

Raw Data To Clean Data Conversion Using Python EDA

```
In [2]: import pandas as pd
```

```
In [3]: emp=pd.read_excel(r"C:\Users\saniy\Downloads\Rawdata.xlsx")
```

```
In [4]: emp
```

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

```
In [5]: emp.columns
```

```
Out[5]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')
```

```
In [6]: emp.shape
```

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

```
In [7]: emp.head()
```

Out[7]:

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascienc#\$	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 [8]:

```
emp.tail()
```

Out[8]:

	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 [9]:

```
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 [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.isnull()
```

```
Out[11]:    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 [12]: emp.isnull().sum()
```

```
Out[12]: Name      0  
Domain     0  
Age       2  
Location   2  
Salary     0  
Exp        1  
dtype: int64
```

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

```
Out[13]: 0      Mike
          1    Teddy^
          2     Uma#r
          3      Jane
          4    Uttam*
          5      Kim
Name: Name, dtype: object
```

```
In [14]: emp['Name'] = emp['Name'].str.replace(r'\W', ' ', regex = True)
```

```
In [15]: emp['Name']
```

```
Out[15]: 0      Mike
          1    Teddy
          2     Umar
          3      Jane
          4    Uttam
          5      Kim
Name: Name, dtype: object
```

```
In [21]: emp['Domain'] = emp['Domain'].str.replace(r'\W', ' ', regex=True)
```

```
In [24]: emp['Domain']
```

```
Out[24]: 0    Datascience
          1      Testing
          2   Dataanalyst
          3    Analytics
          4   Statistics
          5        NLP
Name: Domain, dtype: object
```

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

```
In [27]: emp['Age']
```

```
Out[27]: 0    34
         1    45
         2    NaN
         3    NaN
         4    67
         5    55
Name: Age, dtype: object
```

```
In [32]: emp['Age'] = emp['Age'].str.extract('(\d+)')
<>:1: SyntaxWarning: invalid escape sequence '\d'
<>:1: SyntaxWarning: invalid escape sequence '\d'
C:\Users\sanvi\AppData\Local\Temp\ipykernel_21344\1884116463.py:1: SyntaxWarning: invalid escape sequence '\d'
emp['Age'] = emp['Age'].str.extract('(\d+)')
```

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

```
Out[33]: 0    34
         1    45
         2    NaN
         3    NaN
         4    67
         5    55
Name: Age, dtype: object
```

```
In [47]: emp['Location'] = emp['Location'].str.replace(r'\W', '')
```

```
In [48]: emp['Location']
```

```
Out[48]: 0      Mumbai
         1    Bangalore
         2        NaN
         3   Hyderabad
         4        NaN
         5      Delhi
Name: Location, dtype: object
```

```
In [34]: emp
```

Out[34]:

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascience	34	Mumbai	5^00#0	2+
1	Teddy	Testing	45	Bangalore	10%0000	<3
2	Umar	Dataanalyst	NaN	NaN	1\$5%0000	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 [80]: `emp['Salary']`

Out[80]:

```
0    5^00#0
1    10%0000
2    1$5%000
3    2000^0
4    30000-
5    6000^$0
Name: Salary, dtype: object
```

In [88]: `emp['Salary'] = emp['Salary'].str.replace(r'\W', ' ', regex=True)`

In [89]: `emp['Salary']`

Out[89]:

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

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

```
Out[90]: 0    2  
1    3  
2    4  
3    NaN  
4    5  
5    10  
Name: Exp, dtype: object
```

```
In [91]: emp['Exp']=emp['Exp'].str.extract('(\d+)')
```

```
<>:1: SyntaxWarning: invalid escape sequence '\d'  
<>:1: SyntaxWarning: invalid escape sequence '\d'  
C:\Users\sanvi\AppData\Local\Temp\ipykernel_21344\1466635560.py:1: SyntaxWarning: invalid escape sequence '\d'  
  emp['Exp']=emp['Exp'].str.extract('(\d+)')
```

```
In [92]: emp['Exp']
```

```
Out[92]: 0    2  
1    3  
2    4  
3    NaN  
4    5  
5    10  
Name: Exp, dtype: object
```

```
In [93]: emp
```

```
Out[93]:   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 [94]: clean_data = emp.copy()
```

till now we have raw data we use regex to clean the data and removed all noise characted from the dataset you can also work in same things in sql query as well

- missing values treatment for numerical data

```
In [95]: clean_data
```

```
Out[95]:
```

	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 [96]: clean_data['Age']
```

```
Out[96]:
```

0	34
1	45
2	NaN
3	NaN
4	67
5	55

Name: Age, dtype: object

```
In [97]: import numpy as np
```

```
In [98]: clean_data['Age'] = clean_data['Age'].fillna(np.mean(pd.to_numeric(clean_data['Age'])))
```

```
In [99]: clean_data['Age']
```

```
Out[99]: 0    34  
1    45  
2    50.25  
3    50.25  
4    67  
5    55  
Name: Age, dtype: object
```

```
In [100... clean_data['Exp']
```

```
Out[100... 0    2  
1    3  
2    4  
3    NaN  
4    5  
5    10  
Name: Exp, dtype: object
```

```
In [101... clean_data['Exp'] = clean_data['Exp'].fillna(np.mean(pd.to_numeric(clean_data['Exp'])))
```

```
In [102... clean_data['Exp']
```

```
Out[102... 0    2  
1    3  
2    4  
3    4.8  
4    5  
5    10  
Name: Exp, dtype: object
```

```
In [103... clean_data
```

```
Out[103...]
```

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascienc	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 [104...]
```

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

```
Out[104...]
```

```
2
```

```
In [105...]
```

```
clean_data['Location']
```

```
Out[105...]
```

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

```
In [106...]
```

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

```
In [107...]
```

```
clean_data['Location']
```

```
Out[107...]
```

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

```
In [108... clean_data
```

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

```
In [111... 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 [112... clean_data['Salary'] = clean_data['Salary'].astype(int)
clean_data['Exp'] = clean_data['Exp'].astype(int)
```

```
In [113... 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 [114... 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 [115... 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

Out[116...

	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 [117... clean_data.to_csv('clean_data.csv')

In [119... import os
os.getcwd()

Out[119... 'C:\\\\Users\\\\saniy'

In [120... clean_data

Out[120...]

	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

EDA TECHNIQUE LETS APPLY

In [122...]

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

In [123...]

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

In [125...]

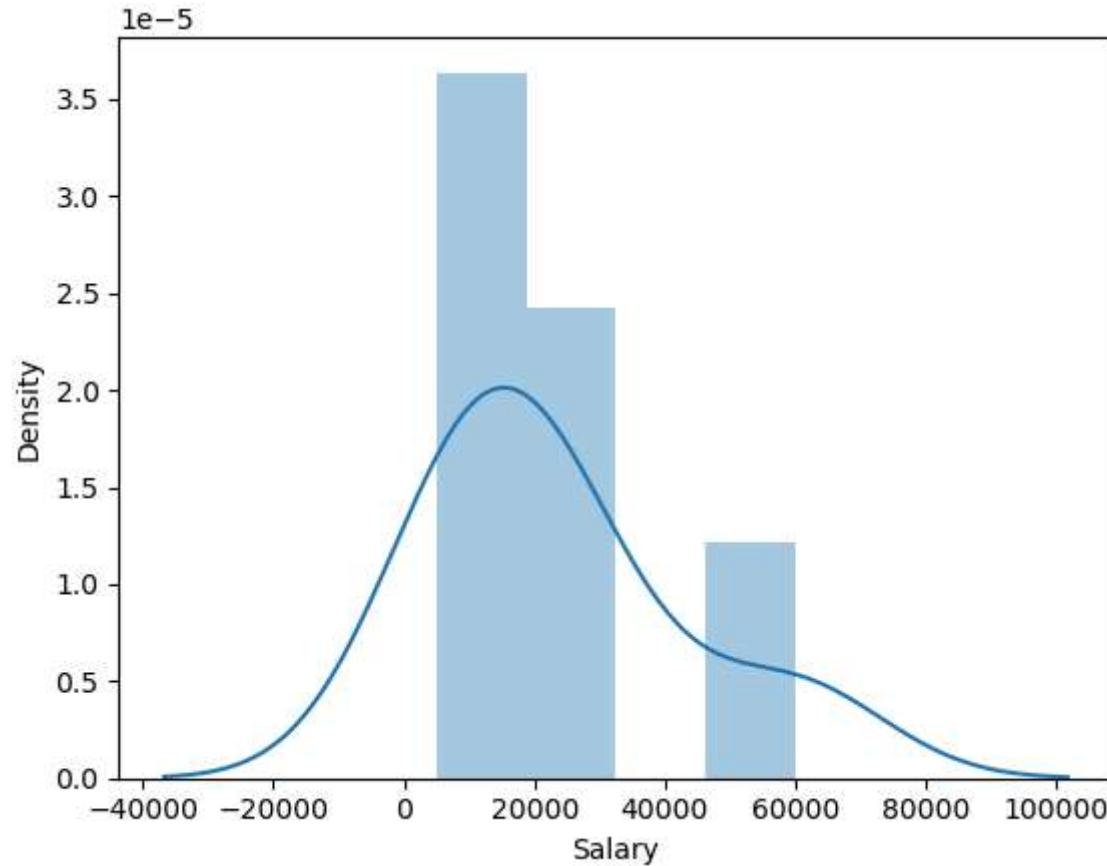
```
clean_data['Salary']
```

Out[125...]

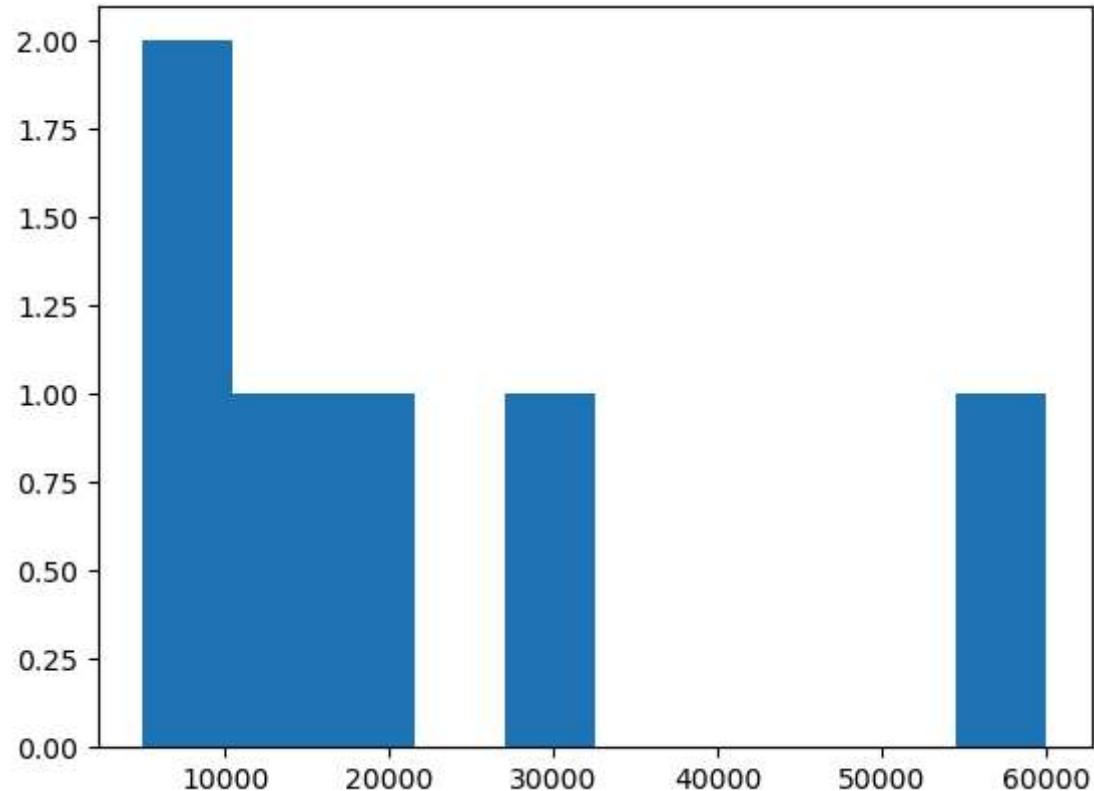
```
0    5000
1   10000
2   15000
3   20000
4   30000
5   60000
Name: Salary, dtype: int32
```

In [126...]

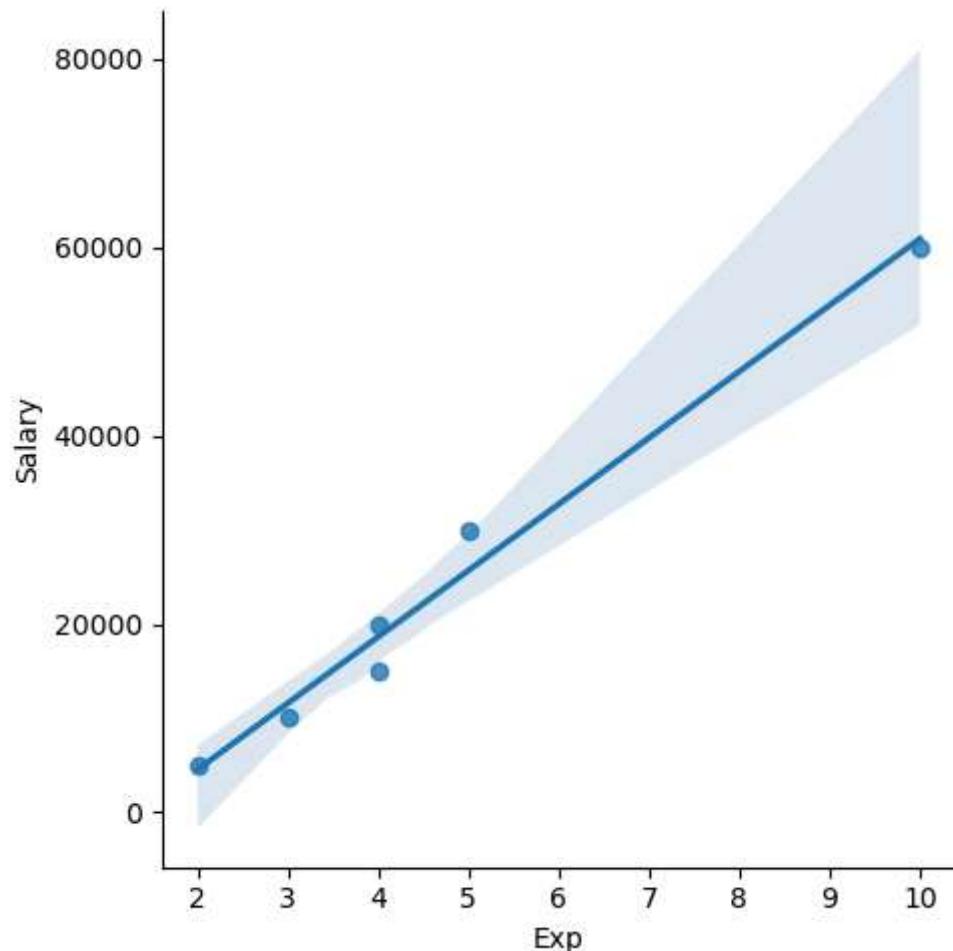
```
vs1= sns.distplot(clean_data['Salary'])
```



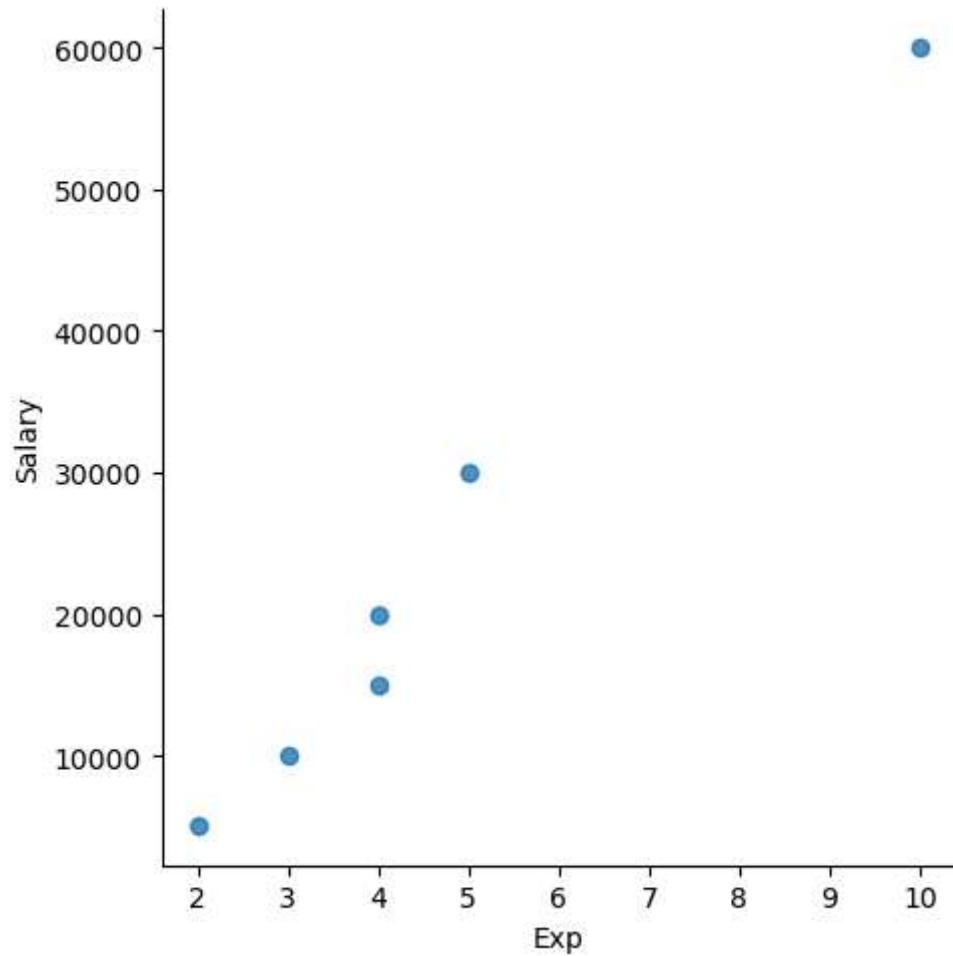
```
In [127]: vis2 = plt.hist(clean_data['Salary'])
```



```
In [128]: vis4 = sns.lmplot(data=clean_data,x = 'Exp', y='Salary')
```



```
In [129]: vis5 = sns.lmplot(data=clean_data,x = 'Exp', y='Salary', fit_reg = False)
```



```
In [131]: clean_data[:]
```

```
Out[131...]
```

	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 [132...]
```

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

```
Out[132...]
```

	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 [133...]
```

```
clean_data[::-1]
```

```
Out[133...]
```

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

```
In [134...]
```

```
clean_data.columns
```

```
Out[134]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')
```

```
In [135]: x_iv=clean_data[['Name','Domain','Age','Location','Exp']]
```

```
In [136]: x_iv
```

```
Out[136]:
```

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

```
In [138]: y_dv
```

```
Out[138]:
```

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

```
In [139]: emp
```

Out[139...]

	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 [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	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 [141...]

```
x_iv
```

Out[141...]

	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 [142...]

y_dv

Out[142...]

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

In [143...]

clean_data

```
Out[143...]
```

	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 [144...]
```

```
imputation = pd.get_dummies(clean_data)
```

```
In [145...]
```

```
imputation
```

```
Out[145...]
```

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

```
In [146...]
```

```
clean_data
```

Out[146...]

	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 [147...]

imputation

Out[147...]

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

raw data with lot of regex, missing, uncleandata

regex, clean

fill missing numerical & cateigroical

clean_dataset (data cleaning)

3 month - 5 montoutlier treatement,

univati, bivariate, corelation

split the data into x_i.v & y_dv

impute cateogrica data to numericaleda part complete

In []:

In []: