

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

```
In [2]: pd.__version__
```

```
Out[2]: '2.3.1'
```

```
In [3]: pip install openpyxl
```

```
Requirement already satisfied: openpyxl in c:\users\dell\appdata\local\programs\python\python313\lib\site-packages (3.1.5)
```

```
Requirement already satisfied: et-xmlfile in c:\users\dell\appdata\local\programs\python\python313\lib\site-packages (from openpyxl) (2.0.0)
```

```
Note: you may need to restart the kernel to use updated packages.
```

```
In [4]: emp = pd.read_excel(r"C:\Users\DELL\Downloads\Rawdata.xlsx")
```

```
In [5]: emp
```

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

```
In [6]: import numpy as np
```

```
In [7]: np.__version__
```

```
Out[7]: '2.3.1'
```

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

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

```
In [9]: len(emp)
```

```
Out[9]: 6
```

```
In [10]: emp.columns
```

```
Out[10]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')
```

```
In [11]: len(emp.columns)
```

Out[11]: 6

In [12]: 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 [13]: emp

	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 [14]: emp['Name']

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

In [15]: emp['Domain']

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

In [16]: emp['Age']

```
Out[16]: 0    34 years
         1    45' yr
         2      NaN
         3      NaN
         4    67-yr
         5    55yr
Name: Age, dtype: object
```

```
In [17]: emp['Location']
```

```
Out[17]: 0      Mumbai
         1    Bangalore
         2      NaN
         3    Hyderabad
         4      NaN
         5      Delhi
Name: Location, dtype: object
```

```
In [18]: emp['Salary']
```

```
Out[18]: 0    5^00#0
         1  10%%000
         2  1$5%000
         3    2000^0
         4    30000-
         5   6000^$0
Name: Salary, dtype: object
```

```
In [19]: emp['Exp']
```

```
Out[19]: 0      2+
         1      <3
         2    4> yrs
         3      NaN
         4    5+ year
         5     10+
Name: Exp, dtype: object
```

```
In [20]: emp[['Domain', 'Domain']]
```

	Domain	Domain
0	Datascience#\$	Datascience#\$
1	Testing	Testing
2	Dataanalyst^^#	Dataanalyst^^#
3	Ana^^lytics	Ana^^lytics
4	Statistics	Statistics
5	NLP	NLP

```
In [21]: emp[['Domain', 'Domain', 'Age', 'Location', 'Salary', 'Exp']]
```

Out[21]:

	Domain	Domain	Age	Location	Salary	Exp
0	Datascience#\$	Datascience#\$	34 years	Mumbai	5^00#0	2+
1	Testing	Testing	45' yr	Bangalore	10%#000	<3
2	Dataanalyst^^#	Dataanalyst^^#	NaN	NaN	1\$5%000	4> yrs
3	Ana^^lytics	Ana^^lytics	NaN	Hyderbad	2000^0	NaN
4	Statistics	Statistics	67-yr	NaN	30000-	5+ year
5	NLP	NLP	55yr	Delhi	6000^\$0	10+

DATA CLEANSING

In [22]: `emp['Name']`

```
Out[22]: 0      Mike
         1    Teddy^
         2     Uma#r
         3      Jane
         4    Uttam*
         5      Kim
Name: Name, dtype: object
```



Correct (for newer pandas versions):

(`regex=True`): IS ADD AND USED

In [23]: `emp['Name'] = emp['Name'].str.replace(r'\W', ' ', regex=True)`In [24]: `emp['Name']`

```
Out[24]: 0      Mike
         1    Teddy
         2     Umar
         3      Jane
         4    Uttam
         5      Kim
Name: Name, dtype: object
```

In [25]: `emp['Age'] = emp['Age'].str.replace(r'\W', ' ', regex=True)`In [26]: `emp['Age']`

```
Out[26]: 0    34years
         1    45yr
         2    NaN
         3    NaN
         4    67yr
         5    55yr
Name: Age, dtype: object
```

```
In [27]: #  Correct (for newer pandas versions):
# .str.extract(r'(\d+)') :- IS USED
emp['Age'] = emp['Age'].str.extract(r'(\d+)')
```

```
In [28]: emp['Age']
```

```
Out[28]: 0    34
         1    45
         2    NaN
         3    NaN
         4    67
         5    55
Name: Age, dtype: object
```

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

	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	Ana^^lytics	NaN	Hyderabad	2000^0	NaN
4	Uttam	Statistics	67	NaN	30000-	5+ year
5	Kim	NLP	55	Delhi	6000^\$0	10+

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

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

```
Out[31]: 0    Datascience
         1    Testing
         2    Dataanalyst
         3    Analytics
         4    Statistics
         5    NLP
Name: Domain, dtype: object
```

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

```
In [33]: emp['Location']
```

```
Out[33]: 0      Mumbai
         1      Bangalore
         2        NaN
         3    Hyderabad
         4        NaN
         5       Delhi
Name: Location, dtype: object
```

```
In [34]: emp['Salary'] = emp['Salary'].str.replace(r'\W', ' ', regex=True)
```

```
In [35]: emp['Salary']
```

```
Out[35]: 0      5000
         1     10000
         2    15000
         3   20000
         4   30000
         5   60000
Name: Salary, dtype: object
```

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

```
In [37]: emp['Exp']
```

```
Out[37]: 0      2
         1      3
         2    4yrs
         3      NaN
         4   5year
         5      10
Name: Exp, dtype: object
```

 **Correct (for newer pandas versions):**

.str.extract(r'(\d+)') :- IS USED

```
In [38]: emp['Exp'] = emp['Exp'].str.extract(r'(\d+)')
```

```
In [39]: emp['Exp']
```

```
Out[39]: 0      2
         1      3
         2      4
         3      NaN
         4      5
         5     10
Name: Exp, dtype: object
```

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

Out[40]:

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

In [42]: `clean_data`

Out[42]:

	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

Clean Missing Value Treatment|

In [43]: `clean_data`

Out[43]:

	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 [44]: `clean_data.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6 entries, 0 to 5
Data columns (total 6 columns):
 #   Column      Non-Null Count  Dtype  
---  --          -----          object 
 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 [45]: `import numpy as np`

In [46]: `clean_data`

Out[46]:

	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 [47]: `clean_data.head()`

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	NaN	NaN	15000	4
3	Jane	Analytics	NaN	Hyderbad	20000	NaN
4	Uttam	Statistics	67	NaN	30000	5

In [48]: `clean_data.head(1)`

Out[48]:

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascience	34	Mumbai	5000	2

In [49]: `clean_data['Age']`

```
Out[49]: 0    34
         1    45
         2    NaN
         3    NaN
         4    67
         5    55
Name: Age, dtype: object
```

```
In [50]: clean_data['Age'] = clean_data['Age'].fillna(np.mean(pd.to_numeric(clean_data['Age'])))
```

```
In [51]: clean_data['Age']
```

```
Out[51]: 0    34
         1    45
         2    50.25
         3    50.25
         4    67
         5    55
Name: Age, dtype: object
```

```
In [52]: emp
```

	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 [53]: clean_data
```

	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	NaN
4	Uttam	Statistics	67	NaN	30000	5
5	Kim	NLP	55	Delhi	60000	10

```
In [54]: clean_data['Exp'] = clean_data['Exp'].fillna(np.mean(pd.to_numeric(clean_data['Exp'])))
```

```
In [55]: clean_data['Exp']
```

```
Out[55]: 0      2
         1      3
         2      4
         3    4.8
         4      5
         5     10
Name: Exp, dtype: object
```

```
In [56]: clean_data
```

```
Out[56]:   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 [61]: clean_data['Location'] = clean_data['Location'].fillna(clean_data['Location'].mode()
```

```
In [62]: clean_data['Location']
```

```
Out[62]: 0      Mumbai
         1  Bangalore
         2  Bangalore
         3  Hyderbad
         4  Bangalore
         5      Delhi
Name: Location, dtype: object
```

```
In [63]: clean_data
```

```
Out[63]:   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 [64]: 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
```

Try TO Convert Object To Cat And Int

```
In [65]: clean_data['Age'] = clean_data['Age'].astype(int)
```

```
In [66]: clean_data['Salary'] = clean_data['Salary'].astype(int)
```

```
In [67]: clean_data['Exp'] = clean_data['Exp'].astype(int)
```

```
In [68]: clean_data
```

```
Out[68]:
```

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascienc	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 [69]: 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 [75]: clean_data['Name'] = clean_data['Name'].astype('category')
```

```
In [76]: clean_data['Domain'] = clean_data['Domain'].astype('category')
```

```
In [77]: clean_data['Location'] = clean_data['Location'].astype('category')
```

```
In [78]: clean_data
```

```
Out[78]:
```

	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 [80]: 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 [81]: clean_data.to_csv('clean_data.csv')
```

```
In [82]: import os
os.getcwd()
```

```
Out[82]: 'C:\\\\Users\\\\DELL'
```

```
In [84]: clean_data.columns
```

```
Out[84]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')
```

```
In [85]: import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [86]: import warnings  
warnings.filterwarnings('ignore')
```

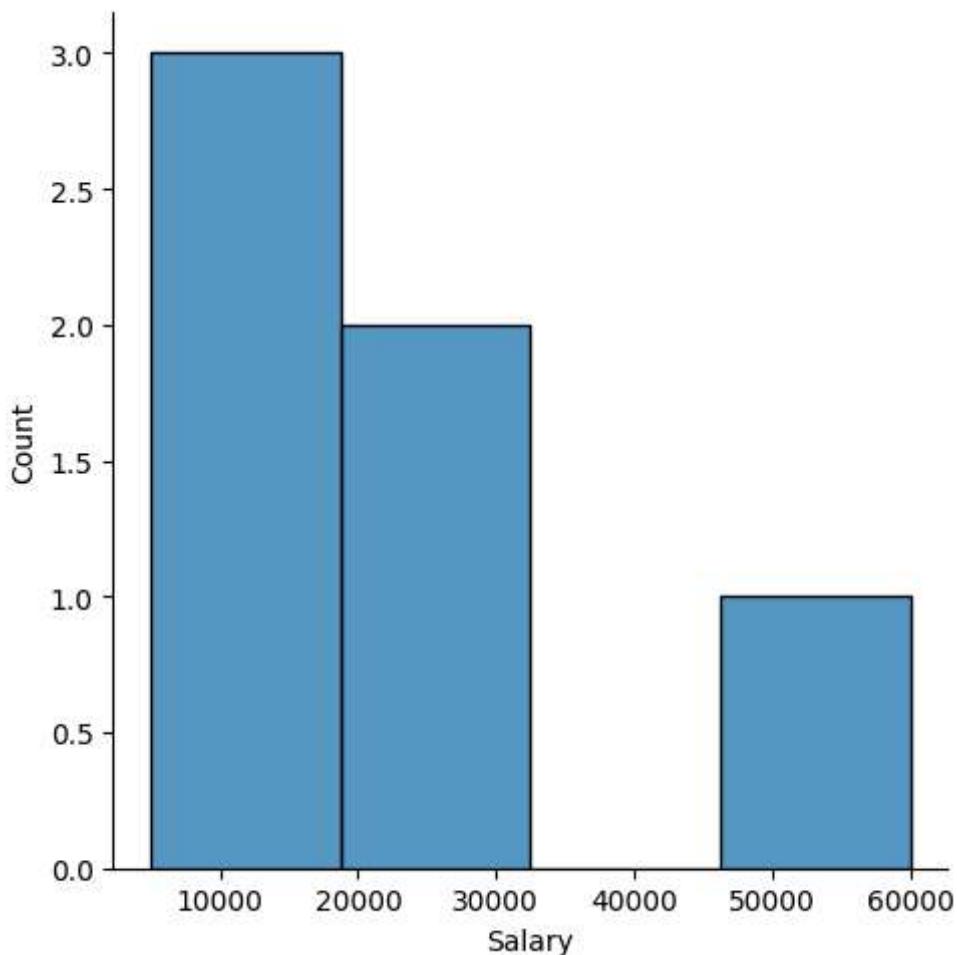
```
In [87]: clean_data
```

```
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 [88]: clean_data['Salary']
```

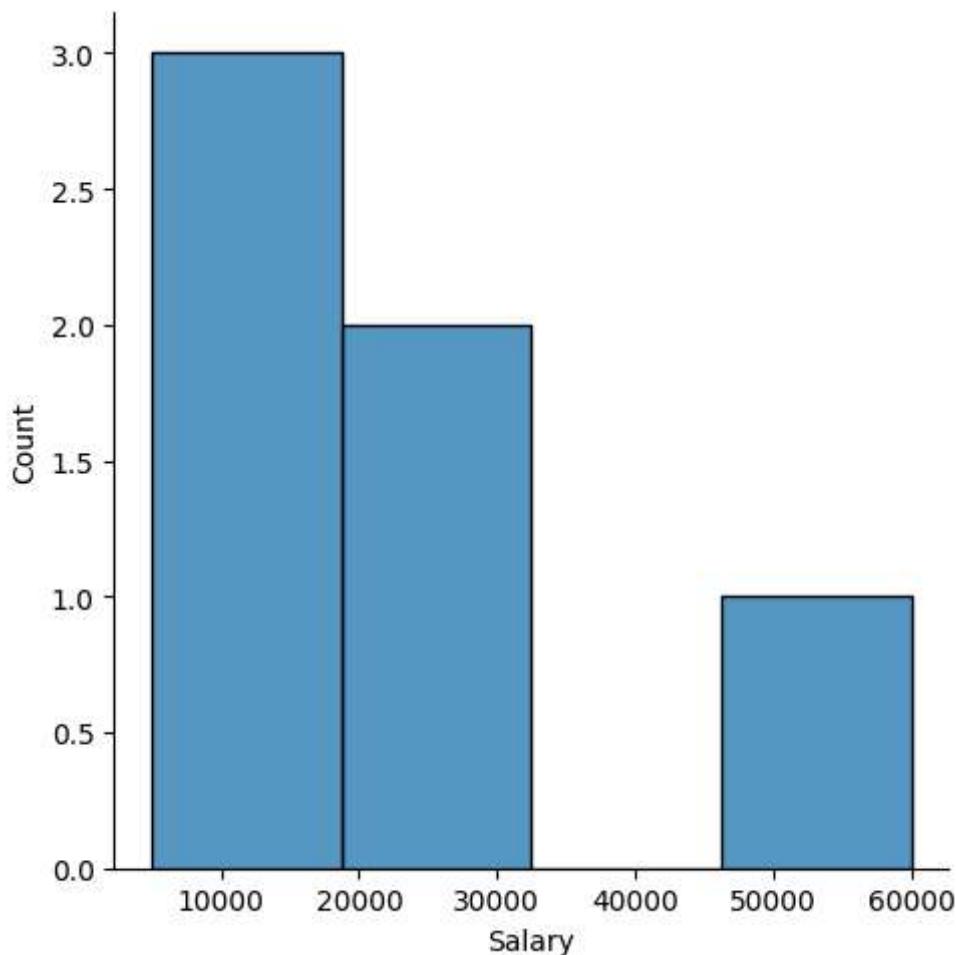
```
Out[88]: 0      5000  
1     10000  
2     15000  
3     20000  
4     30000  
5     60000  
Name: Salary, dtype: int64
```

```
In [89]: vis1 = sns.displot(clean_data['Salary'])
```

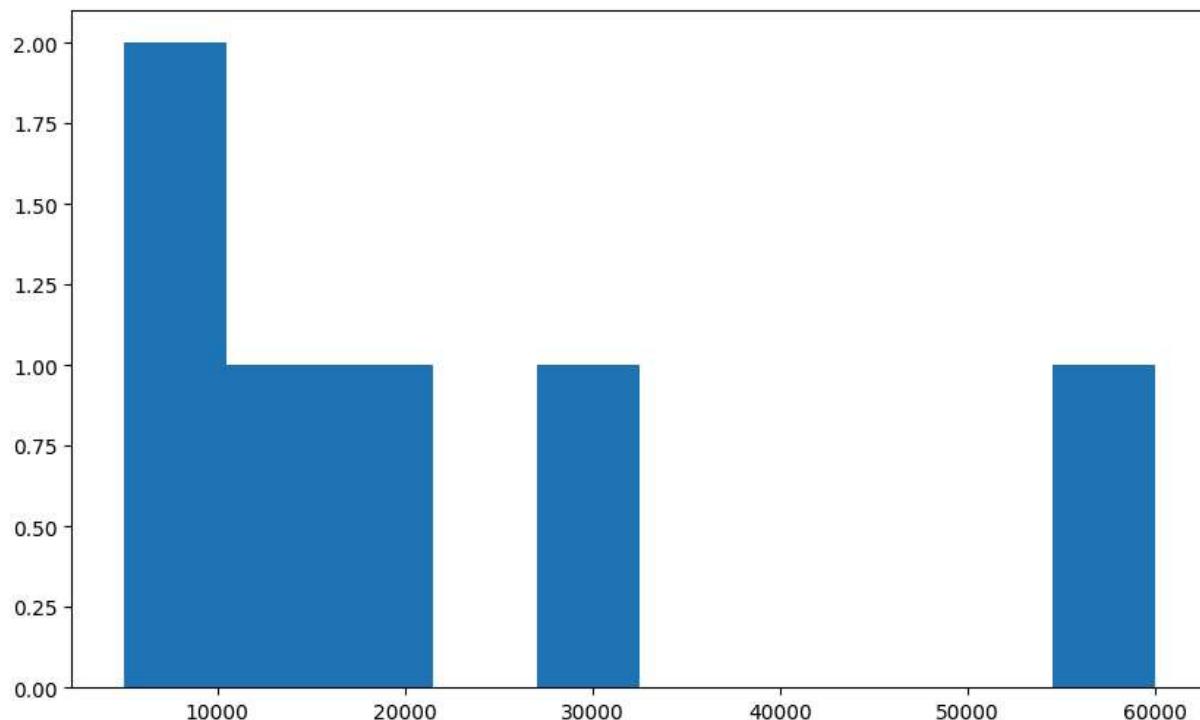


```
In [95]: plt.rcParams['figure.figsize'] = 10,6
```

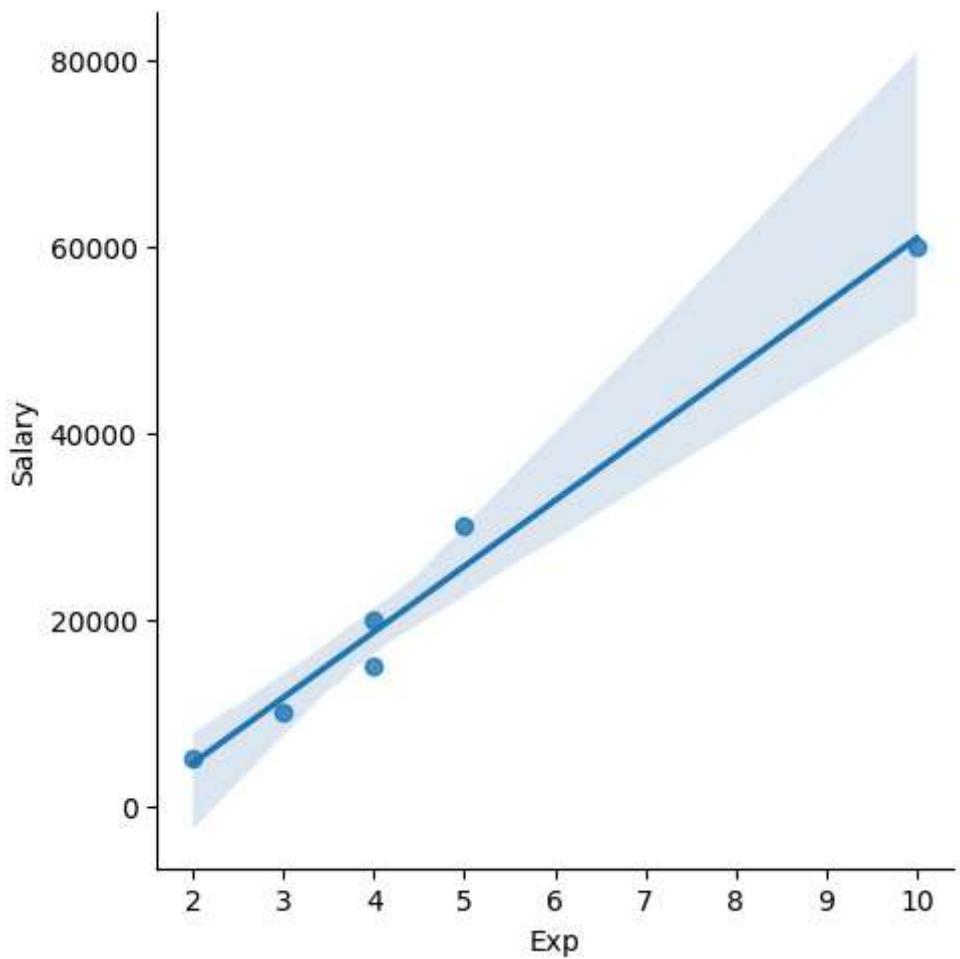
```
In [96]: vis1 = sns.displot(clean_data['Salary'])
```



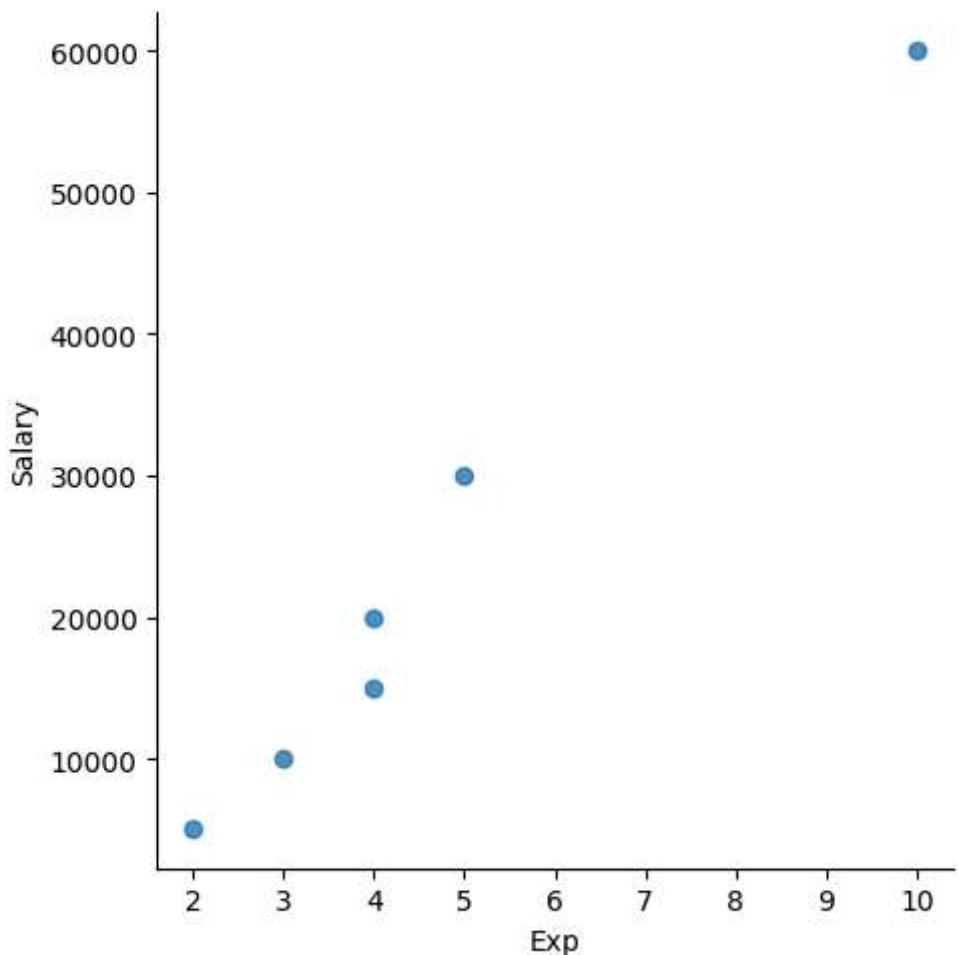
```
In [97]: vis2 = plt.hist(clean_data['Salary'])
```



```
In [99]: vis4 = sns.lmplot(data=clean_data,x = 'Exp', y = 'Salary')
```



```
In [105]: vis6 = sns.lmplot(data=clean_data,x = 'Exp' , y = 'Salary' , fit_reg = False)
```



```
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	Hyderabad	20000	4
4	Uttam	Statistics	67	Bangalore	30000	5
5	Kim	NLP	55	Delhi	60000	10

```
In [107...]: clean_data[:]
```

Out[107...]

	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 [108...]

clean_data[:2]

Out[108...]

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

In [109...]

clean_data[2:]

Out[109...]

	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 [110...]

clean_data[:]

Out[110...]

	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 [111...]

clean_data[0:1]

Out[111...]

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascience	34	Mumbai	5000	2

In [114...]

```
x_iv = clean_data.drop(['Salary'], axis=1)
```

In [115...]

```
x_iv
```

Out[115...]

	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 [116...]

```
x_iv.columns
```

Out[116...]

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

In [118...]

```
clean_data.columns
```

Out[118...]

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

In [120...]

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

In [121...]

```
y_dv
```

Out[121...]

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

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	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 [123...]

x_iv

Out[123...]

	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 [124...]

y_dv

Out[124...]

Salary

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

In [125...]

imputation = pd.get_dummies(clean_data)

In [126...]

imputation

Out[126...]

	Age	Salary	Exp	Name_Jane	Name_Kim	Name_Mike	Name_Teddy	Name_Umar	Name_Vincent
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	False
5	55	60000	10	False	True	False	False	False	False



In []: