

# Raw Data to clean Data conversion using python EDA

```
In [200...]: import pandas as pd
import warnings
warnings.filterwarnings('ignore')

In [202...]: emp=pd.read_excel(r"D:\DA ALL NOTES\DAY24\Rawdata.xlsx")

In [204...]: emp
```

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascience#\$	34 years	Mumbai	5^00#0	2+
1	Teddy^	Testing	45' yr	Bangalore	10%6000	<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 [206...]: emp.shape

Out[206...]: (6, 6)

In [208...]: len(emp)

Out[208...]: 6

In [210...]: emp.columns

Out[210...]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')

In [212...]: 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 [214...]: emp.head()
```

```
Out[214...]
```

	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 [216...]
```

```
emp.tail()
```

```
Out[216...]
```

	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 [218...]
```

```
emp.describe()
```

```
Out[218...]
```

	Name	Domain	Age	Location	Salary	Exp
count	6	6	4	4	6	5
unique	6	6	4	4	6	5
top	Mike	Datascience#\$	34 years	Mumbai	5^00#0	2+
freq	1	1	1	1	1	1

```
In [220...]
```

```
emp.isnull().sum()
```

```
Out[220...]
```

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

```
In [222...]
```

```
emp['Name']
```

```
Out[222...]
```

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

```
In [224...]
```

```
emp['Domain']
```

```
Out[224... 0    Datascienc#$  
1          Testing  
2  Dataanalyst^^#  
3      Ana^^lytics  
4      Statistics  
5          NLP  
Name: Domain, dtype: object
```

```
In [226... emp[['Name', 'Domain']]
```

```
Out[226...   Name      Domain  
0   Mike    Datascienc#$  
1  Teddy        Testing  
2  Uma#r  Dataanalyst^^#  
3   Jane    Ana^^lytics  
4  Uttam*      Statistics  
5     Kim        NLP
```

```
In [228... emp['Name']=emp['Name'].str.replace(r'\W', ' ', regex=True)
```

```
In [230... emp['Name']
```

```
Out[230... 0    Mike  
1  Teddy  
2   Umar  
3   Jane  
4  Uttam  
5     Kim  
Name: Name, dtype: object
```

```
In [232... emp['Domain']=emp['Domain'].str.replace(r'\W', ' ', regex=True)
```

```
In [234... emp['Domain']
```

```
Out[234... 0    Datascienc  
1      Testing  
2  Dataanalyst  
3    Analytics  
4  Statistics  
5      NLP  
Name: Domain, dtype: object
```

```
In [236... emp['Name']=emp['Name'].str.replace(r'\W', ' ', regex=True)
```

```
In [238... emp
```

```
Out[238...]
```

	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 [240...]
```

```
emp['Age']=emp['Age'].str.replace(r'\W',' ',regex=True)  
emp['Age']=emp['Age'].str.extract('(\d+)')
```

```
In [242...]
```

```
emp['Age']
```

```
Out[242...]
```

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

```
In [244...]
```

```
emp
```

```
Out[244...]
```

	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 [246...]
```

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

```
In [248...]
```

```
emp['Location']
```

```
Out[248...]
```

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

```
In [250...]
```

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

```
In [252...]
```

```
emp['Salary']
```

```

Out[252... 0      5000
          1      10000
          2      15000
          3      20000
          4      30000
          5      60000
Name: Salary, dtype: object

In [254... emp['Exp']=emp['Exp'].str.replace(r'\W',' ',regex=True)

In [256... emp['Exp']

Out[256... 0      2
          1      3
          2      4yrs
          3      NaN
          4      5year
          5      10
Name: Exp, dtype: object

In [258... emp['Exp']=emp['Exp'].str.extract('(\d+)')

In [260... emp['Exp']

Out[260... 0      2
          1      3
          2      4
          3      NaN
          4      5
          5      10
Name: Exp, dtype: object

In [262... emp

Out[262...    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 [264... clean_data=emp.copy()

```

**Till now we have raw data we use replace to clean data and remove all noise character from the dataset .**

**you can also work in same thing in sql query as well**

## -Missing values treatment for numerical data

In [269...]: clean\_data

```
Out[269...]:
```

	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 [271...]: clean\_data['Age']

```
Out[271...]:
```

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

Name: Age, dtype: object

In [273...]: import numpy as np

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

In [277...]: clean\_data['Age']

```
Out[277...]:
```

0	34
1	45
2	50.25
3	50.25
4	67
5	55

Name: Age, dtype: object

In [279...]: clean\_data['Exp']

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

```
In [281... clean_data['Exp']=clean_data['Exp'].fillna(np.mean(pd.to_numeric(clean_data['Exp'])))
```

```
In [283... clean_data['Exp']
```

```
Out[283... 0      2  
1      3  
2      4  
3    4.8  
4      5  
5     10  
Name: Exp, 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  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 [287... clean_data['Location'].isna().sum()
```

```
Out[287... 2
```

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

```
Out[289... 0      Mumbai  
1      Bangalore  
2      NaN  
3      Hyderabad  
4      NaN  
5      Delhi  
Name: Location, dtype: object
```

```
In [291... clean_data['Location']=clean_data['Location'].fillna(clean_data['Location'].mode())
```

```
In [293... clean_data['Location']
```

```
Out[293...]: 0      Mumbai  
1      Bangalore  
2      Bangalore  
3      Hyderabad  
4      Bangalore  
5      Delhi  
Name: Location, dtype: object
```

```
In [295...]: clean_data
```

```
Out[295...]:
```

	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 [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      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 [299...]: clean_data.isnull().sum()
```

```
Out[299...]:
```

Name	0
Domain	0
Age	0
Location	0
Salary	0
Exp	0

dtype: int64

```
In [301...]: clean_data.columns
```

```
Out[301...]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')
```

```
In [303...]: clean_data['Salary']=clean_data['Salary'].astype(int)  
clean_data['Exp']=clean_data['Exp'].astype(int)  
clean_data['Age']=clean_data['Age'].astype(int)
```

```
In [305...]: clean_data
```

```
Out[305...]
```

	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 [307...]
```

```
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 [309...]
```

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

```
In [311...]
```

```
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 [313...]
```

```
clean_data
```

```
Out[313...]
```

	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...]
```

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

```
In [317...]
```

```
import os  
os.getcwd()
```

```
Out[317...]
```

```
'C:\\\\Users\\\\subhra kanta sahoo'
```

```
In [319...]
```

```
clean_data
```

```
Out[319...]
```

	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 [322...]
```

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

```
In [324...]
```

```
clean_data
```

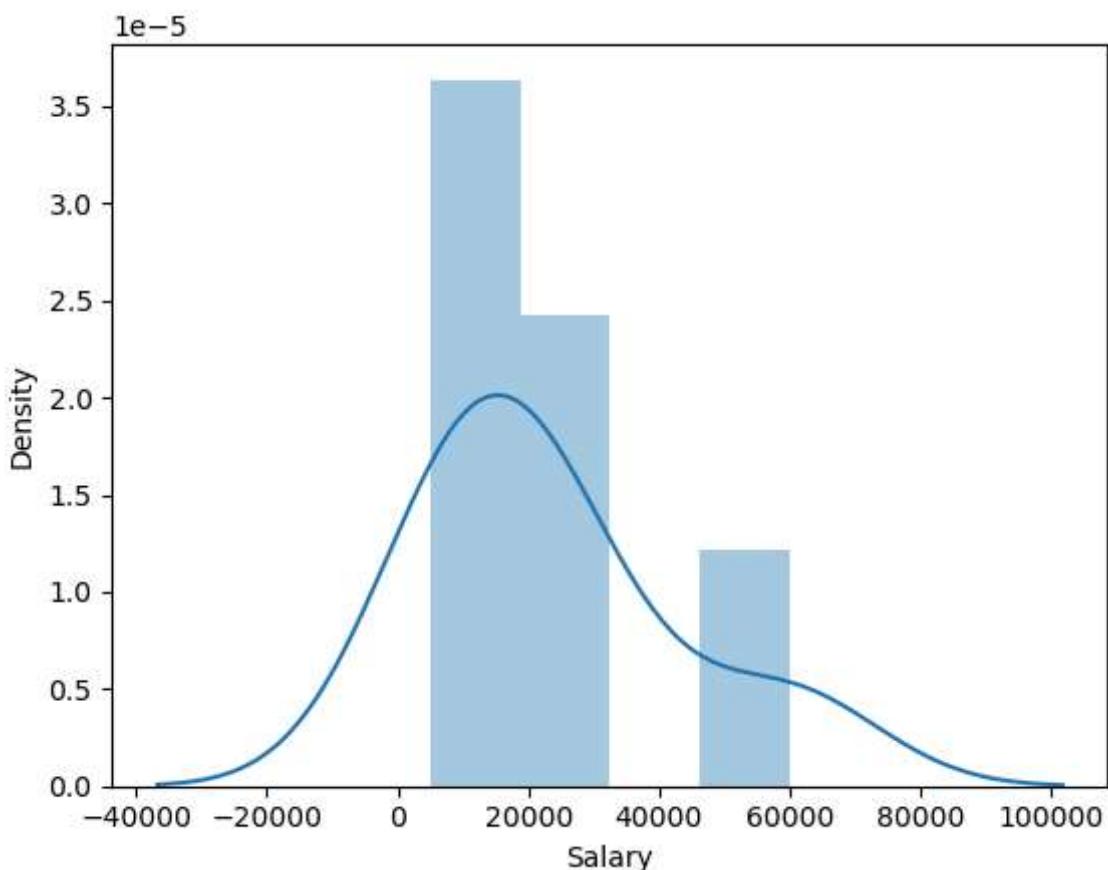
```
Out[324...]
```

	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

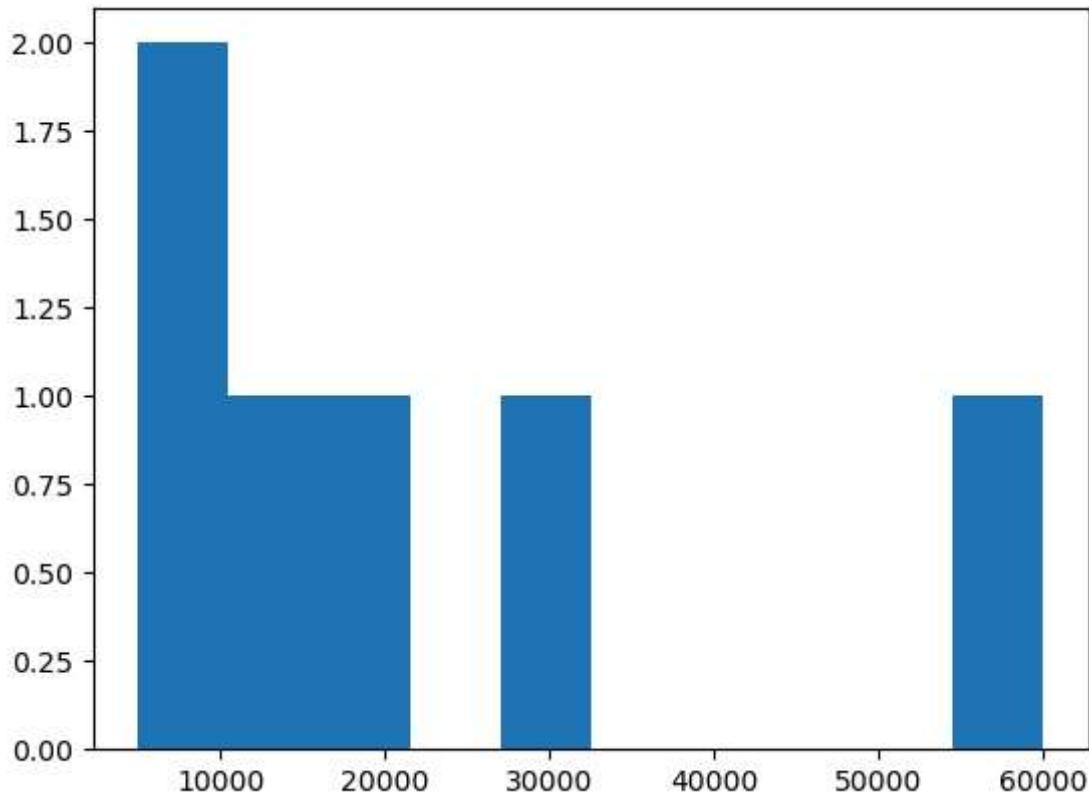
```
In [326...]: clean_data['Salary']
```

```
Out[326...]: 0      5000
 1     10000
 2    15000
 3   20000
 4   30000
 5   60000
Name: Salary, dtype: int32
```

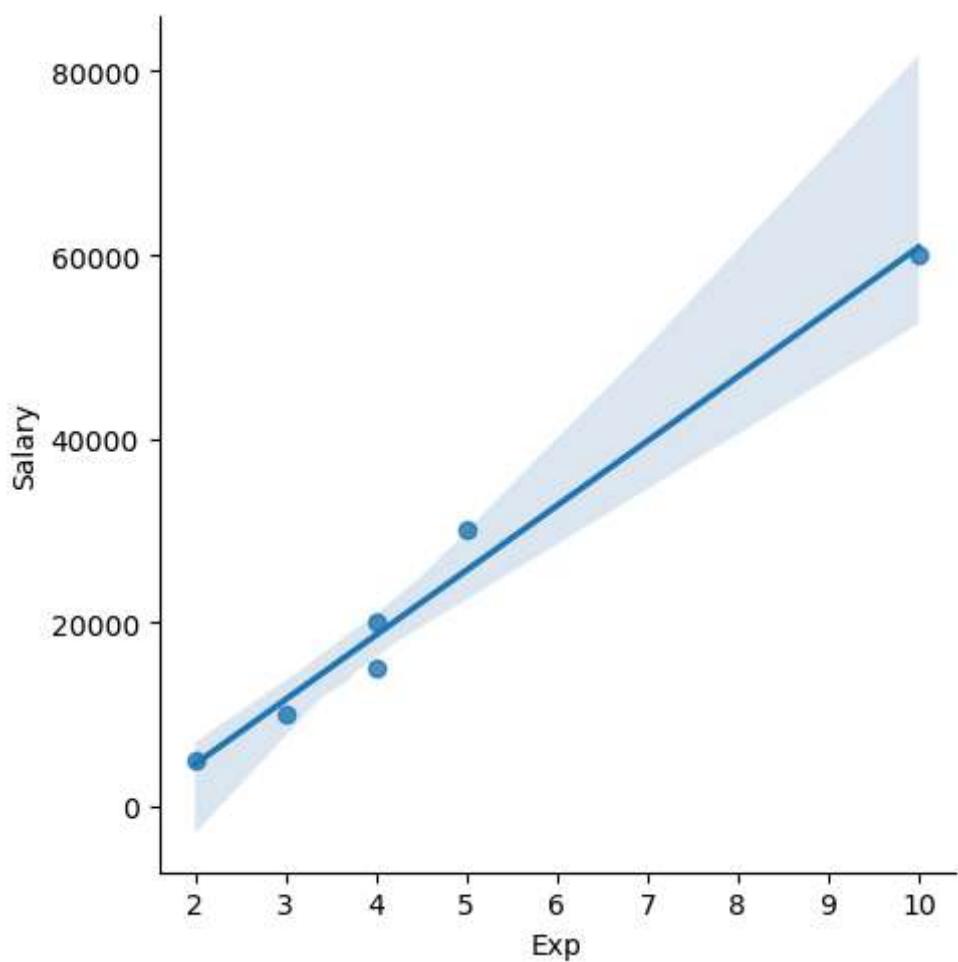
```
In [328...]: vis1=sns.distplot(clean_data['Salary'])
```



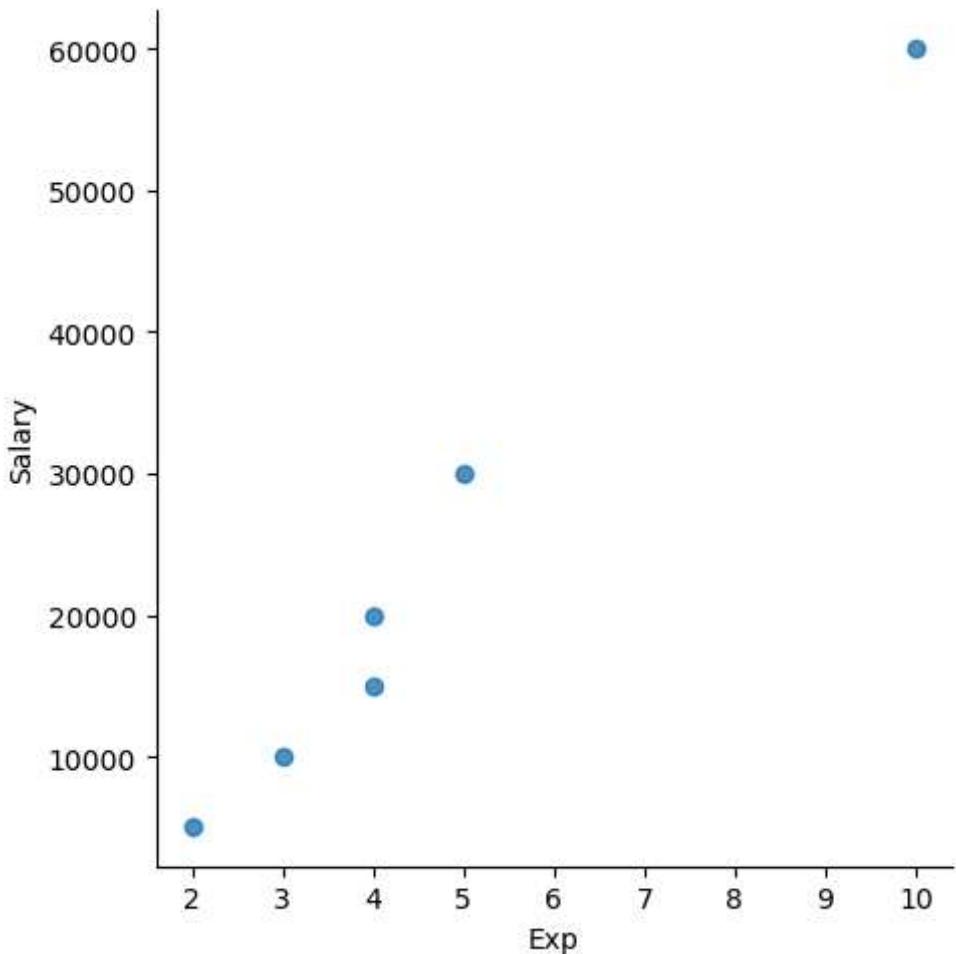
```
In [330...]: vis2=plt.hist(clean_data['Salary'])
```



```
In [332]: vis3 = sns.lmplot(data=clean_data,x='Exp',y='Salary')
```



```
In [334]: vis4=sns.lmplot(data=clean_data,x='Exp',y='Salary',fit_reg=False)
```



```
In [336...]: clean_data[:]
```

```
Out[336...]:
```

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

```
Out[338...]:
```

	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 [340...]: clean_data[::-1]
```

```
Out[340...]
```

	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 [342...]
```

```
clean_data.columns
```

```
Out[342...]
```

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

```
In [344...]
```

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

```
In [346...]
```

```
X_iv
```

```
Out[346...]
```

	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 [352...]
```

```
Y_dv=clean_data[['Salary']]
```

```
In [354...]
```

```
Y_dv
```

```
Out[354...]
```

```
Salary
```

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

```
In [356...]
```

```
emp
```

```
Out[356...]
```

	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 [358...]
```

```
clean_data
```

```
Out[358...]
```

	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 [360...]
```

```
X_iv
```

```
Out[360...]
```

	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 [362...]
```

```
Y_dv
```

```
Out[362...]
```

**Salary**

<b>0</b>	5000
<b>1</b>	10000
<b>2</b>	15000
<b>3</b>	20000
<b>4</b>	30000
<b>5</b>	60000

```
In [364...]
```

```
clean_data
```

```
Out[364...]
```

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

```
In [382...]
```

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

```
In [384...]
```

```
imputation
```

```
Out[384...]
```

	Age	Salary	Exp	Name_Jane	Name_Kim	Name_Mike	Name_Teddy	Name_Umar
<b>0</b>	34	5000	2	0	0	1	0	0
<b>1</b>	45	10000	3	0	0	0	1	0
<b>2</b>	50	15000	4	0	0	0	0	1
<b>3</b>	50	20000	4	1	0	0	0	0
<b>4</b>	67	30000	5	0	0	0	0	0
<b>5</b>	55	60000	10	0	1	0	0	0



```
In [386...]
```

```
imputation1 = pd.get_dummies(clean_data)
```

```
In [388...]
```

```
imputation1
```

Out[388...]

	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 [390...]

clean\_data

Out[390...]

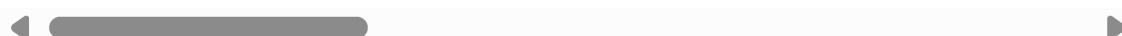
	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascience	34	Mumbai	5000	2
1	Teddy	Testing	45	Bangalore	10000	3
2	Umar	Dataanalyst	50	Bangalore	15000	4
3	Jane	Analytics	50	Hyderbad	20000	4
4	Uttam	Statistics	67	Bangalore	30000	5
5	Kim	NLP	55	Delhi	60000	10

In [392...]

imputation

Out[392...]

	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 [394...]

imputation1

Out[394...]

	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

raw data with lot of regex, missing, uncleandata

regex, clean

fill missing numerical & cateigroica

clean\_dataset ( data cleaning) 3 month - 5mont

outlier treatment, univati, bivariate, corelation

split the data into x\_iv & y\_dv

impute cateogrica data to numerical

eda part complete

Next step

- we splitn x\_iv -- x\_train, x\_test
- we split y\_dv -- y\_train, y\_test
- build the ml model with x\_train & y\_train

COMPLETED