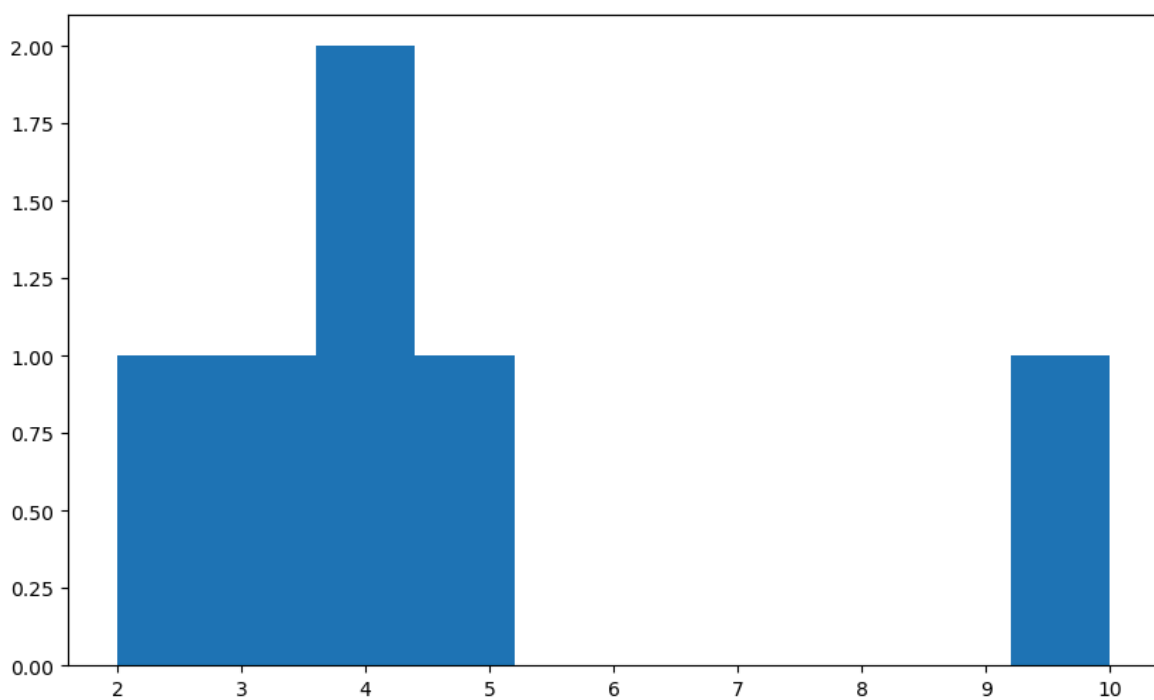
```
In [118... vis1 = sns.distplot(clean_data['Salary'])
```



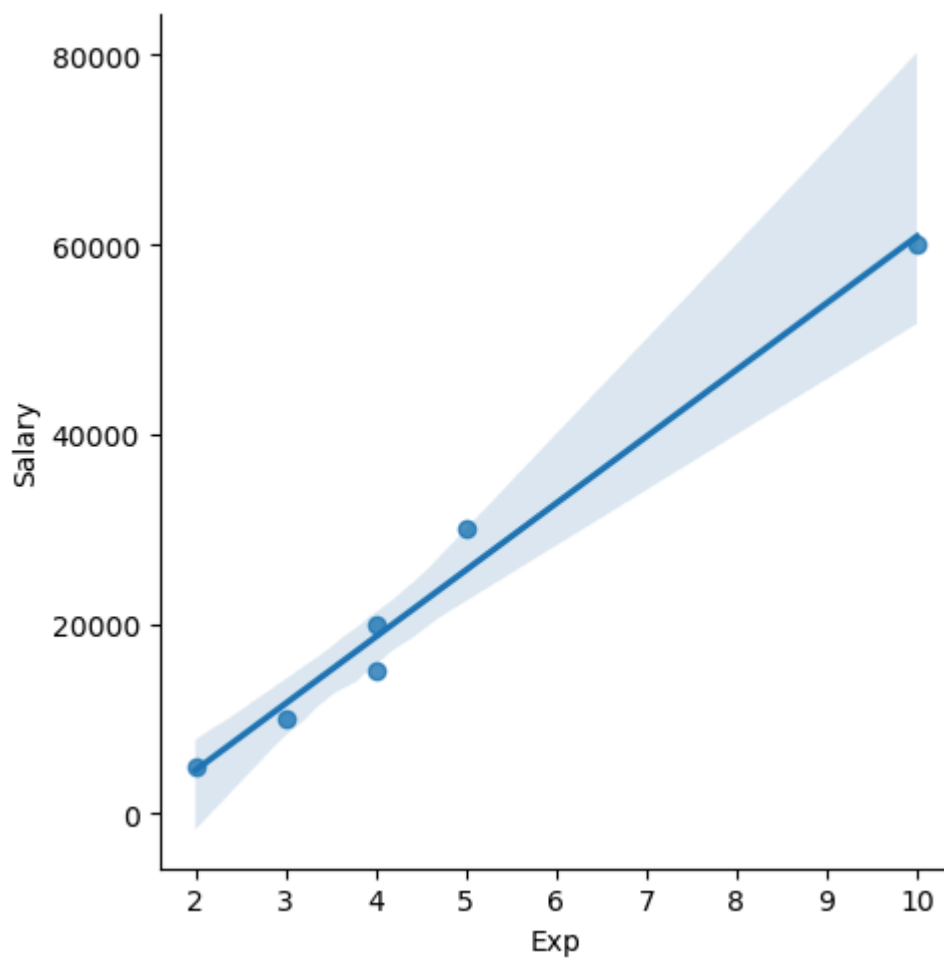
```
In [119... vis2 = plt.hist(clean_data['Salary'])
```



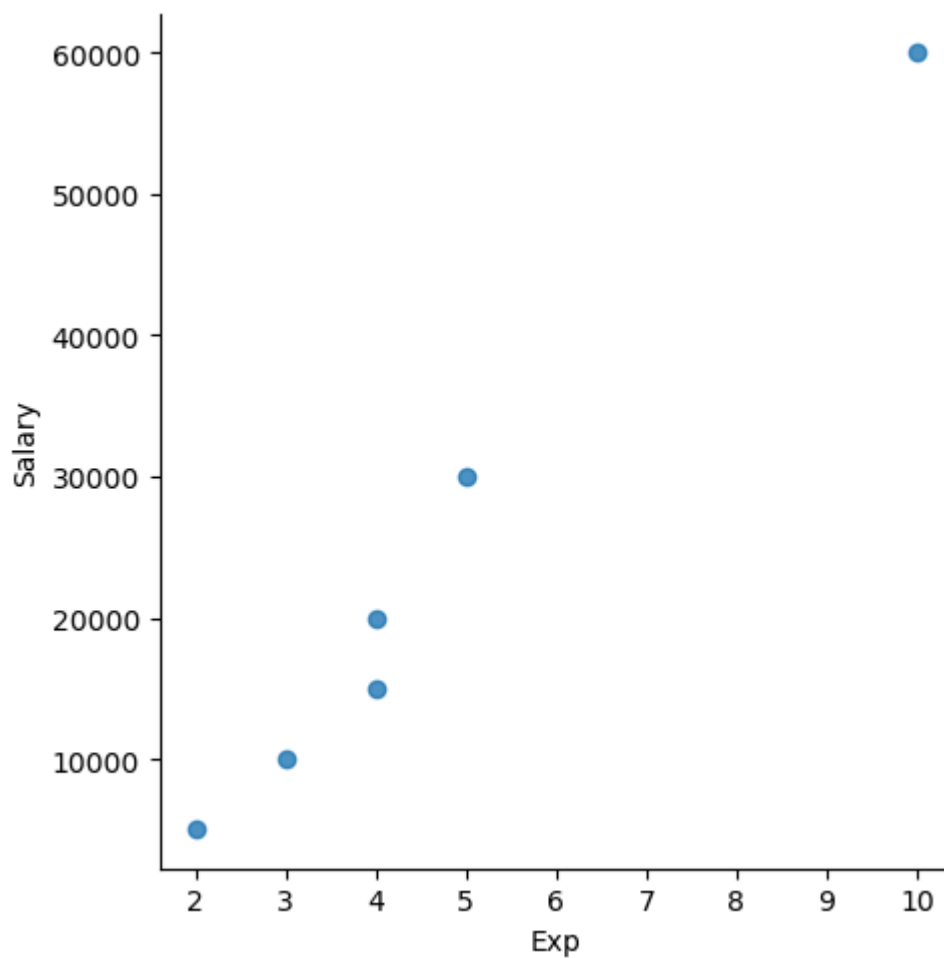
```
In [121...] vis3 = plt.hist(clean_data['Exp'])
```



```
In [122...] vis4 = sns.lmplot(data=clean_data, x='Exp', y = 'Salary')
```

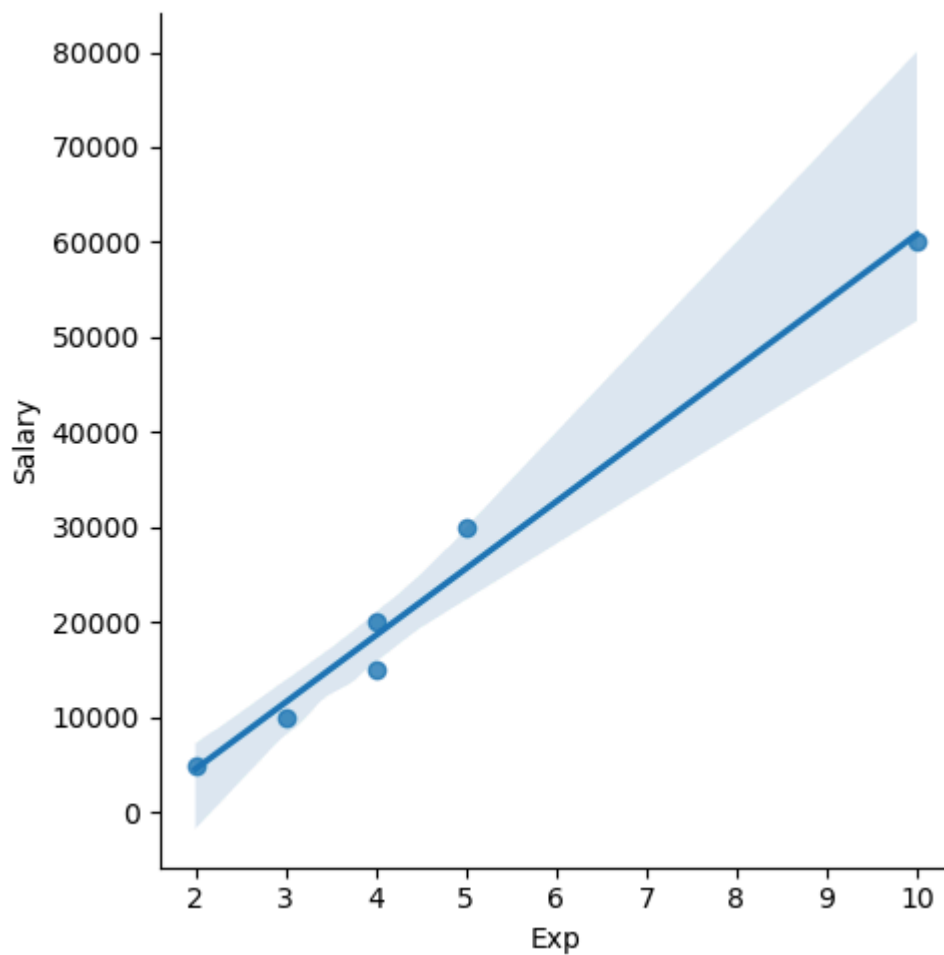


```
In [123... vis5 = sns.lmplot(data = clean_data, x = 'Exp', y = 'Salary', fit_reg=False)
```



In [124...

```
vis5 = sns.lmplot(data = clean_data, x = 'Exp', y = 'Salary', fit_reg=True)
```



In [125... `clean_data[:]`

Out[125...

	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 [126... `clean_data[:2]`

Out[126...

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

In [127... `clean_data`

Out[127...

	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 [128...

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

In [129...

```
x_iv
```

Out[129...

	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 [130...

```
x_iv.columns
```

Out[130...

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

In [131...

```
clean_data.columns
```

Out[131...

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

In [136...

```
x_dv = clean_data.drop(['Name', 'Domain', 'Age', 'Location', 'Exp'],axis=1)
```

In [138...

```
x_dv
```

Out[138...

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

In [139...

```
clean_data
```


Out[139...

	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 [140...

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

In [141...

```
imputation
```

Out[141...

	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 []: