

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

```
In [3]: pd.__version__
```

```
Out[3]: '2.2.2'
```

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

```
In [7]: emp
```

```
Out[7]:
```

	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 [9]: id(emp)
```

```
Out[9]: 2139623741232
```

```
In [32]: emp.shape
```

```
Out[32]: (6, 6)
```

```
In [36]: len(emp)
```

```
Out[36]: 6
```

```
In [26]: emp.columns
```

```
Out[26]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')
```

```
In [38]: len(emp.columns)
```

```
Out[38]: 6
```

```
In [40]: 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 [42]: emp

```

Out[42]:
   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+

```

In [44]: emp['Name']

```

Out[44]:
0      Mike
1      Teddy^
2      Uma#r
3      Jane
4      Uttam*
5      Kim
Name: Name, dtype: object

```

In [46]: emp['Domain']

```

Out[46]:
0      Datascience#$
1      Testing
2      Dataanalyst^^#
3      Ana^^lytics
4      Statistics
5      NLP
Name: Domain, dtype: object

```

In [48]: emp['Age']

```

Out[48]:
0      34 years
1      45' yr
2      NaN
3      NaN
4      67-yr
5      55yr
Name: Age, dtype: object

```

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

```
Out[50]: 0      Mumbai
          1      Bangalore
          2         NaN
          3      Hyderbad
          4         NaN
          5       Delhi
          Name: Location, dtype: object
```

```
In [52]: emp['Salary']
```

```
Out[52]: 0      5^00#0
          1     10%%000
          2     1$5%000
          3     2000^0
          4     30000-
          5     6000^$0
          Name: Salary, dtype: object
```

```
In [54]: emp['Exp']
```

```
Out[54]: 0      2+
          1     <3
          2     4> yrs
          3         NaN
          4     5+ year
          5     10+
          Name: Exp, dtype: object
```

```
In [56]: emp[['Name', 'Domain']]
```

```
Out[56]:
```

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

```
In [60]: emp[['Name', 'Domain', 'Age']]
```

Out[60]:

	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 [62]: emp[['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp']]

Out[62]:

	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 [64]: emp

Out[64]:

	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 [28]: emp.head()

Out[28]:

	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]: emp.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]: emp.isnull()
```

```
Out[17]:
```

	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 [19]: emp.isna()
```

```
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]: emp
```

Out[21]:

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

Out[23]:

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

Cleaning

In [66]: `emp`

Out[66]:

	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 [68]: `emp['Name']`

Out[68]:

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

In [70]: `emp['Name'] = emp['Name'].str.replace(r'\W', '', regex=True)`

In [72]: `emp['Name']`

```
Out[72]: 0    Mike
         1    Teddy
         2    Umar
         3    Jane
         4    Uttam
         5    Kim
         Name: Name, dtype: object
```

```
In [74]: emp['Domain'] = emp['Domain'].str.replace(r'\W', '', regex=True)
```

```
In [76]: emp['Domain']
```

```
Out[76]: 0    Datascience
         1      Testing
         2    Dataanalyst
         3      Analytics
         4      Statistics
         5           NLP
         Name: Domain, dtype: object
```

```
In [78]: emp['Age'] = emp['Age'].str.replace(r'\W', '', regex=True)
```

```
In [80]: emp['Age']
```

```
Out[80]: 0    34years
         1     45yr
         2      NaN
         3      NaN
         4     67yr
         5     55yr
         Name: Age, dtype: object
```

```
In [82]: emp['Age'] = emp['Age'].str.extract('(\d+)')
```

```
In [84]: emp['Age']
```

```
Out[84]: 0     34
         1     45
         2    NaN
         3    NaN
         4     67
         5     55
         Name: Age, dtype: object
```

```
In [86]: emp
```

```
Out[86]:
```

	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 [88]: emp['Location']
```

```
Out[88]: 0      Mumbai
1      Bangalore
2         NaN
3      Hyderabad
4         NaN
5         Delhi
Name: Location, dtype: object
```

```
In [90]: emp['Salary']
```

```
Out[90]: 0      5^00#0
1      10%%000
2      1$5%000
3      2000^0
4      30000-
5      6000^$0
Name: Salary, dtype: object
```

```
In [96]: emp['Salary'] = emp['Salary'].str.replace(r'\W', '', regex=True)
```

```
In [98]: emp['Salary']
```

```
Out[98]: 0      5000
1      10000
2      15000
3      20000
4      30000
5      60000
Name: Salary, dtype: object
```

```
In [100... emp
```

```
Out[100...
   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 [102... emp['Exp']
```

```
Out[102... 0      2+
1      <3
2      4> yrs
3      NaN
4      5+ year
5      10+
Name: Exp, dtype: object
```

```
In [104... emp['Exp'] = emp['Exp'].str.extract('(\d+)')
```


In [106...

```
emp['Exp']
```

Out[106...

```
0      2
1      3
2      4
3     NaN
4      5
5     10
Name: Exp, dtype: object
```

In [108...

```
emp
```

Out[108...

	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 [110...

```
clean_data=emp.copy()
```

In [112...

```
emp
```

Out[112...

	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 [169...

```
clean_data
```

Out[169...

	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

Missing Value Treatment

```
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  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 [171...] 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 [116...] clean_data['Age']
```

```
Out[116...]
0    34
1    45
2    NaN
3    NaN
4    67
5    55
Name: Age, dtype: object
```

```
In [118...] import numpy as np
```

```
In [120...] clean_data['Age'] = clean_data['Age'].fillna(np.mean(pd.to_numeric(clean_data['A
```

```
In [122...] clean_data['Age']
```

```
Out[122...]
0    34
1    45
2    50.25
3    50.25
4    67
5    55
Name: Age, dtype: object
```

```
In [124...] clean_data['Exp']
```

```
Out[124...] 0      2
            1      3
            2      4
            3    NaN
            4      5
            5     10
            Name: Exp, dtype: object
```

```
In [126...] clean_data['Exp'] = clean_data['Exp'].fillna(np.mean(pd.to_numeric(clean_data['E
```

```
In [128...] clean_data['Exp']
```

```
Out[128...] 0      2
            1      3
            2      4
            3    4.8
            4      5
            5     10
            Name: Exp, dtype: object
```

```
In [130...] clean_data
```

```
Out[130...]
   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 [134...] clean_data['Location'].isnull().sum()
```

```
Out[134...] 2
```

```
In [136...] clean_data['Location'] = clean_data['Location'].fillna(clean_data['Location'].mo
```

```
In [138...] clean_data['Location']
```

```
Out[138...] 0      Mumbai
            1    Bangalore
            2    Bangalore
            3    Hyderbad
            4    Bangalore
            5      Delhi
            Name: Location, dtype: object
```

```
In [140...] clean_data
```

Out[140...

	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 [142...

```
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 [144...

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

In [146...

```
clean_data['Age']
```

Out[146...

```
0    34
1    45
2    50
3    50
4    67
5    55
Name: Age, dtype: int32
```

In [148...

```
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 [150...

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

In [152... `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 [154... `clean_data['Name'] = clean_data['Name'].astype('category')`
`clean_data['Domain'] = clean_data['Domain'].astype('category')`
`clean_data['Location'] = clean_data['Location'].astype('category')`

In [156... `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 [158... `clean_data`

Out[158...

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

In [162... `import os`
`os.getcwd()`

Out[162... `'C:\\Users\\DELL'`

In [164...

```
clean_data
```

Out[164...

	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

EDA Techniques Let's Apply

In [173...

```
import matplotlib.pyplot as plt
import seaborn as sns
```

In [175...

```
import warnings
warnings.filterwarnings('ignore')
```

In [177...

```
clean_data
```

Out[177...

	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 [179...

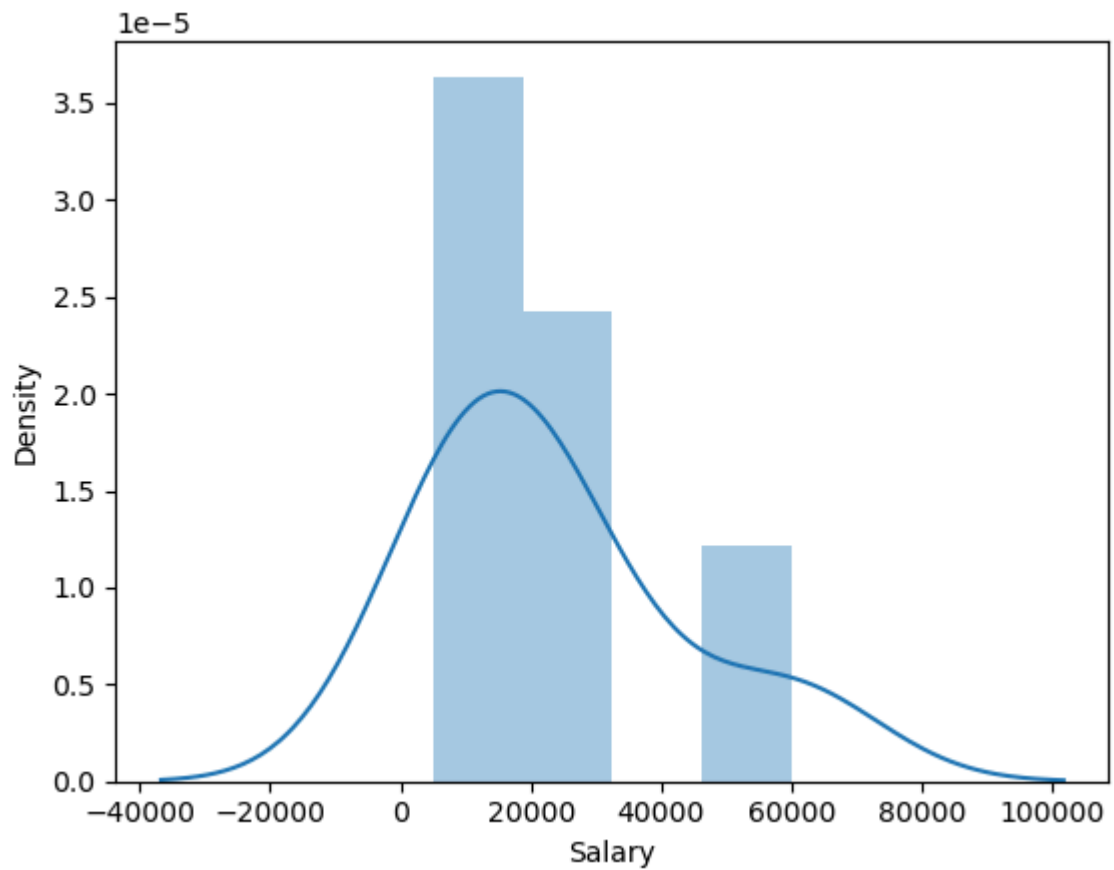
```
clean_data['Salary']
```

Out[179...

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

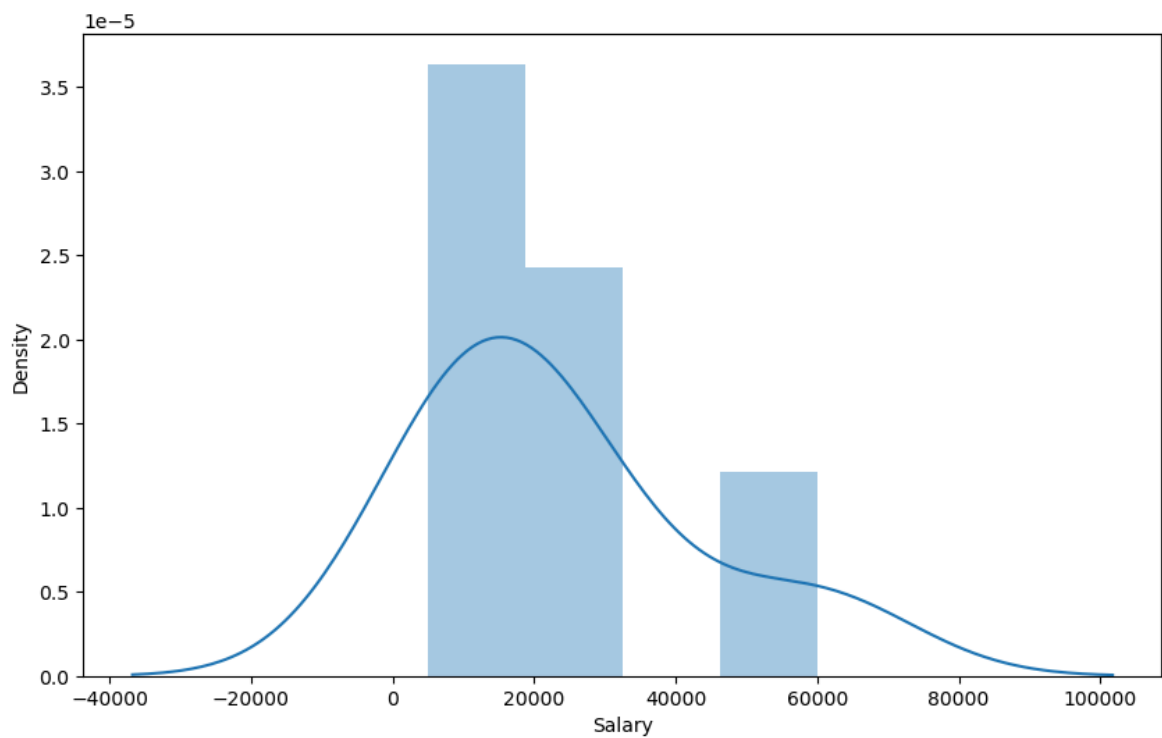
In [185...

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

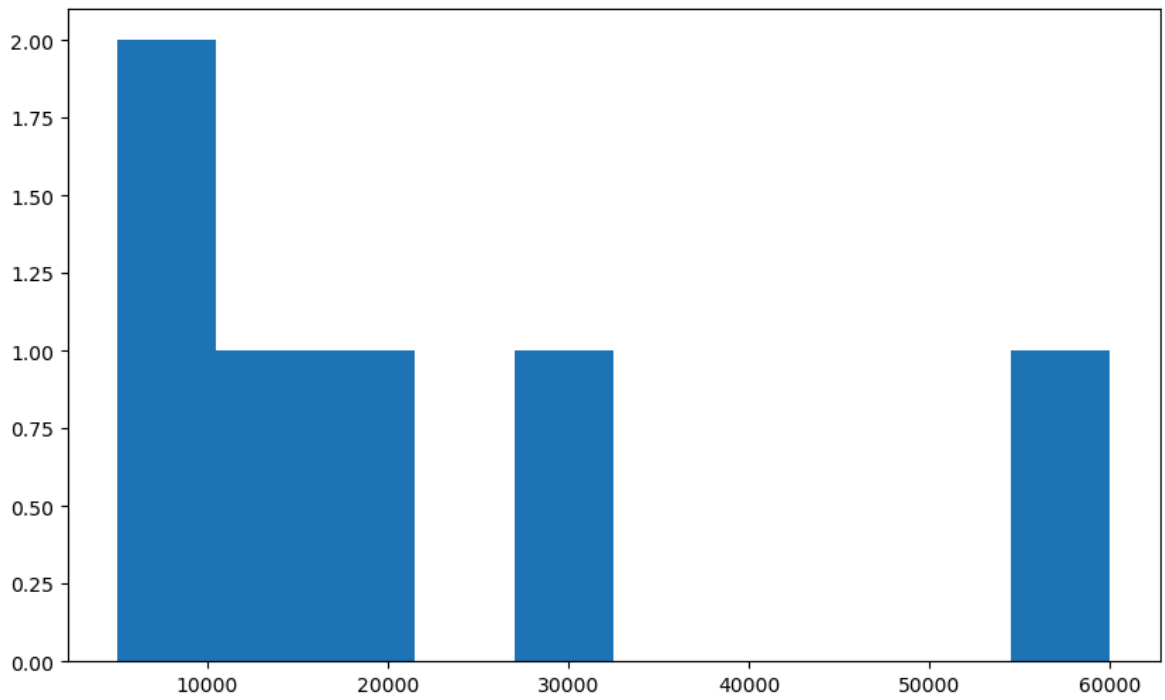


```
In [187...] plt.rcParams['figure.figsize'] = 10,6
```

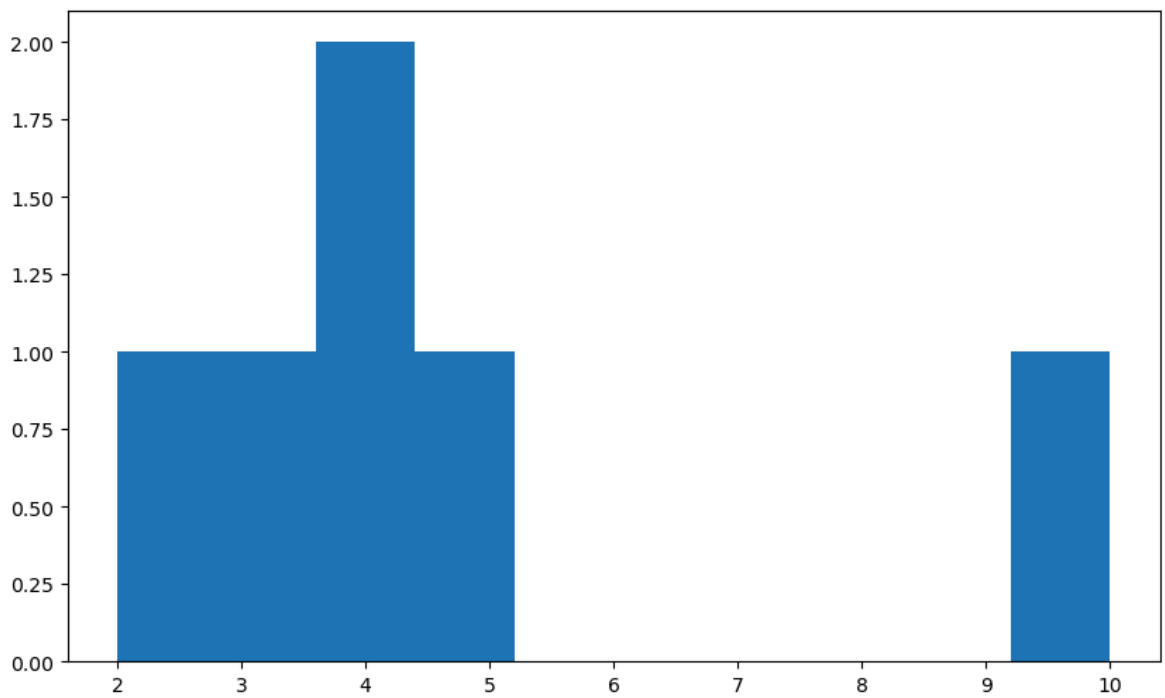
```
In [189...] vis1 = sns.distplot(clean_data['Salary'])
```



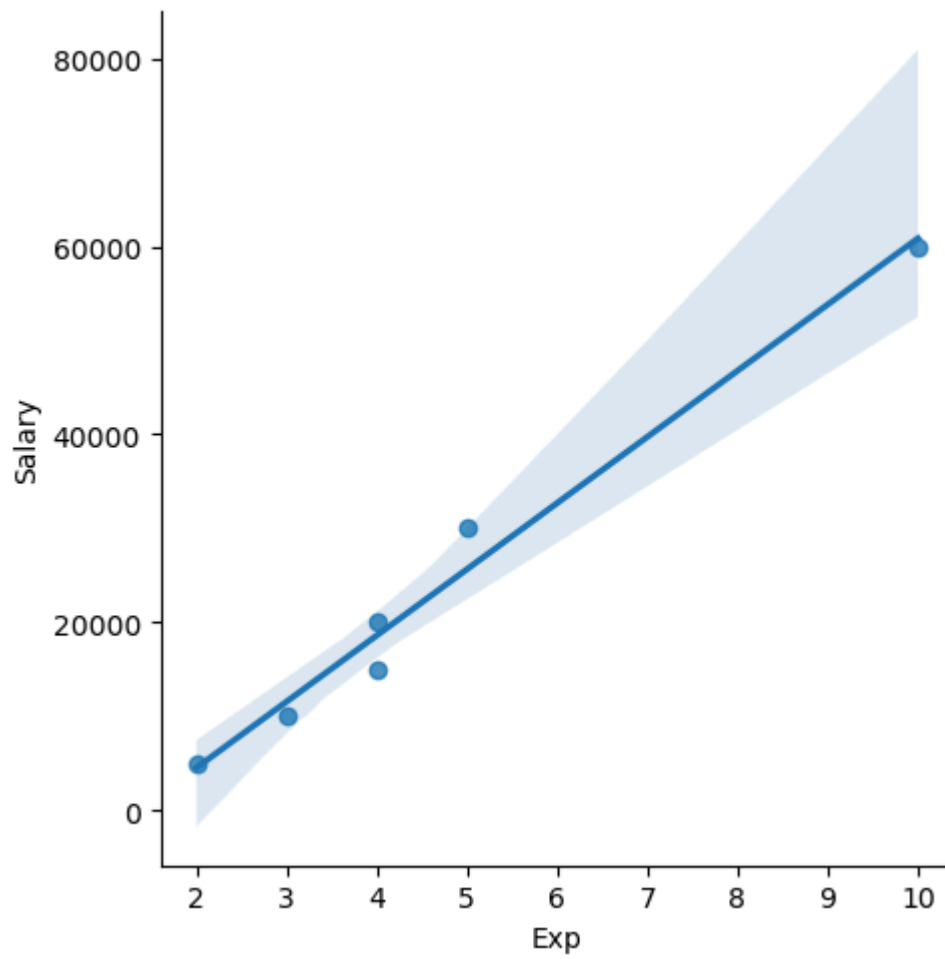
```
In [191...] vis2 = plt.hist(clean_data['Salary'])
```



In [193... `vis3 = plt.hist(clean_data['Exp'])`

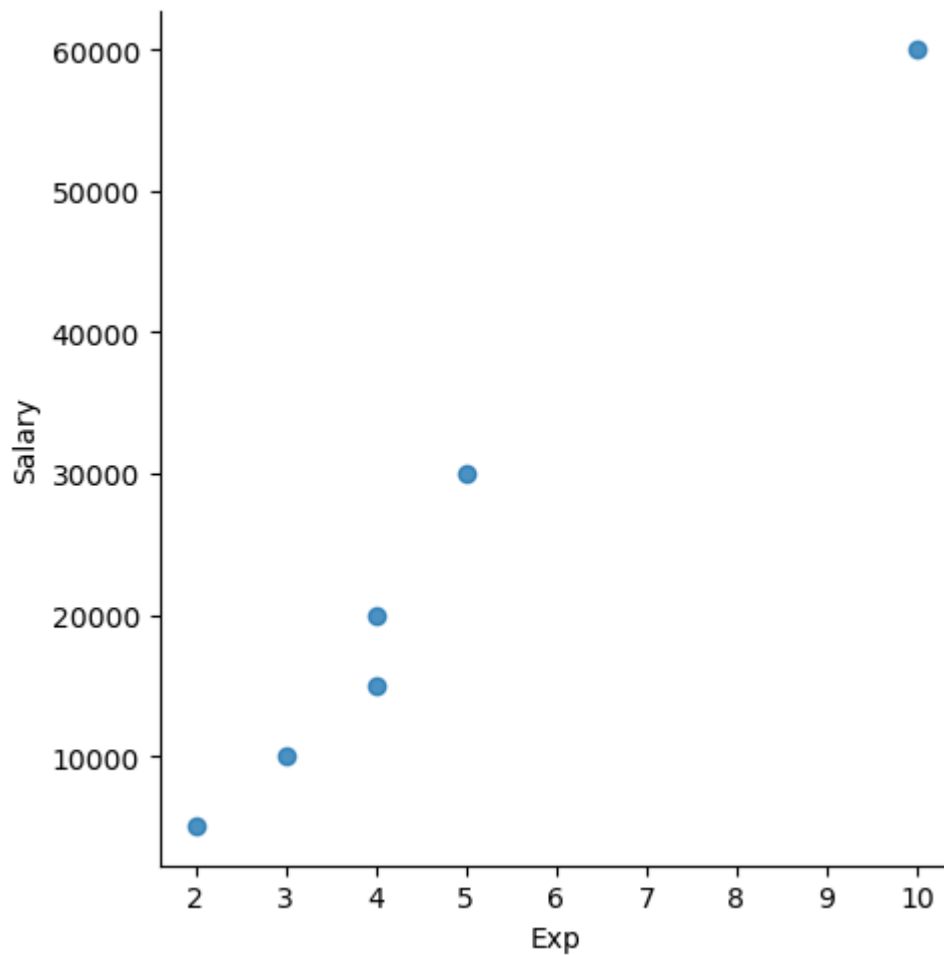


In [195... `vis4 = sns.lmplot(data=clean_data, x = 'Exp', y='Salary')`



In [197...

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



In [199... `clean_data[:]`

Out[199...

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

Out[201...

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

Out[203...

	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 [205...

```
clean_data.columns
```

Out[205...

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

In [207...

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

In [211...

```
X_iv
```

Out[211...

	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 [213...

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

In [215...

```
y_dv
```

Out[215...

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

In [217...

```
emp
```

Out[217...

	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 [219...

clean_data

Out[219...

	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 [221...

X_iv

Out[221...

	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 [225...

y_dv

Out[225...

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

In [227...

clean_data

Out[227...

	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 [229...

imputation = pd.get_dummies(clean_data)

In [231...

imputation

Out[231...

	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 [233...

clean_data

Out[233...

	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 [235...

imputation

Out[235...

	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 []: