

Day 24 - Exploratory Data Analysis: Practical

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')
```

```
In [2]: emp = pd.read_excel(r'C:\Users\Arman\Downloads\dataset\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.head()
```

```
Out[4]:
```

	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 [5]: emp.tail()
```

```
Out[5]:
```

	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 [6]: emp.shape
```

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

```
In [7]: emp.columns
```

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

```
In [8]: 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 [9]: emp['Name']
```

```
Out[9]: 0      Mike
1      Teddy^
2      Uma#r
3      Jane
4      Uttam*
5      Kim
Name: Name, dtype: object
```

```
In [10]: emp['Domain']
```

```
Out[10]: 0      Datascience#$
1          Testing
2      Dataanalyst^^#
3      Ana^^lytics
4          Statistics
5          NLP
Name: Domain, dtype: object
```

```
In [11]: emp['Age']
```

```
Out[11]: 0      34 years
1      45' yr
2      NaN
3      NaN
4      67-yr
5      55yr
Name: Age, dtype: object
```

```
In [12]: emp['Location']
```

```
Out[12]: 0      Mumbai
1      Bangalore
2      NaN
3      Hyderbad
4      NaN
5      Delhi
Name: Location, dtype: object
```

```
In [13]: emp['Salary']
```

```
Out[13]: 0      5^00#0
1      10%%000
2      1$5%000
3      2000^0
4      30000-
5      6000^$0
Name: Salary, dtype: object
```

```
In [14]: emp['Exp']
```

```
Out[14]: 0      2+
         1      <3
         2      4> yrs
         3      NaN
         4      5+ year
         5      10+
         Name: Exp, dtype: object
```

```
In [15]: emp
```

```
Out[15]:
```

	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+

Data Cleaning

```
In [16]: emp['Name'] = emp['Name'].str.replace(r'\W', '', regex=True)
emp['Name']
```

```
Out[16]: 0      Mike
         1      Teddy
         2      Umar
         3      Jane
         4      Uttam
         5      Kim
         Name: Name, dtype: object
```

```
In [17]: emp['Domain'] = emp['Domain'].str.replace(r'\W', '', regex=True)
emp['Age'] = emp['Age'].str.replace(r'\W', '', regex=True)
emp['Salary'] = emp['Salary'].str.replace(r'\W', '', regex=True)
emp['Location'] = emp['Location'].str.replace(r'\W', '', regex=True)
```

```
In [18]: emp
```

```
Out[18]:
```

	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 [19]: emp['Age'] = emp['Age'].str.extract(r'(\d+)')
emp['Exp'] = emp['Exp'].str.extract(r'(\d+)')
```

```
In [20]: emp
```

Out[20]:

	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 [21]: `clean_data = emp.copy()`

In [22]: `clean_data`

Out[22]:

	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

Missing Value Treatment

In [23]: `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 [24]: `clean_data`

Out[24]:

	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 [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['Age'])))
         clean_data['Exp'] = clean_data['Exp'].fillna(np.mean(pd.to_numeric(clean_data['Exp'])))
```

```
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]: clean_data['Exp']
```

```
Out[28]: 0     2
         1     3
         2     4
         3   4.8
         4     5
         5    10
         Name: Exp, dtype: object
```

```
In [29]: emp
```

```
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	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 [30]: clean_data
```

```
Out[30]:
```

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

```
Out[31]: 0      Mumbai
1      Bangalore
2      Bangalore
3      Hyderabad
4      Bangalore
5      Delhi
Name: Location, dtype: object
```

```
In [32]: 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 [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	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 [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      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 [35]: 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 [36]: 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 [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	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

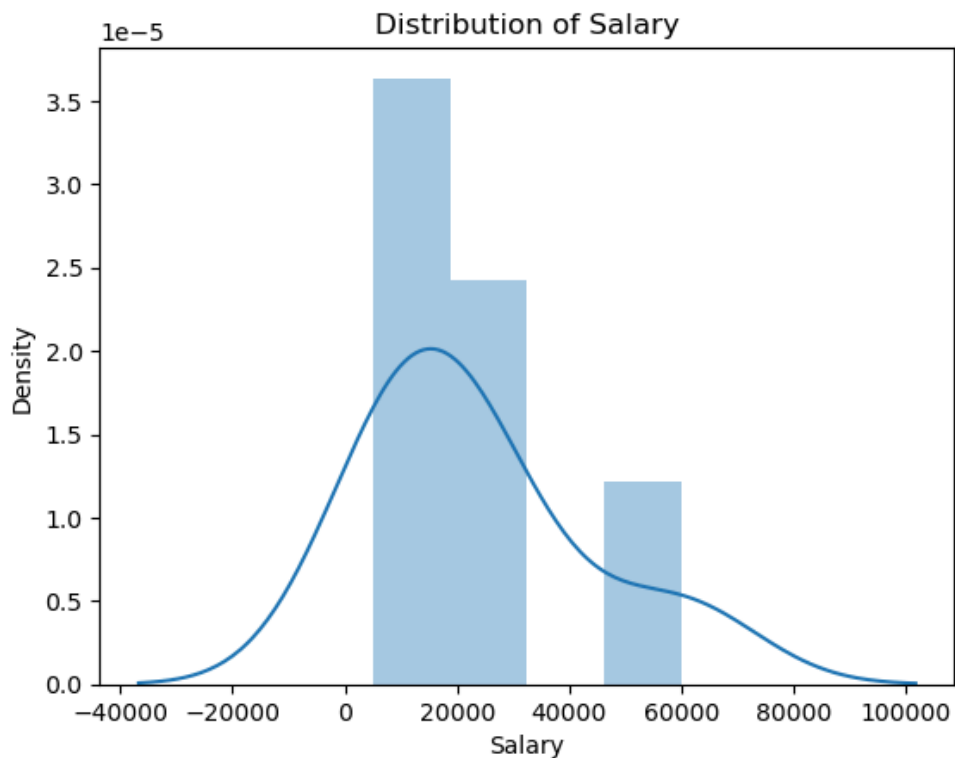
```
In [38]: clean_data.to_csv('clean_data.csv')
```

```
In [39]: clean_data['Salary']
```

```
Out[39]: 0    5000
1   10000
2   15000
3   20000
4   30000
5   60000
Name: Salary, dtype: int64
```

```
In [40]: vis1 = sns.distplot(clean_data['Salary'])
plt.title("Distribution of Salary")
```

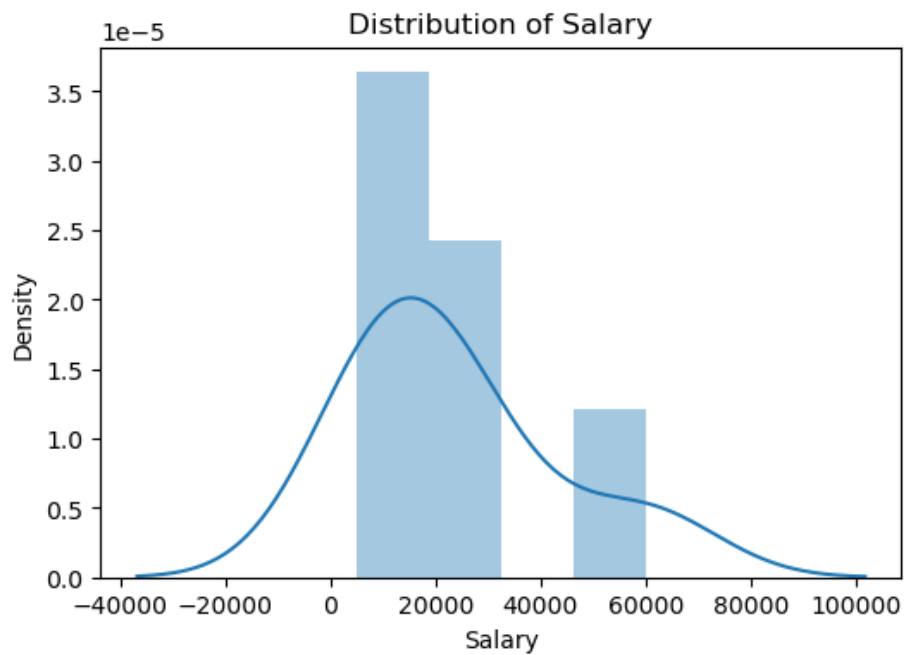
```
Out[40]: Text(0.5, 1.0, 'Distribution of Salary')
```



```
In [41]: plt.rcParams['figure.figsize'] = 6,4
```

```
In [42]: vis1 = sns.distplot(clean_data['Salary'])
plt.title("Distribution of Salary")
```

```
Out[42]: Text(0.5, 1.0, 'Distribution of Salary')
```



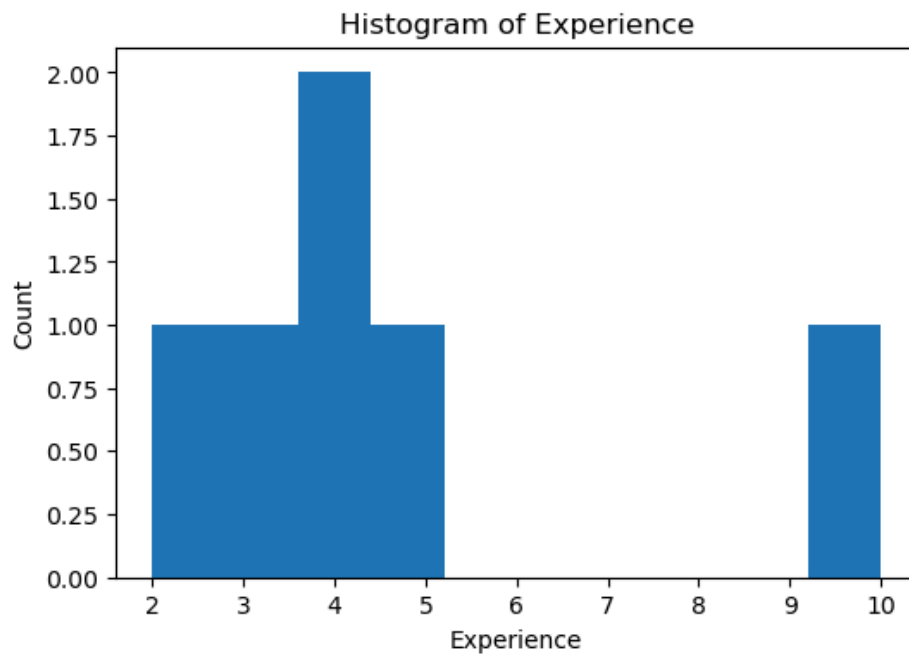
```
In [43]: vis2 = plt.hist(clean_data['Salary'])  
plt.title("Histogram of Salary")  
plt.xlabel("Salary")  
plt.ylabel("Count")
```

Out[43]: Text(0, 0.5, 'Count')

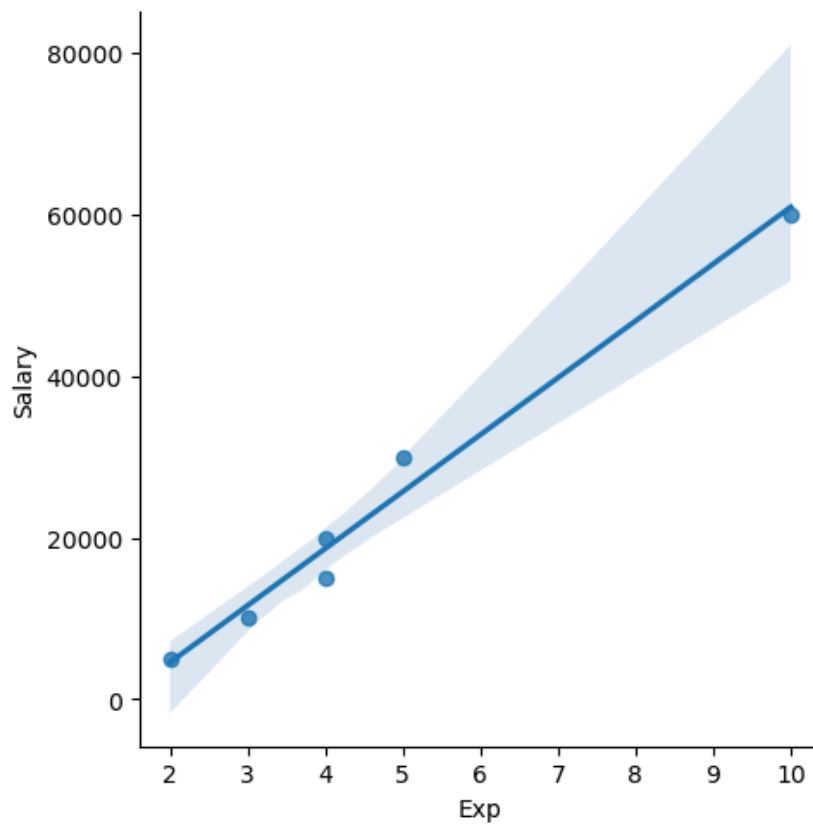


```
In [44]: vis3 = plt.hist(clean_data['Exp'])  
plt.title("Histogram of Experience")  
plt.xlabel("Experience")  
plt.ylabel("Count")
```

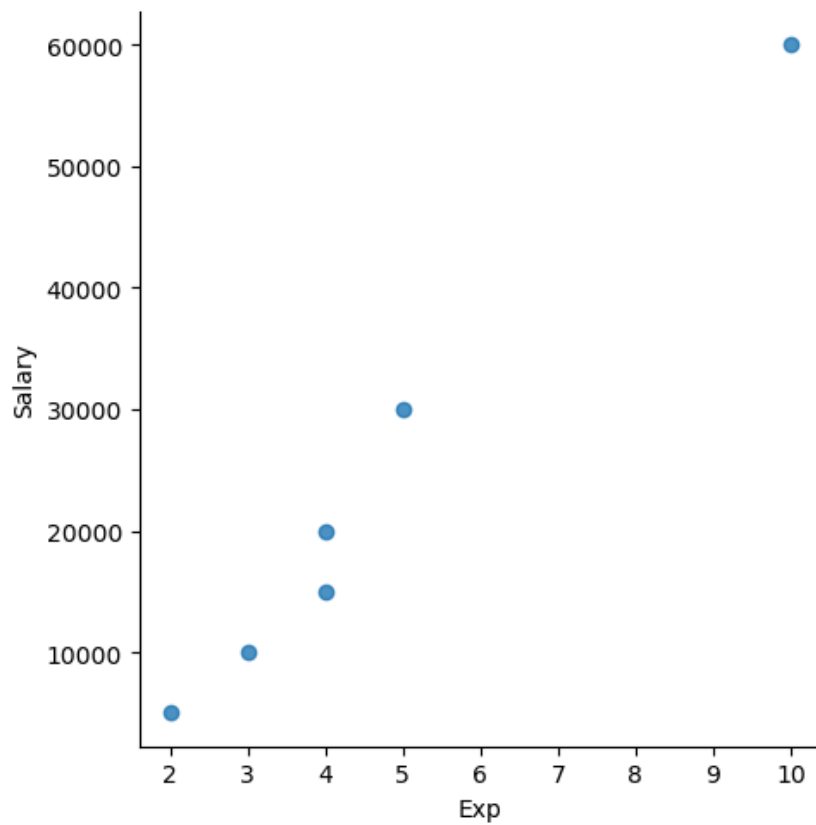
Out[44]: Text(0, 0.5, 'Count')



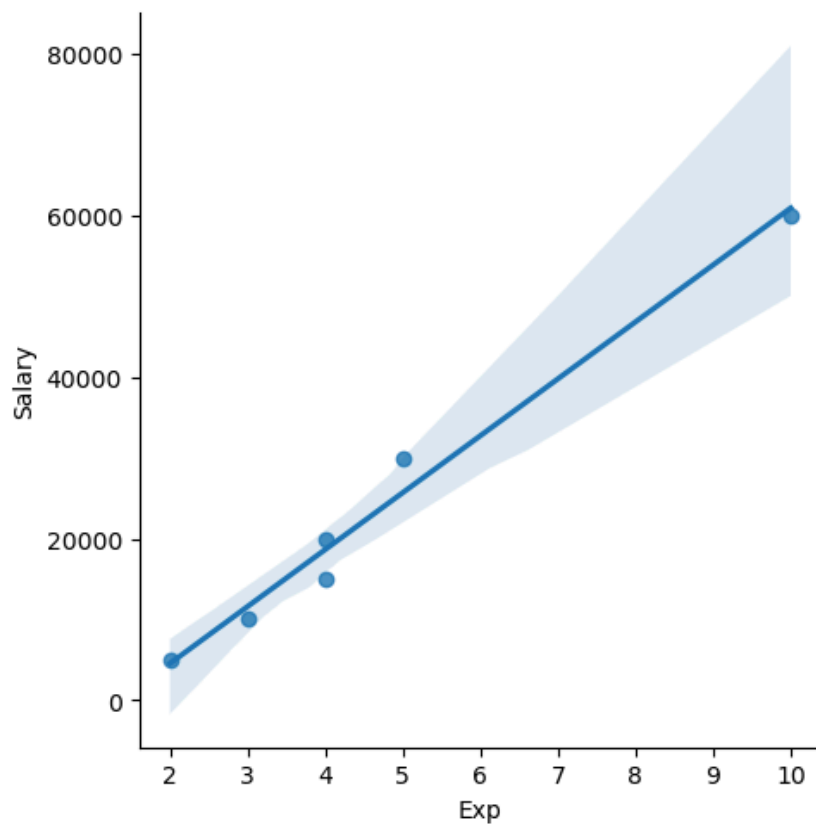
```
In [45]: vis4 = sns.lmplot(data=clean_data, x = 'Exp', y='Salary')
```



```
In [46]: vis5 = sns.lmplot(data=clean_data, x = 'Exp', y='Salary', fit_reg = False)
```



```
In [47]: vis6 = sns.lmplot(data=clean_data, x = 'Exp', y='Salary', fit_reg = True)
```



```
In [48]: clean_data
```

Out[48]:

	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 [49]: `clean_data[:2]`

Out[49]:

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

In [50]: `clean_data[2:]`

Out[50]:

	Name	Domain	Age	Location	Salary	Exp
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 [51]: `clean_data[0:6:2]`

Out[51]:

	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 [52]: `clean_data`

Out[52]:

	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 [53]: `x_iv = emp[['Name', 'Domain', 'Age', 'Location', 'Exp']] # Independent variables`
`y_dv = emp[['Salary']] # Dependent variable`

In [54]: `x_iv`

Out[54]:

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

In [55]: y_dv

Out[55]:

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

In [56]: imputation = pd.get_dummies(clean_data)

In [57]: imputation

Out[57]:

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

