

In [3]: `import pandas as pd`

In [5]: `raw_data=pd.read_excel(r'C:\Users\user\Documents\Rawdata.xlsx')`
`raw_data`

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 [7]: `id(raw_data)`

Out[7]: 2258201618016

In [9]: `raw_data.columns`

Out[9]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')

In [11]: `raw_data.shape`

Out[11]: (6, 6)

In [13]: `raw_data.head()`

Out[13]:

	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 [15]: `raw_data.tail()`

Out[15]:

	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 [17]: `raw_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 [19]: `raw_data.isnull()`

Out[19]:

	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]: `raw_data.isna() # same as isnull()`

Out[21]:

	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 [23]: raw_data.isnull().sum()
```

```
Out[23]: Name      0
        Domain    0
        Age       2
        Location   2
        Salary     0
        Exp       1
        dtype: int64
```

```
In [25]: raw_data
```

```
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	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+

DATA CLEANING AND DATA CLEASING

```
In [28]: raw_data['Name']=raw_data['Name'].str.replace(r'\W','',regex=True) # \W capture
raw_data['Name'] # '' This is the replacement string. Here, it's an empty string
```

```
Out[28]: 0      Mike
        1      Teddy
        2      Umar
        3      Jane
        4      Uttam
        5      Kim
        Name: Name, dtype: object
```

```
In [30]: raw_data['Domain']
```

```
Out[30]: 0      Datascience#$
        1      Testing
        2      Dataanalyst^^#
        3      Ana^^lytics
        4      Statistics
        5      NLP
        Name: Domain, dtype: object
```

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

```
In [34]: raw_data['Domain']
```

```
Out[34]: 0    Datascience
         1      Testing
         2    Dataanalyst
         3      Analytics
         4    Statistics
         5         NLP
         Name: Domain, dtype: object
```

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

```
Out[36]: 0    34years
         1     45yr
         2      NaN
         3      NaN
         4     67yr
         5     55yr
         Name: Age, dtype: object
```

```
In [68]: raw_data['Age']=raw_data['Age'].str.extract('(\d+)') # \d+: \d matches any digit
         raw_data['Age']
```

```
<>:1: SyntaxWarning: invalid escape sequence '\d'
<>:1: SyntaxWarning: invalid escape sequence '\d'
C:\Users\user\AppData\Local\Temp\ipykernel_1852\2837976887.py:1: SyntaxWarning: i
nvalid escape sequence '\d'
    raw_data['Age']=raw_data['Age'].str.extract('(\d+)') # \d+: \d matches any digi
t (0-9), and the + means "one or more" digits.
```

```
Out[68]: 0     34
         1     45
         2    NaN
         3    NaN
         4     67
         5     55
         Name: Age, dtype: object
```

```
In [70]: raw_data
```

```
Out[70]:
```

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

Out[72]:

	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 [74]: raw_data['Salary']=raw_data['Salary'].str.replace(r'\W', '', regex=True)
raw_data
```

Out[74]:

	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 [76]: raw_data['Exp']=raw_data['Exp'].str.extract('(\d+)')
raw_data
```

```
<>:1: SyntaxWarning: invalid escape sequence '\d'
<>:1: SyntaxWarning: invalid escape sequence '\d'
C:\Users\user\AppData\Local\Temp\ipykernel_1852\206955048.py:1: SyntaxWarning: in
valid escape sequence '\d'
raw_data['Exp']=raw_data['Exp'].str.extract('(\d+)')
```

Out[76]:

	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 [78]: clean_data=raw_data.copy()
```

```
In [80]: clean_data
```

Out[80]:

	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

EDA Techniques

In [83]: `clean_data.isnull().sum()`

Out[83]:

Name	0
Domain	0
Age	2
Location	2
Salary	0
Exp	1

dtype: int64

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

Out[85]:

0	34
1	45
2	NaN
3	NaN
4	67
5	55

Name: Age, dtype: object

In [87]: `import numpy as np`

In [89]: `clean_data['Age']=clean_data['Age'].fillna(np.mean(pd.to_numeric(clean_data['Age'], errors='coerce')))`

Out[89]:

0	34
1	45
2	50.25
3	50.25
4	67
5	55

Name: Age, dtype: object

In [91]: `clean_data['Exp']=clean_data['Exp'].fillna(np.mean(pd.to_numeric(clean_data['Exp'], errors='coerce')))`

```
Out[91]: 0      2
         1      3
         2      4
         3    4.8
         4      5
         5     10
         Name: Exp, dtype: object
```

```
In [93]: clean_data['Location'].isnull().sum()
```

```
Out[93]: 2
```

```
In [105... clean_data['Location']=clean_data['Location'].fillna(clean_data['Location'].mode
clean_data['Location'] # The [0] is used to select the first mode from the Serie
```

```
Out[105... 0      Mumbai
         1    Bangalore
         2    Hyderabad
         3    Hyderabad
         4    Hyderabad
         5        Delhi
         Name: Location, dtype: object
```

```
In [122... clean_data
```

```
Out[122...   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 [124... 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 [129... clean_data['Age']=clean_data['Age'].astype(int)
clean_data['Age']
```

```
Out[129...] 0    34
            1    45
            2    50
            3    50
            4    67
            5    55
            Name: Age, dtype: int32
```

```
In [131...] clean_data['Salary']=clean_data['Salary'].astype(int)
            clean_data['Salary']
```

```
Out[131...] 0    5000
            1   10000
            2   15000
            3   20000
            4   30000
            5   60000
            Name: Salary, dtype: int32
```

```
In [133...] clean_data['Exp']=clean_data['Exp'].astype(int)
            clean_data['Exp']
```

```
Out[133...] 0     2
            1     3
            2     4
            3     4
            4     5
            5    10
            Name: Exp, dtype: int32
```

```
In [135...] 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 [143...] clean_data['Name']=clean_data['Name'].astype('category')
            clean_data['Location']=clean_data['Location'].astype('category')
            clean_data['Domain']=clean_data['Domain'].astype('category')
            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 [146... `clean_data`

Out[146...

	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 [148... `import os`
`os.getcwd()`

Out[148... `'C:\\Users\\user'`

In [150... `import matplotlib.pyplot as plt # visualization`
`import seaborn as sns`

In [152... `import warnings`
`warnings.filterwarnings('ignore')`

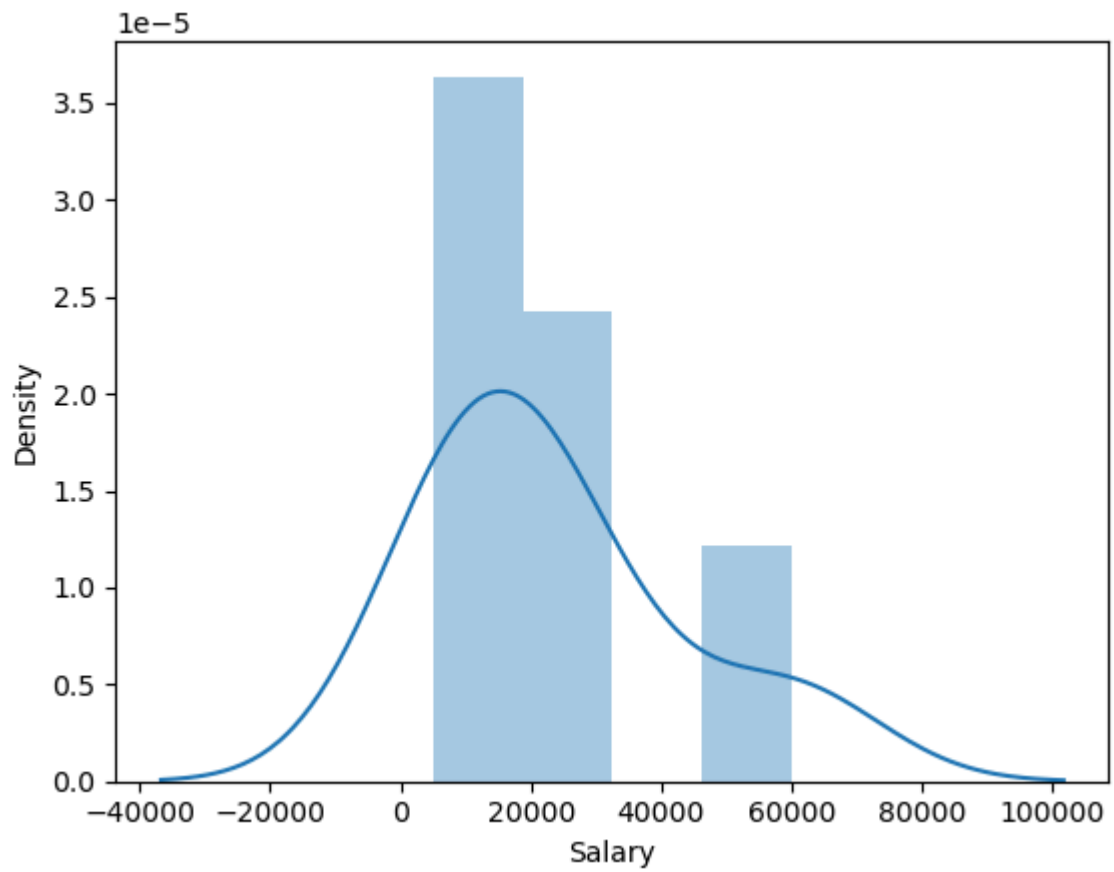
In [154... `clean_data['Salary']`

Out[154...

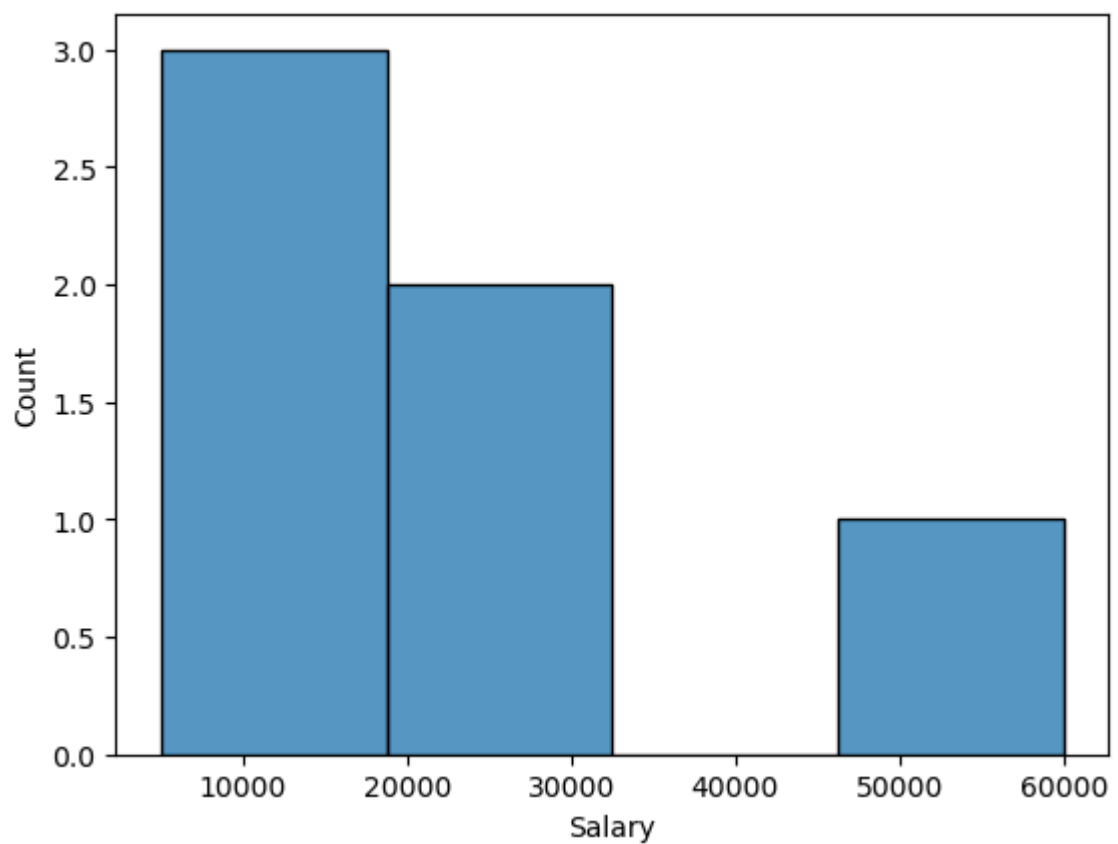
0	5000
1	10000
2	15000
3	20000
4	30000
5	60000

Name: Salary, dtype: int32

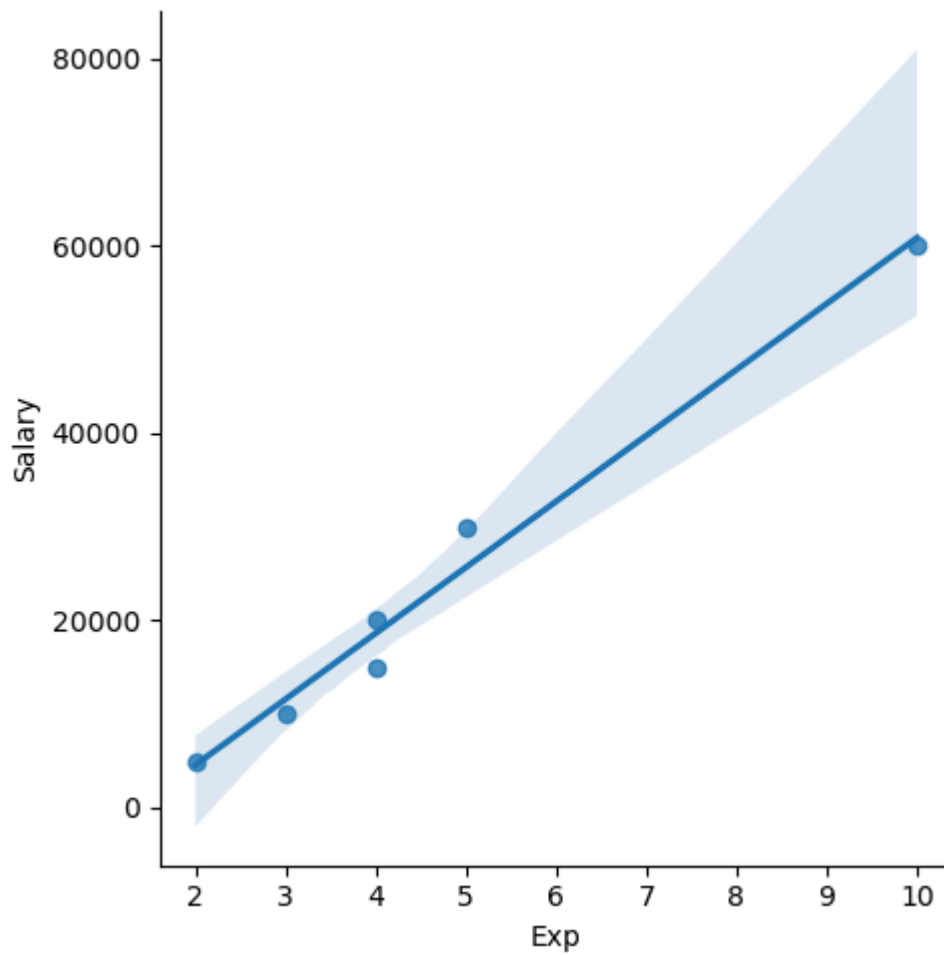
In [156... `vis1=sns.distplot(clean_data['Salary'])`



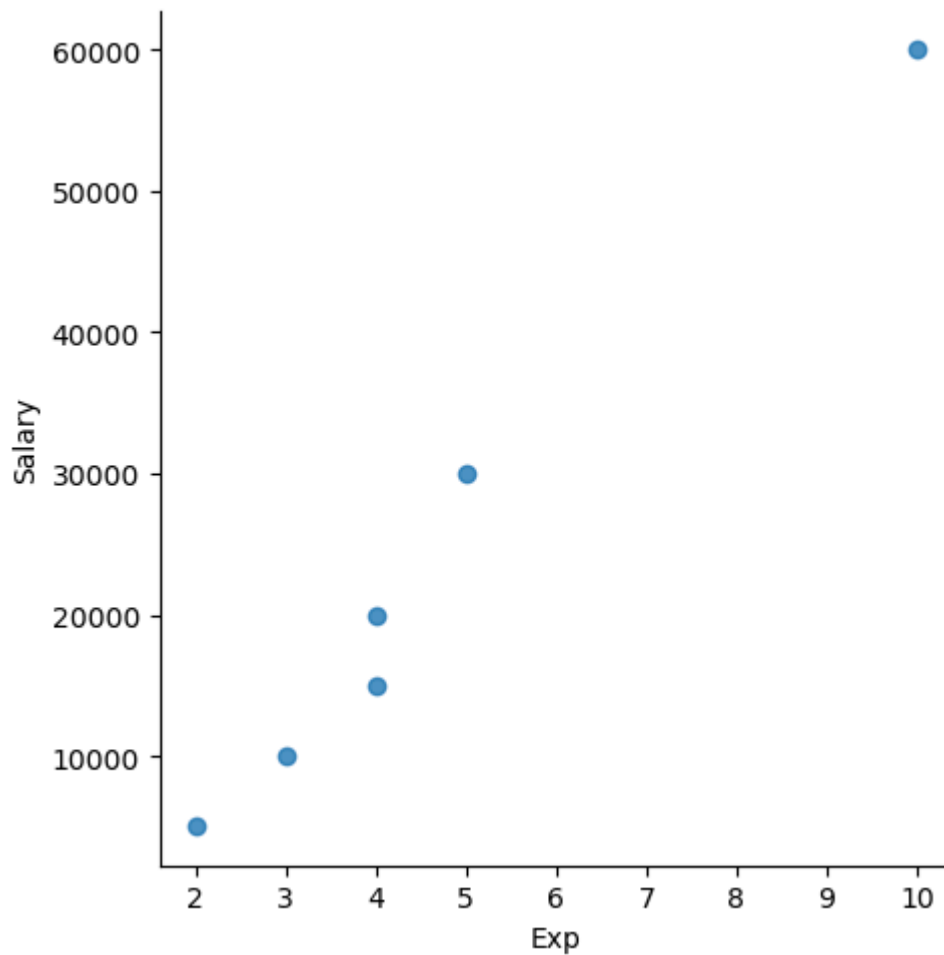
```
In [158... vis2=sns.histplot(clean_data['Salary'])
```



```
In [160... vis3=sns.lmplot(data=clean_data,x='Exp',y='Salary')
```



```
In [162... vis4=sns.lmplot(data=clean_data,x='Exp',y='Salary',fit_reg=False)
```



```
In [166... clean_data[:]
```

```
Out[166... 
```

	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 [168... clean_data[0:6:2]
```

```
Out[168... 
```

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

Out[170...

	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 [172...

```
clean_data.columns
```

Out[172...

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

In [176...

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

Out[176...

	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 [180...

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

Out[180...

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

In [182...

```
raw_data
```

Out[182...

	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 [184...

clean_data

Out[184...

	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 [186...

x_iv

Out[186...

	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 [188...

y_dv

Out[188...

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

In [192...

```
imputation=pd.get_dummies(clean_data) # variable creation
imputation
```

Out[192...

	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 []: