

In [234... `import pandas as pd`

In [235... `pd.__version__`

Out[235... `'2.3.0'`

In [236... `emp = pd.read_excel(r"C:\Users\hp\Downloads\Rawdata.xlsx")`
`emp`

Out[236...

	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 [237... `id(emp)`

Out[237... `3057123670544`

In [238... `emp.columns`

Out[238... `Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')`

In [239... `emp.shape`

Out[239... `(6, 6)`

In [240... `emp.head()`

Out[240...

	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 [241... `emp.tail()`

Out[241...

	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 [242...

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 [243...

emp.isnull()

Out[243...

	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 [244...

emp.isna()

Out[244...

	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 [245... `emp.isnull().sum()`

Out[245...
 Name 0
 Domain 0
 Age 2
 Location 2
 Salary 0
 Exp 1
 dtype: int64

DATA CLEANING OR DATA CLEANSING

In [246... `emp`

Out[246...

	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 [247... `emp['Name']`

Out[247...
 0 Mike
 1 Teddy^
 2 Uma#r
 3 Jane
 4 Uttam*
 5 Kim
 Name: Name, dtype: object

In [248... `emp['Name'] = emp['Name'].str.replace(r'\W', '', regex=True)`

In [249... `emp['Name']`

Out[249...
 0 Mike
 1 Teddy
 2 Umar
 3 Jane
 4 Uttam
 5 Kim
 Name: Name, dtype: object

In [250... `emp`

Out[250...

	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 [251...

emp['Domain']

Out[251...

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

In [252...

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

In [253...

emp['Domain']

Out[253...

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

In [254...

emp

Out[254...

	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	Analytics	NaN	Hyderbad	2000^0	NaN
4	Uttam	Statistics	67-yr	NaN	30000-	5+ year
5	Kim	NLP	55yr	Delhi	6000^\$0	10+

In [255...

emp['Age']

```
Out[255...] 0    34 years
            1    45' yr
            2      NaN
            3      NaN
            4    67-yr
            5    55yr
            Name: Age, dtype: object
```

```
In [256...] emp['Age'] = emp['Age'].str.extract('(\d+)')
```

```
In [257...] emp['Age']
```

```
Out[257...] 0    34
            1    45
            2    NaN
            3    NaN
            4    67
            5    55
            Name: Age, dtype: object
```

```
In [258...] emp
```

```
Out[258...]
   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  Hyderabad  2000^0  NaN
4  Uttam  Statistics  67      NaN  30000-  5+ year
5  Kim     NLP      55      Delhi  6000^$0  10+
```

```
In [259...] emp['Location']
```

```
Out[259...] 0    Mumbai
            1  Bangalore
            2      NaN
            3  Hyderabad
            4      NaN
            5    Delhi
            Name: Location, dtype: object
```

```
In [260...] emp['Location'] = emp['Location'].str.replace(r'\W', '', regex=True)
```

```
In [261...] emp['Location']
```

```
Out[261...] 0    Mumbai
            1  Bangalore
            2      NaN
            3  Hyderabad
            4      NaN
            5    Delhi
            Name: Location, dtype: object
```

```
In [262...] emp
```

Out[262...

	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 [263...

emp['Salary']

Out[263...

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

In [264...

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

In [265...

emp['Salary']

Out[265...

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

In [266...

emp

Out[266...

	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 [267...

emp['Exp']

```
Out[267...] 0      2+
            1      <3
            2      4> yrs
            3      NaN
            4      5+ year
            5      10+
            Name: Exp, dtype: object
```

```
In [268...] emp['Exp'] = emp['Exp'].str.extract('(\d+)')
```

```
In [269...] emp['Exp']
```

```
Out[269...] 0      2
            1      3
            2      4
            3      NaN
            4      5
            5      10
            Name: Exp, dtype: object
```

```
In [270...] clean_data = emp.copy()
```

```
In [271...] clean_data
```

```
Out[271...]
   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 [272...] clean_data.isnull().sum()
```

```
Out[272...] Name      0
            Domain    0
            Age       2
            Location   2
            Salary     0
            Exp       1
            dtype: int64
```

```
In [273...] clean_data['Age']
```

```
Out[273...] 0      34
            1      45
            2      NaN
            3      NaN
            4      67
            5      55
            Name: Age, dtype: object
```

```
In [274...] import numpy as np
```

```
In [275... clean_data['Age'] = clean_data['Age'].fillna(np.mean(pd.to_numeric(clean_data['A
```

```
In [276... clean_data['Age']
```

```
Out[276... 0      34
1      45
2     50.25
3     50.25
4      67
5      55
Name: Age, dtype: object
```

```
In [277... clean_data['Exp']
```

```
Out[277... 0      2
1      3
2      4
3     NaN
4      5
5     10
Name: Exp, dtype: object
```

```
In [278... clean_data['Exp'] = clean_data['Exp'].fillna(np.mean(pd.to_numeric(clean_data['E
```

```
In [279... clean_data['Exp']
```

```
Out[279... 0      2
1      3
2      4
3     4.8
4      5
5     10
Name: Exp, dtype: object
```

```
In [280... clean_data
```

```
Out[280...
```

	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 [281... clean_data['Location'].isnull().sum()
```

```
Out[281... np.int64(2)
```

```
In [282... clean_data['Location']
```



```
Out[282...] 0      Mumbai
            1      Bangalore
            2      NaN
            3      Hyderabad
            4      NaN
            5      Delhi
            Name: Location, dtype: object
```

```
In [283...] clean_data['Location'] = clean_data['Location'].fillna(clean_data['Location'].mode[0])
```

```
In [284...] clean_data['Location']
```

```
Out[284...] 0      Mumbai
            1      Bangalore
            2      Bangalore
            3      Hyderabad
            4      Bangalore
            5      Delhi
            Name: Location, dtype: object
```

```
In [285...] clean_data
```

```
Out[285...]
   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 [286...] 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 [287...] 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 [288... clean_data['Age']=clean_data['Age'].astype(int)
```

```
In [289... clean_data['Age']
```

```
Out[289... 0    34
1    45
2    50
3    50
4    67
5    55
Name: Age, dtype: int64
```

```
In [290... 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     int64
3   Location    6 non-null     object
4   Salary      6 non-null     object
5   Exp         6 non-null     object
dtypes: int64(1), object(5)
memory usage: 420.0+ bytes
```

```
In [291... clean_data['Salary'] = clean_data['Salary'].astype(int)
```

```
In [292... clean_data['Salary']
```

```
Out[292... 0    5000
1   10000
2   15000
3   20000
4   30000
5   60000
Name: Salary, dtype: int64
```

```
In [293... clean_data['Exp'] = clean_data['Exp'].astype(int)
```

```
In [294... clean_data['Exp']
```

```
Out[294...] 0      2
            1      3
            2      4
            3      4
            4      5
            5     10
            Name: Exp, dtype: int64
```

```
In [295...] 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     int64
3   Location    6 non-null     object
4   Salary      6 non-null     int64
5   Exp         6 non-null     int64
dtypes: int64(3), object(3)
memory usage: 420.0+ bytes
```

```
In [296...] 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 [297...] 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     int64
3   Location    6 non-null     category
4   Salary      6 non-null     int64
5   Exp         6 non-null     int64
dtypes: category(3), int64(3)
memory usage: 938.0 bytes
```

```
In [298...] clean_data
```

```
Out[298...]   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 [299... clean_data.to_csv('clean_data.csv')
```

```
In [300... import os  
os.getcwd()
```

```
Out[300... 'C:\\Users\\hp'
```

```
In [301... clean_data
```

```
Out[301... 
```

	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

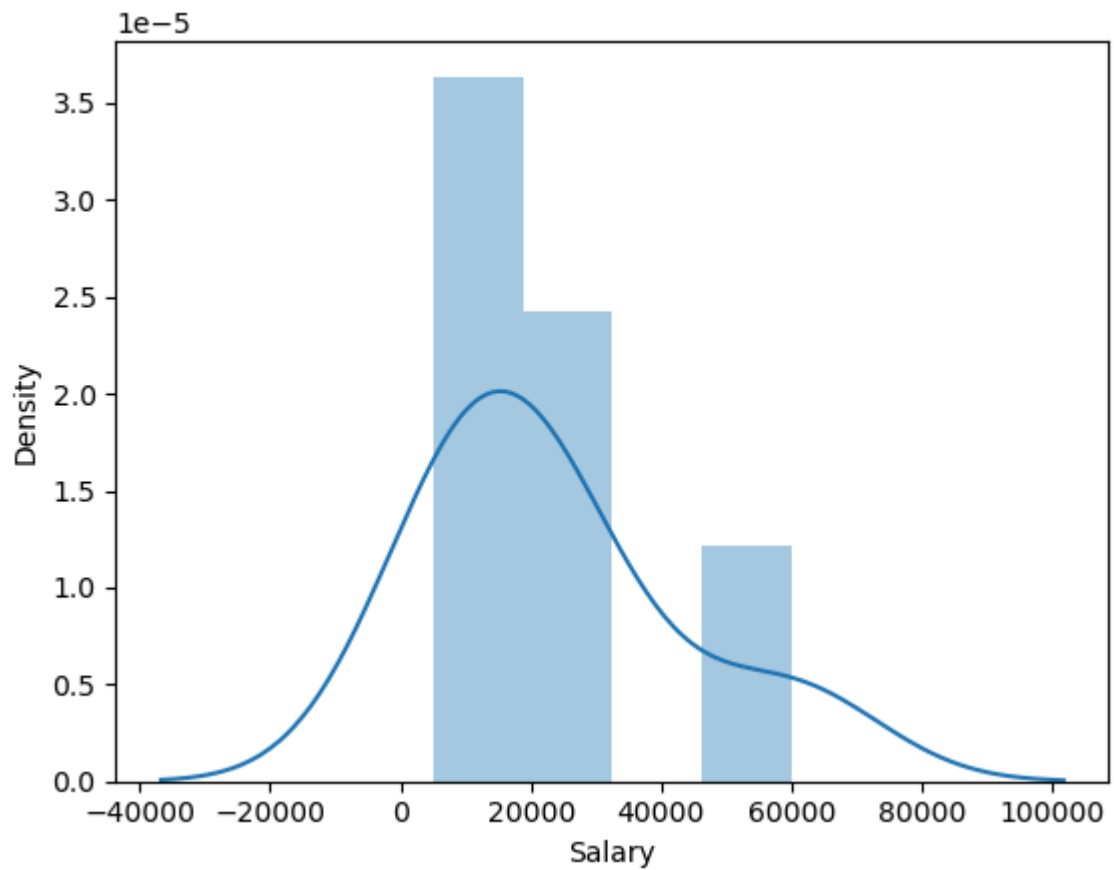
```
In [302... import matplotlib.pyplot as plt  
import seaborn as sns
```

```
In [303... import warnings  
warnings.filterwarnings('ignore')
```

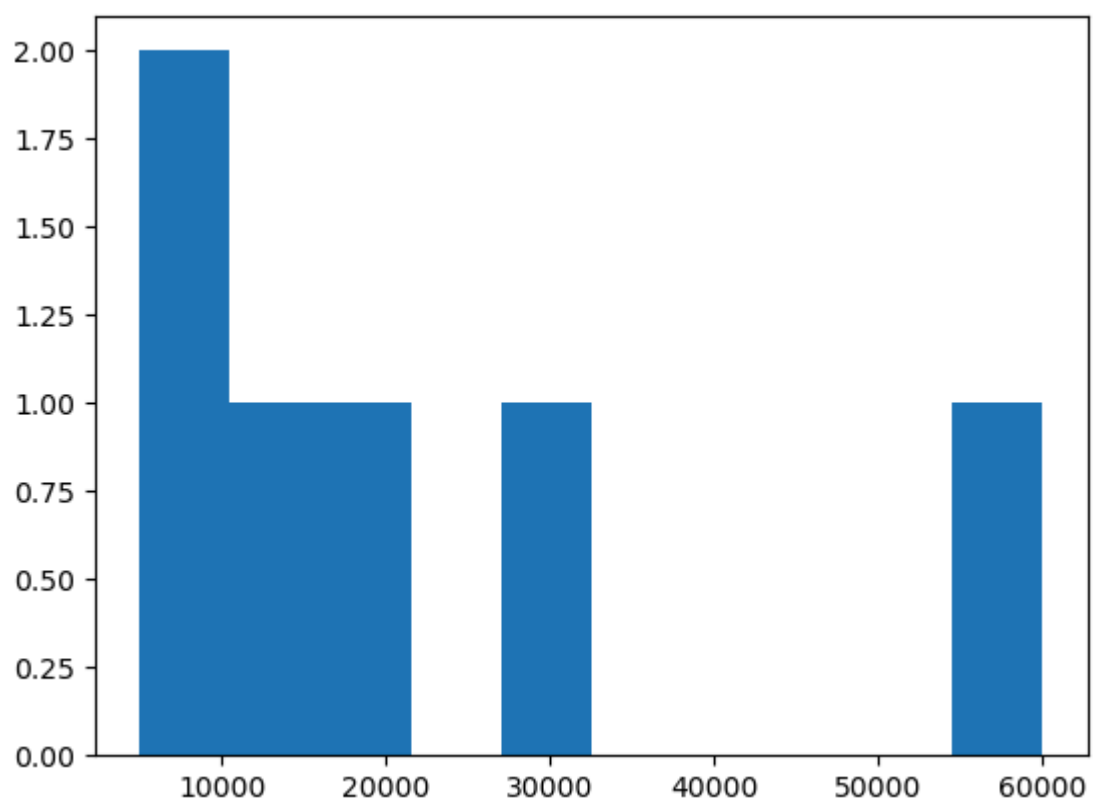
```
In [304... clean_data['Salary']
```

```
Out[304... 0    5000  
1   10000  
2   15000  
3   20000  
4   30000  
5   60000  
Name: Salary, dtype: int64
```

```
In [305... vis1 = sns.distplot(clean_data['Salary'])
```



```
In [306... vis2 = plt.hist(clean_data['Salary'])
```



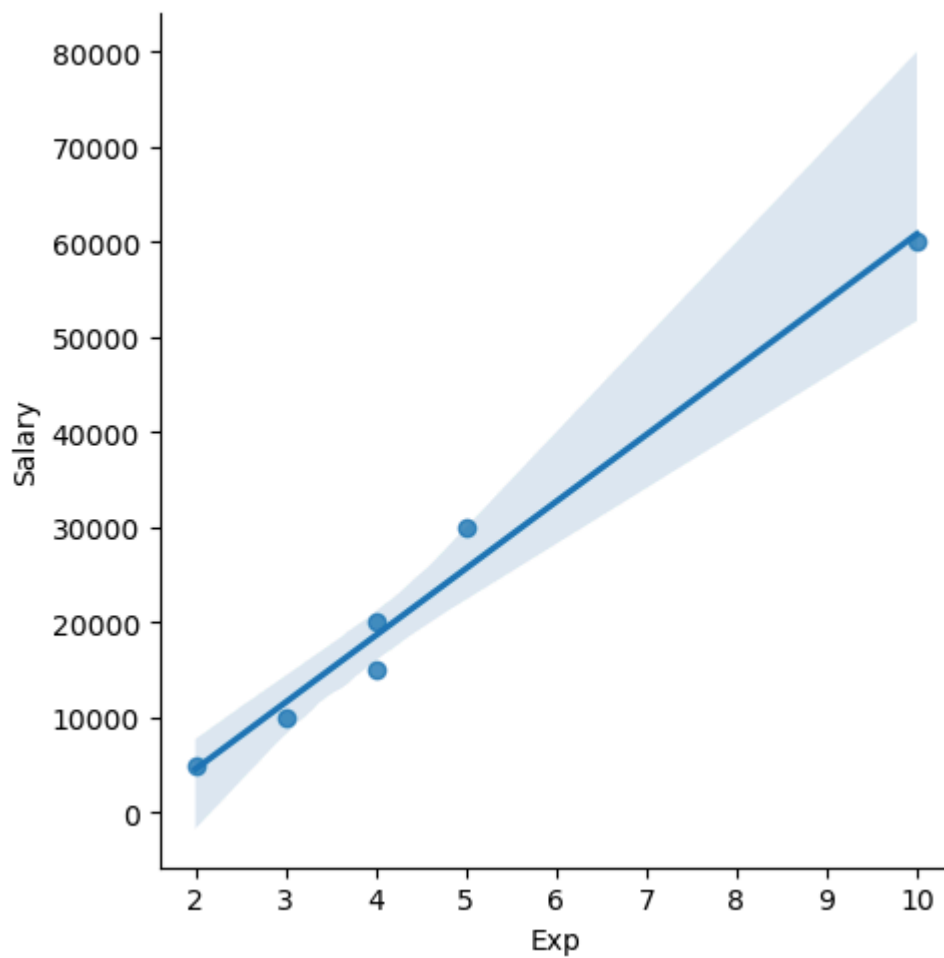
```
In [307... clean_data
```

Out[307...

	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

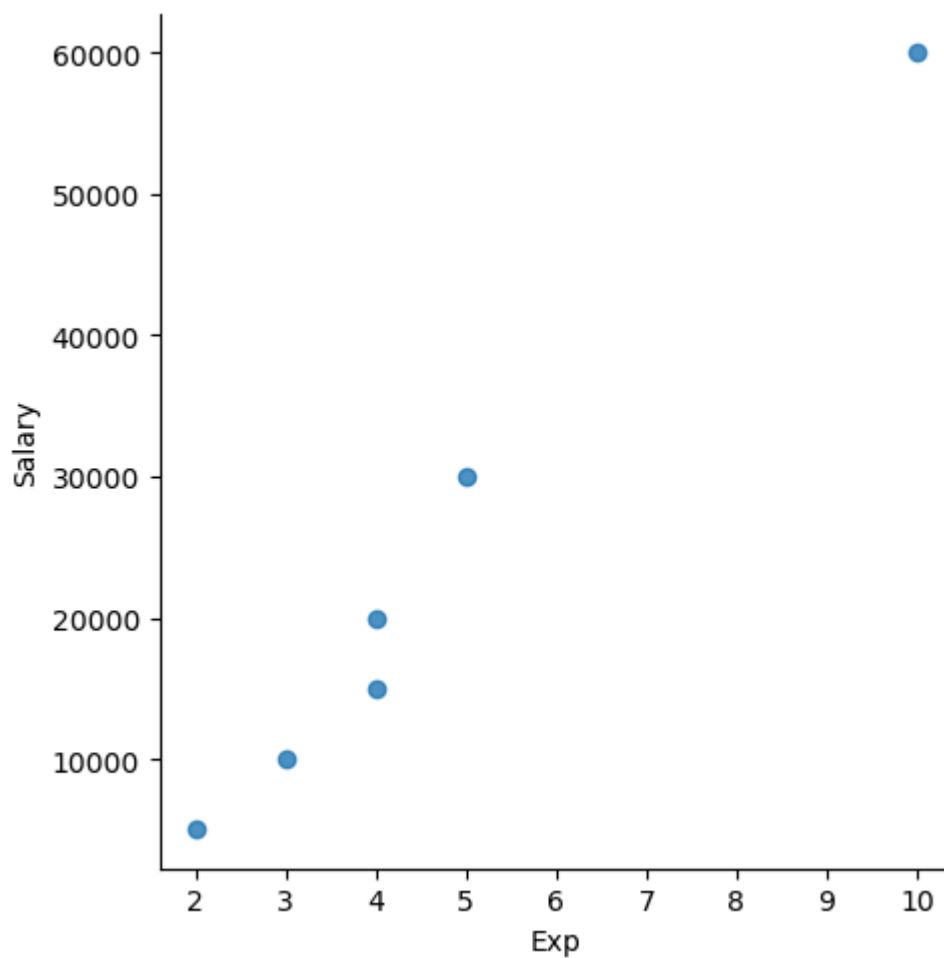
In [308...

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



In [309...

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



In [310... `clean_data[:]`

Out[310...

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

Out[311...

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

Out[312...

	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 [313...

`clean_data.columns`

Out[313...

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

In [314...

`clean_data`

Out[314...

	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 [315...

`emp`

Out[315...

	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 [316...

`clean_data`

Out[316...

	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 [320...

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

Out[320...

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

In [321...

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

Out[321...

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

In [322...

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

Out[322...

	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 [323...

clean_data

Out[323...

	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 [324...

len(clean_data)

Out[324...

6

In [325...

imputation.columns

Out[325...

```
Index(['Age', 'Salary', 'Exp', 'Name_Jane', 'Name_Kim', 'Name_Mike',
      'Name_Teddy', 'Name_Umar', 'Name_Uttam', 'Domain_Analytics',
      'Domain_Dataanalyst', 'Domain_Datascience', 'Domain_NLP',
      'Domain_Statistics', 'Domain_Testing', 'Location_Bangalore',
      'Location_Delhi', 'Location_Hyderabad', 'Location_Mumbai'],
      dtype='object')
```

In [326...

id(emp)

Out[326...

3057123670544

In [327...

emp.columns

Out[327...

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

In [328...

emp.shape

Out[328...

(6, 6)

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []: