

```
In [3]: import pandas as pd
```

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

```
Out[4]: '2.1.4'
```

```
In [5]: emp = pd.read_excel(r'C:\Users\ABHILASH REDDY\Downloads\EDA.xlsx')
```

```
In [6]: emp
```

```
Out[6]:
```

	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]: emp.columns
```

```
Out[7]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')
```

```
In [8]: emp.shape
```

```
Out[8]: (6, 6)
```

```
In [9]: emp.head()
```

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

In [10]: `emp.tail()`

Out[10]:

	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 [11]: `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 [12]: emp

Out[12]:

	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 [13]: emp.isnull()

```
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]: # data cleaning
```

```
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)
```

```
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
```

```
Out[19]:
```

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

```
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
```

Out[28]:

	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 [29]: `emp['Salary']`

Out[29]:

```

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

```

In [30]: `emp['Salary'] = emp['Salary'].str.replace(r'\W','',regex=True)`

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

Out[31]:

```

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

```

In [32]: `emp`

Out[32]:

	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 [33]: 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()`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

## EDA TECHNIQUES

In [38]: `clean_data`

```
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.isnull().sum()
```

```
Out[39]: Name      0
Domain    0
Age       2
Location  2
Salary    0
Exp       1
dtype: int64
```

```
In [40]: clean_data['Age']
```

```
Out[40]: 0      34
1      45
2     NaN
3     NaN
4      67
5      55
Name: Age, dtype: object
```

```
In [41]: import numpy as np
```

```
In [42]: clean_data['Age'] = clean_data['Age'].fillna(np.mean(pd.to_numeric(clean_data['Age'])))
```

```
In [43]: clean_data['Age']
```

```
Out[43]: 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['Exp'])))
```

```
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	Hyderabad	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']=clean_data['Location'].fillna(clean_data['Location'].mode()[0])
```

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

```
Out[50]: 0      Mumbai
1    Bangalore
2    Bangalore
3     Hyderabad
4    Bangalore
5         Delhi
Name: Location, 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	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 [52]: clean_data['Age'] = clean_data['Age'].astype(int)
clean_data['Salary'] = clean_data['Salary'].astype(int)
clean_data['Exp'] = clean_data['Exp'].astype(int)
```

```
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      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 [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      int32
 5   Exp         6 non-null      int32
dtypes: int32(3), object(3)
memory usage: 348.0+ bytes
```

```
In [56]: 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 [57]: clean_data['Salary'] = clean_data['Salary'].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 [60]: clean_data['Name'] = clean_data['Name'].astype('category')
clean_data['Domain'] = clean_data['Domain'].astype('category')
clean_data['Location'] = clean_data['Location'].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 [59]: clean_data.to_csv('clean_data.csv')
```

```
In [61]: import os
os.getcwd()
```

```
Out[61]: 'C:\\Users\\ABHILASH REDDY'
```

```
In [62]: clean_data.columns
```

```
Out[62]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')
```

```
In [63]: import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [64]: import warnings
warnings.filterwarnings('ignore')
```

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

```
Out[65]:
```

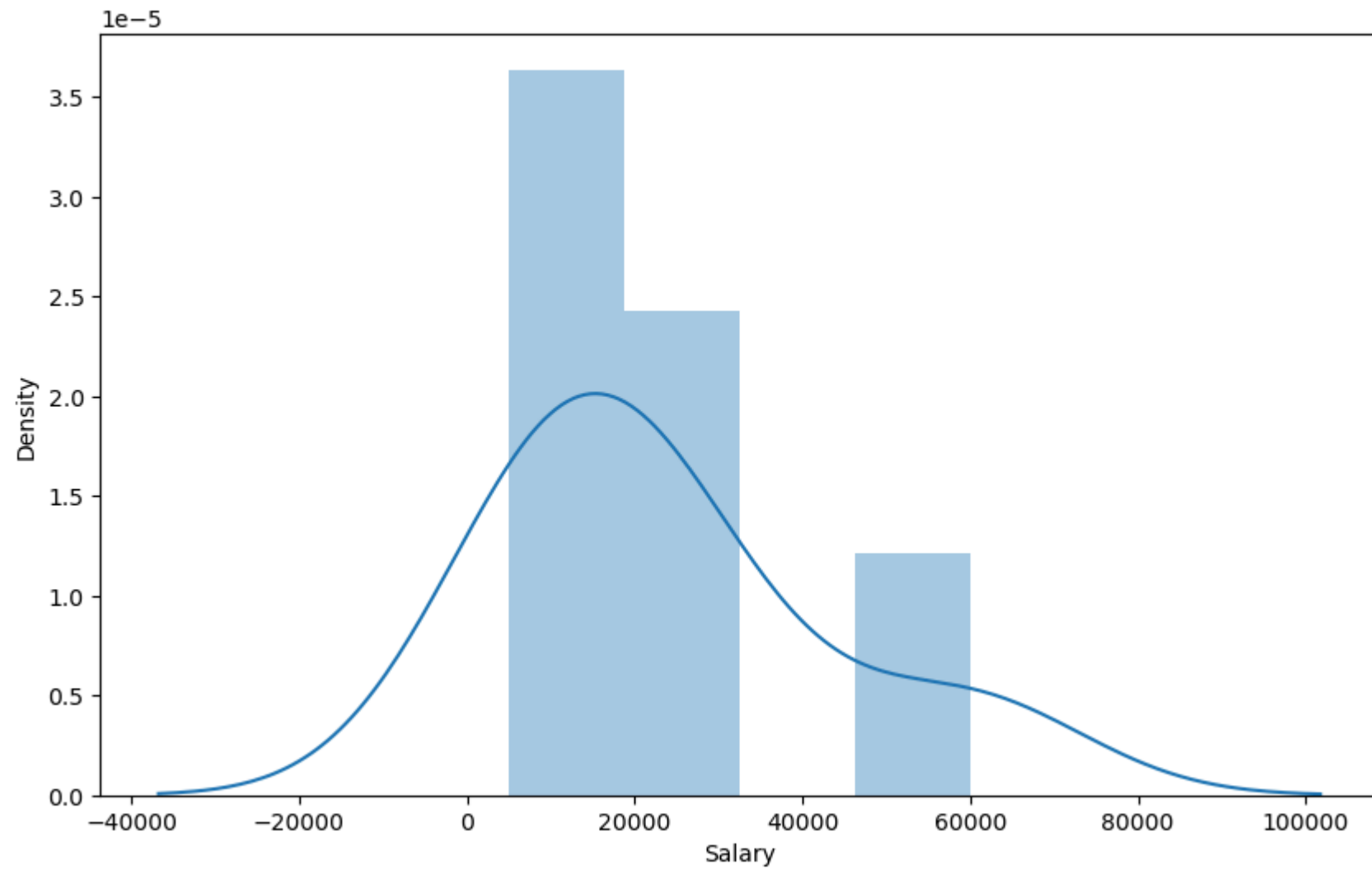
	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascience	34	Mumbai	5000	2
1	Teddy	Testing	45	Bangalore	10000	3
2	Umar	Dataanalyst	50	Bangalore	15000	4
3	Jane	Analytics	50	Hyderabad	20000	4
4	Uttam	Statistics	67	Bangalore	30000	5
5	Kim	NLP	55	Delhi	60000	10

```
In [66]: clean_data['Salary']
```

```
Out[66]: 0    5000
1   10000
2   15000
3   20000
4   30000
5   60000
Name: Salary, dtype: int32
```

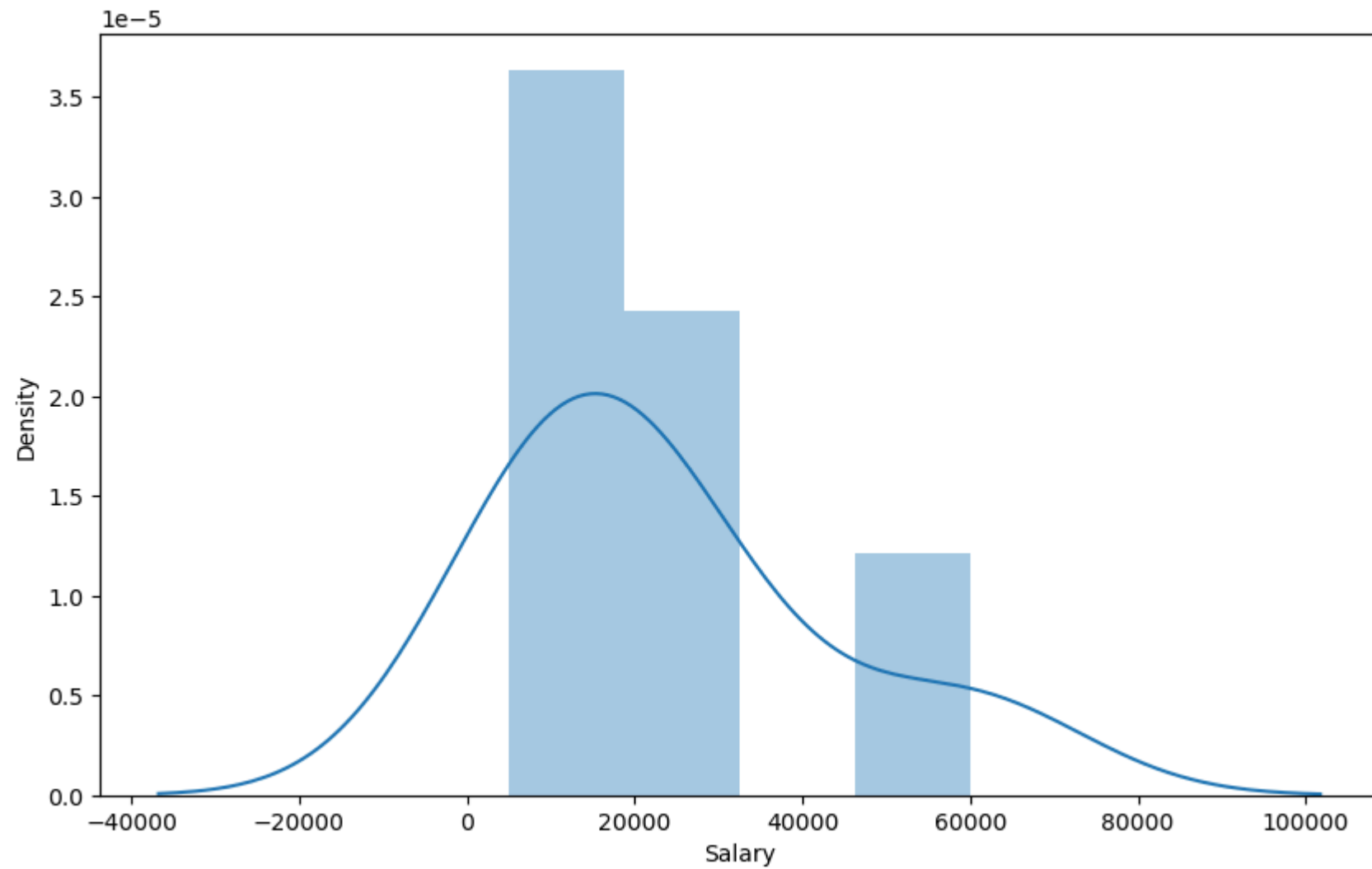
```
In [69]: visualization1 =sns.distplot(clean_data['Salary'])
```



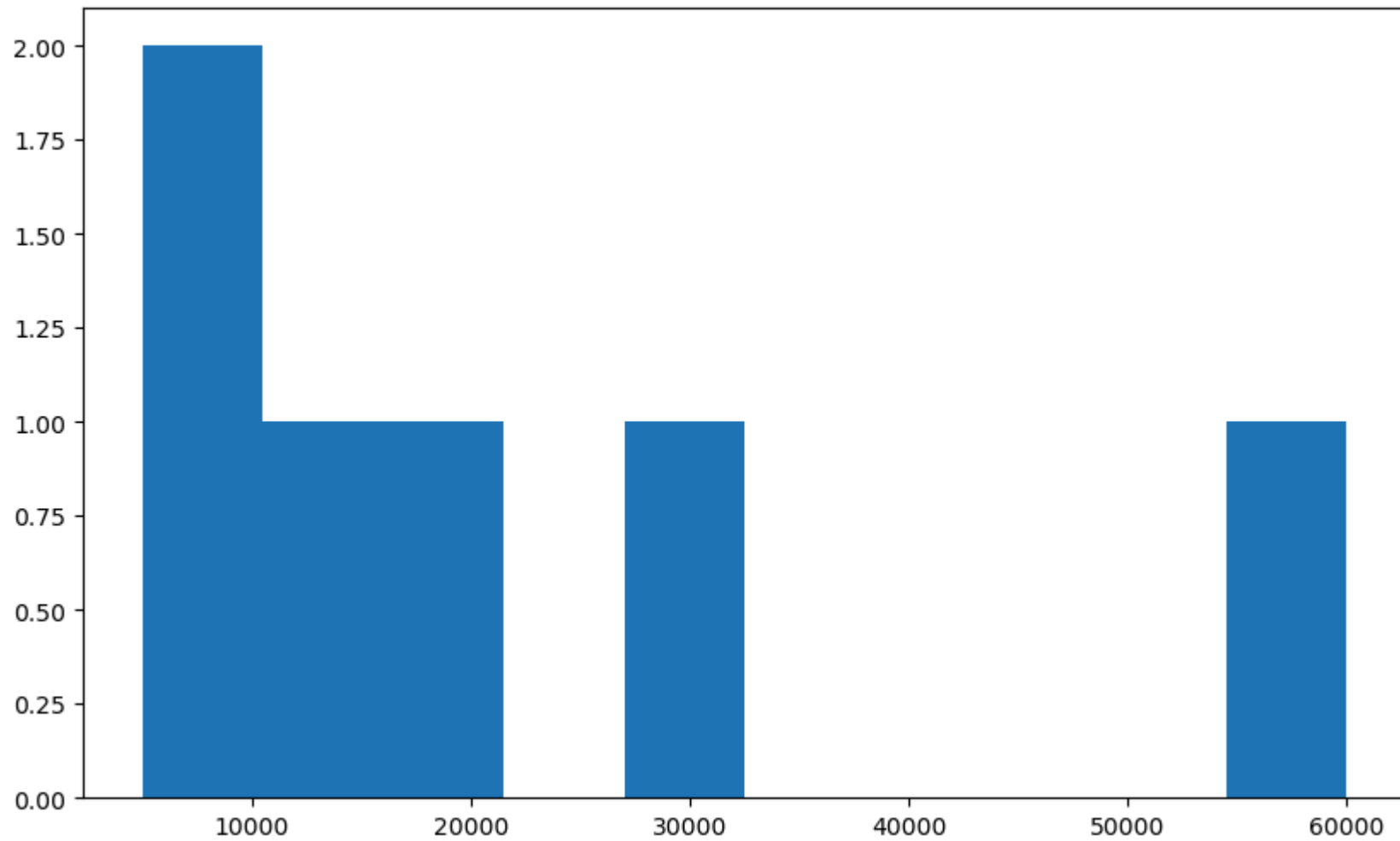


```
In [68]: plt.rcParams['figure.figsize'] = 10,6
```

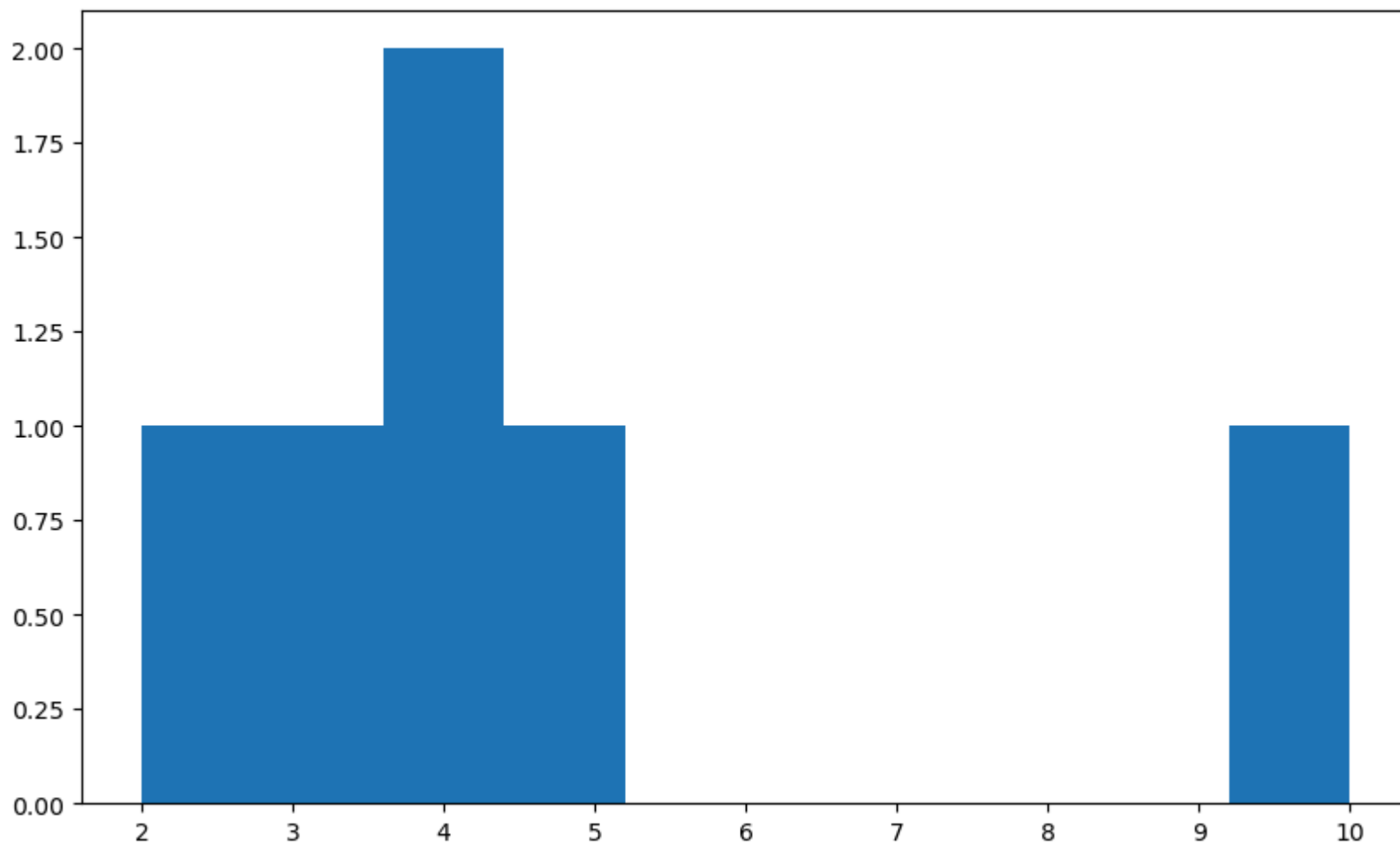
```
In [70]: visualization1 = sns.distplot(clean_data['Salary'])
```



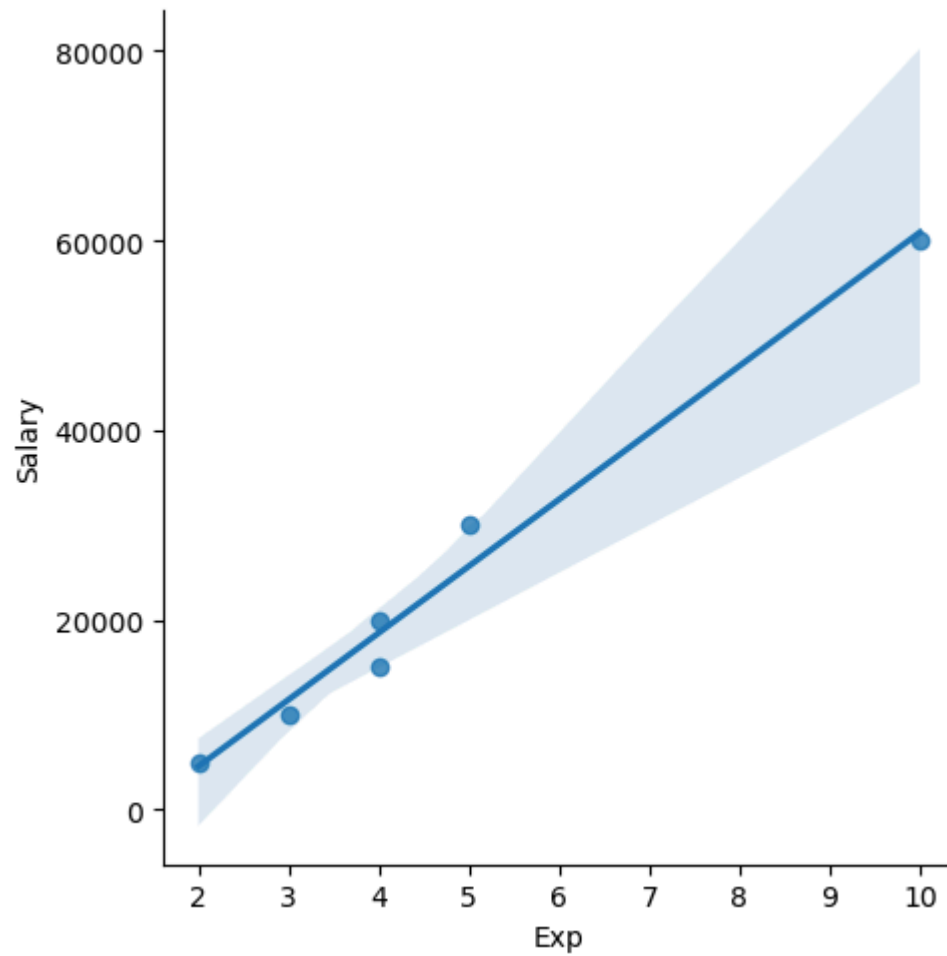
```
In [73]: visualization2 = plt.hist(clean_data['Salary'])
```



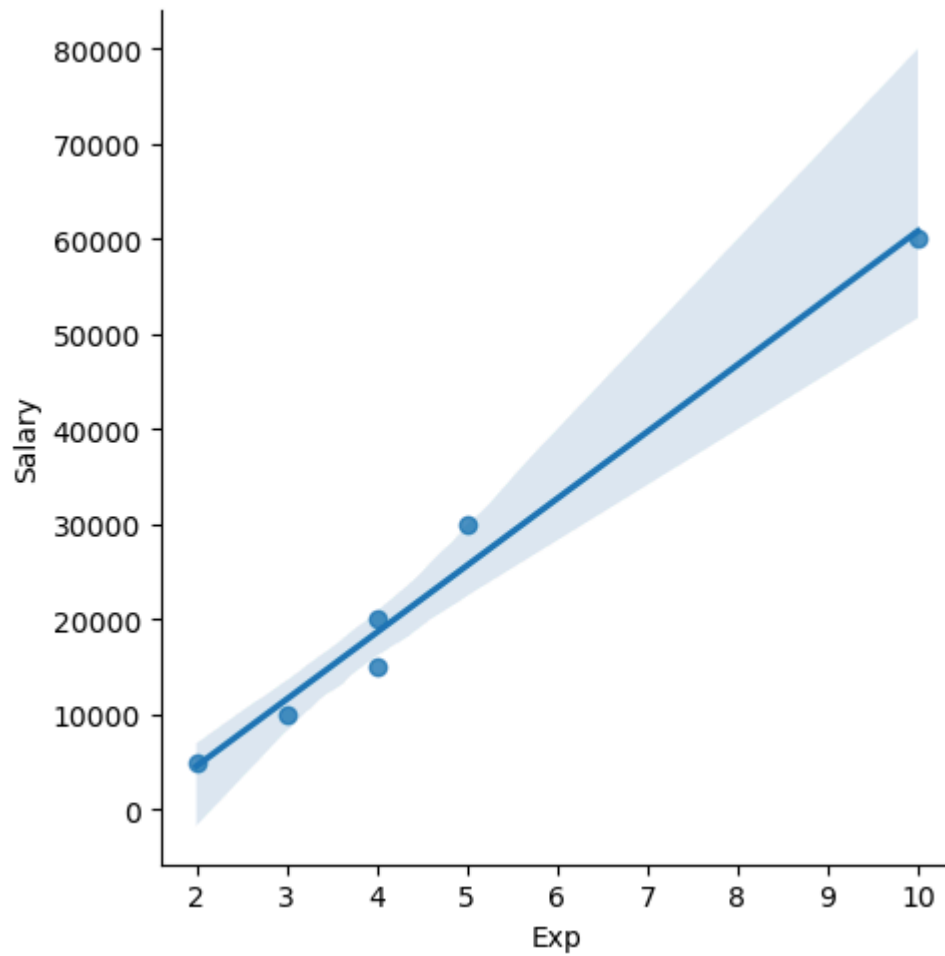
```
In [75]: visualization3 = plt.hist(clean_data['Exp'])
```



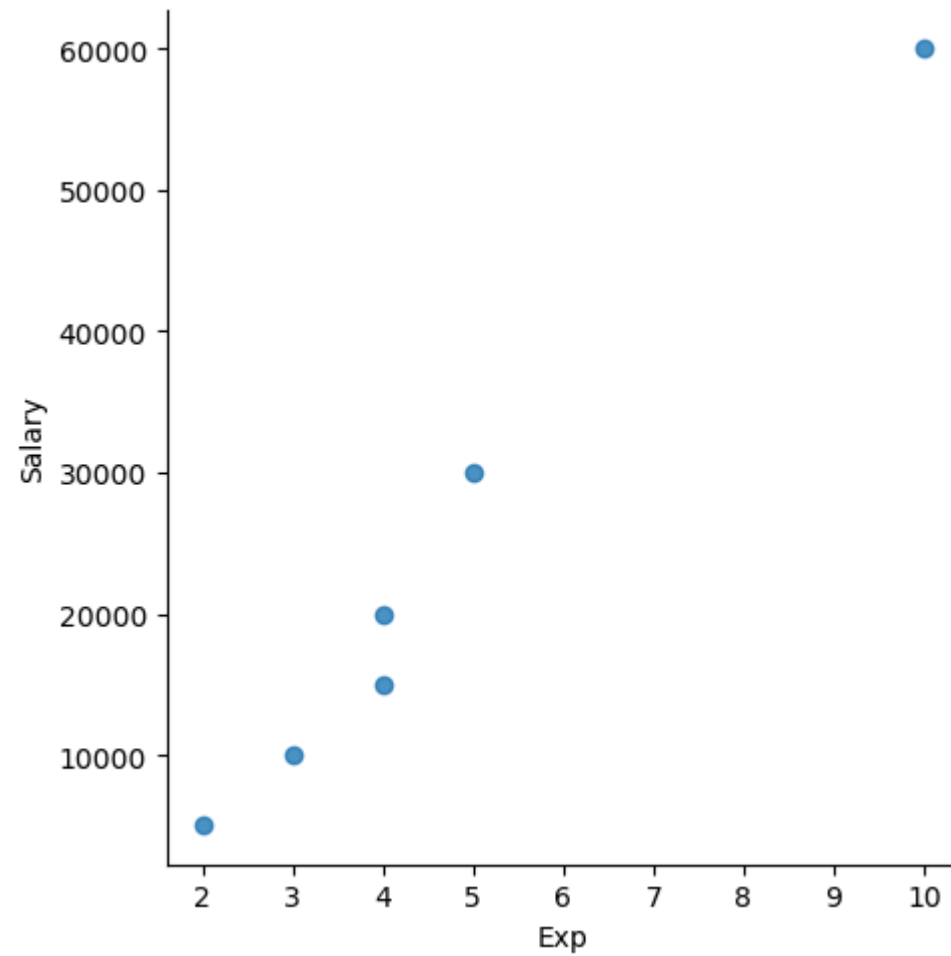
```
In [77]: visualization3 = sns.lmplot(data=clean_data, x='Exp',y='Salary')
```



```
In [78]: visualization3 = sns.lmplot(data=clean_data, x='Exp',y='Salary',fit_reg=True)
```



```
In [79]: visualization3 = sns.lmplot(data=clean_data, x='Exp',y='Salary',fit_reg=False)
```



```
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	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 [81]: `clean_data[:]`

Out[81]:

	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 [82]: `clean_data[:2]`

Out[82]:

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

In [83]: `clean_data[2:]`



```
Out[83]:
```

	<b>Name</b>	<b>Domain</b>	<b>Age</b>	<b>Location</b>	<b>Salary</b>	<b>Exp</b>
<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 [84]: clean_data[:]
```

```
Out[84]:
```

	<b>Name</b>	<b>Domain</b>	<b>Age</b>	<b>Location</b>	<b>Salary</b>	<b>Exp</b>
<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 [85]: clean_data[0:1]
```

```
Out[85]:
```

	<b>Name</b>	<b>Domain</b>	<b>Age</b>	<b>Location</b>	<b>Salary</b>	<b>Exp</b>
<b>0</b>	Mike	Datascience	34	Mumbai	5000	2

```
In [86]: clean_data[0,2]
```

```

-----
KeyError                                Traceback (most recent call last)
File ~\anaconda3\Lib\site-packages\pandas\core\indexes\base.py:3791, in Index.get_loc(self, key)
    3790 try:
-> 3791     return self._engine.get_loc(casted_key)
    3792 except KeyError as err:

File index.pyx:152, in pandas._libs.index.IndexEngine.get_loc()

File index.pyx:181, in pandas._libs.index.IndexEngine.get_loc()

File pandas\_libs\hashtable_class_helper.pxi:7080, in pandas._libs.hashtable.PyObjectHashTable.get_item()

File pandas\_libs\hashtable_class_helper.pxi:7088, in pandas._libs.hashtable.PyObjectHashTable.get_item()

KeyError: (0, 2)

```

The above exception was the direct cause of the following exception:

```

KeyError                                Traceback (most recent call last)
Cell In[86], line 1
----> 1 clean_data[0,2]

File ~\anaconda3\Lib\site-packages\pandas\core\frame.py:3893, in DataFrame.__getitem__(self, key)
    3891 if self.columns.nlevels > 1:
    3892     return self._getitem_multilevel(key)
-> 3893 indexer = self.columns.get_loc(key)
    3894 if is_integer(indexer):
    3895     indexer = [indexer]

File ~\anaconda3\Lib\site-packages\pandas\core\indexes\base.py:3798, in Index.get_loc(self, key)
    3793     if isinstance(casted_key, slice) or (
    3794         isinstance(casted_key, abc.Iterable)
    3795         and any(isinstance(x, slice) for x in casted_key)
    3796     ):
    3797         raise InvalidIndexError(key)
-> 3798     raise KeyError(key) from err
    3799 except TypeError:
    3800     # If we have a listlike key, _check_indexing_error will raise
    3801     # InvalidIndexError. Otherwise we fall through and re-raise

```

```

3802     # the TypeError.
3803     self._check_indexing_error(key)

```

**KeyError:** (0, 2)

In [87]: clean\_data

Out[87]:

	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 [93]: x\_iv # Salary got deleted using drop function

Out[93]:

	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 [94]: x\_iv

Out[94]:

	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 [95]: `x_iv.columns`Out[95]: `Index(['Name', 'Domain', 'Age', 'Location', 'Exp'], dtype='object')`In [97]: `clean_data.columns`Out[97]: `Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')`In [98]: `clean_data`

Out[98]:

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

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 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 [102...

x\_iv

Out[102...

	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 [103...

y\_dv

Out[103...

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

In [ ]:

In [106...

clean\_data

Out[106...

	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 [104...

```
imputation = pd.get_dummies(clean_data)# imputation = is also called as TRANSFOMER
```

In [105...

```
imputation
```

Out[105...

	Age	Salary	Exp	Name_Jane	Name_Kim	Name_Mike	Name_Teddy	Name_Umar	Name_Uttam	Domain_Analytics	Domain_Dataanaly
0	34	5000	2	False	False	True	False	False	False	False	Fals
1	45	10000	3	False	False	False	True	False	False	False	Fals
2	50	15000	4	False	False	False	False	True	False	False	Tru
3	50	20000	4	True	False	False	False	False	False	True	Fals
4	67	30000	5	False	False	False	False	False	True	False	Fals
5	55	60000	10	False	True	False	False	False	False	False	Fals

In [ ]: