

20th May

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

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

```
Out[2]: '2.2.2'
```

```
In [3]: # !pip install---upgrade openpyxl
```

```
In [4]: emp=pd.read_excel(r"C:\Users\J. Rajesh\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	Hyderbad	2000^0	NaN
4	Uttam*	Statistics	67-yr	NaN	30000-	5+ year
5	Kim	NLP	55yr	Delhi	6000^\$0	10+

```
In [6]: id(emp)
```

```
Out[6]: 2388361076400
```

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

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

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

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

```
In [9]: emp.head()
```

```
Out[9]:
```

	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 [10]: `emp.tail()`

Out[10]:

	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 [11]: `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 [12]: `emp`

Out[12]:

	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 [13]: `emp.isnull()` *# it will give the true in the missing values*

Out[13]:

	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 [14]: `emp.isnull().sum()`

Out[14]:

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

In [15]: `emp.columns`

Out[15]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')

In [16]: `emp`

Out[16]:

	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 or Data Cleansing

In [17]: `emp`

Out[17]:

	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 [18]: emp['Name']

Out[18]:

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

In [19]: emp['Name']=emp['Name'].str.replace(r'\W','',regex=True) # to remove all the sp

In [20]: emp['Name']

Out[20]:

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

In [21]: emp

Out[21]:

	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 [22]: emp['Domain']

```
Out[22]: 0    Datascience#$
        1      Testing
        2  Dataanalyst^^#
        3    Ana^^lytics
        4    Statistics
        5        NLP
        Name: Domain, dtype: object
```

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

```
Out[23]: 0    Datascience
        1      Testing
        2    Dataanalyst
        3      Analytics
        4    Statistics
        5        NLP
        Name: Domain, dtype: object
```

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

```
Out[24]: 0    34years
        1    45yr
        2      NaN
        3      NaN
        4    67yr
        5    55yr
        Name: Age, dtype: object
```

```
In [25]: emp['Age']=emp['Age'].str.extract('(\d+)') #to remove the irregular format
        emp['Age']
```

```
<>:1: SyntaxWarning: invalid escape sequence '\d'
<>:1: SyntaxWarning: invalid escape sequence '\d'
C:\Users\J. Rajesh\AppData\Local\Temp\ipykernel_13564\1973738011.py:1: SyntaxWarning: invalid escape sequence '\d'
    emp['Age']=emp['Age'].str.extract('(\d+)') #to remove the irregular format
```

```
Out[25]: 0    34
        1    45
        2    NaN
        3    NaN
        4    67
        5    55
        Name: Age, dtype: object
```

```
In [26]: emp
```

Out[26]:

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

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

Out[31]:

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

In [32]: `emp`

Out[32]:

	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 [33]: `emp['Exp']`

Out[33]:

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

In [34]: `emp['Exp']=emp['Exp'].str.extract('(\d+)')`
`emp['Exp']`

```
<>:1: SyntaxWarning: invalid escape sequence '\d'
<>:1: SyntaxWarning: invalid escape sequence '\d'
C:\Users\J. Rajesh\AppData\Local\Temp\ipykernel_13564\2863867557.py:1: SyntaxWarn
ing: invalid escape sequence '\d'
emp['Exp']=emp['Exp'].str.extract('(\d+)')
```

```
Out[34]: 0      2
         1      3
         2      4
         3      NaN
         4      5
         5     10
         Name: Exp, dtype: object
```

```
In [35]: emp
```

```
Out[35]:
```

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

```
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	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

%% md

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

```
In [38]: clean_data['Age']
```

```
Out[38]: 0      34
         1      45
         2     NaN
         3     NaN
         4      67
         5      55
         Name: Age, dtype: object
```

```
In [40]: import numpy as np
```

```
In [41]: clean_data['Age']=clean_data['Age'].fillna(np.mean(pd.to_numeric(clean_data['Age
```

```
In [42]: clean_data['Age']
```

```
Out[42]: 0      34
         1      45
         2    50.25
         3    50.25
         4      67
         5      55
         Name: Age, dtype: object
```

```
In [43]: clean_data['Exp']
```

```
Out[43]: 0      2
         1      3
         2      4
         3     NaN
         4      5
         5     10
         Name: Exp, dtype: object
```

```
In [44]: clean_data['Exp']=clean_data['Exp'].fillna(np.mean(pd.to_numeric(clean_data['Exp
```

```
Out[44]: 0      2
         1      3
         2      4
         3    4.8
         4      5
         5     10
         Name: Exp, dtype: object
```

```
In [45]: clean_data
```

```
Out[45]:
```

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

```
Out[46]: 2
```

```
In [47]: clean_data['Location']
```

```
Out[47]: 0      Mumbai
1    Bangalore
2         NaN
3    Hyderabad
4         NaN
5        Delhi
Name: Location, dtype: object
```

```
In [48]: clean_data['Location'] = clean_data['Location'].fillna(clean_data['Location'].mode()[0])
```

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

```
In [49]: clean_data
```

```
Out[49]:
```

	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 [50]: 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 [51]: clean_data['Age'] = clean_data['Age'].astype(int)
```

In [52]: `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 [53]: `clean_data['Salary'] = clean_data['Salary'].astype(int)`
`clean_data['Exp'] = clean_data['Exp'].astype(int)`

In [54]: `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 [56]: `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 [57]: `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 [58]: `clean_data`

Out[58]:

	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 [59]: `clean_data.to_csv('clean_data.csv')`

In [60]: `import os`
`os.getcwd()`

Out[60]: 'C:\\Users\\J. Rajesh'

In [61]: `clean_data`

Out[61]:

	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

LETS APPLY EDA TECHNIQUE

In [62]: `import matplotlib.pyplot as plt # for visualization`
`import seaborn as sns`

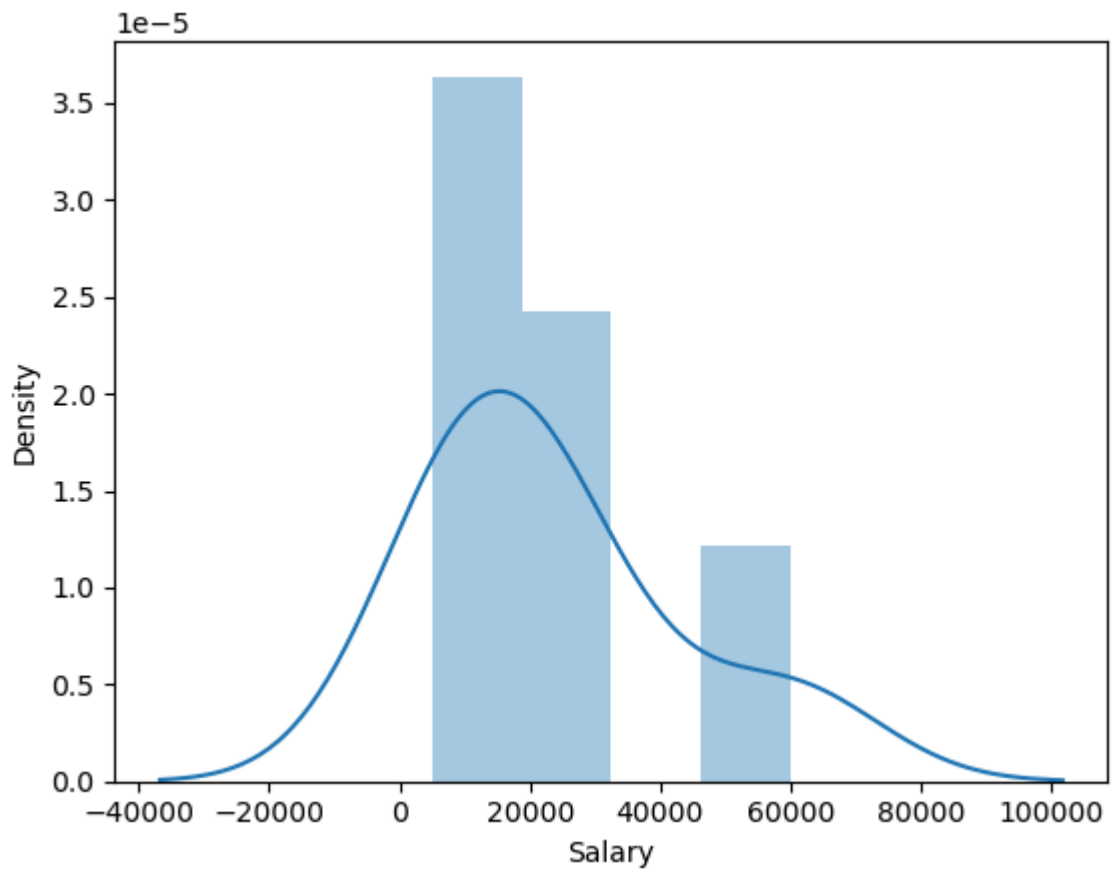
In [63]: `import warnings`
`warnings.filterwarnings('ignore')`

In [64]: `clean_data['Salary']`

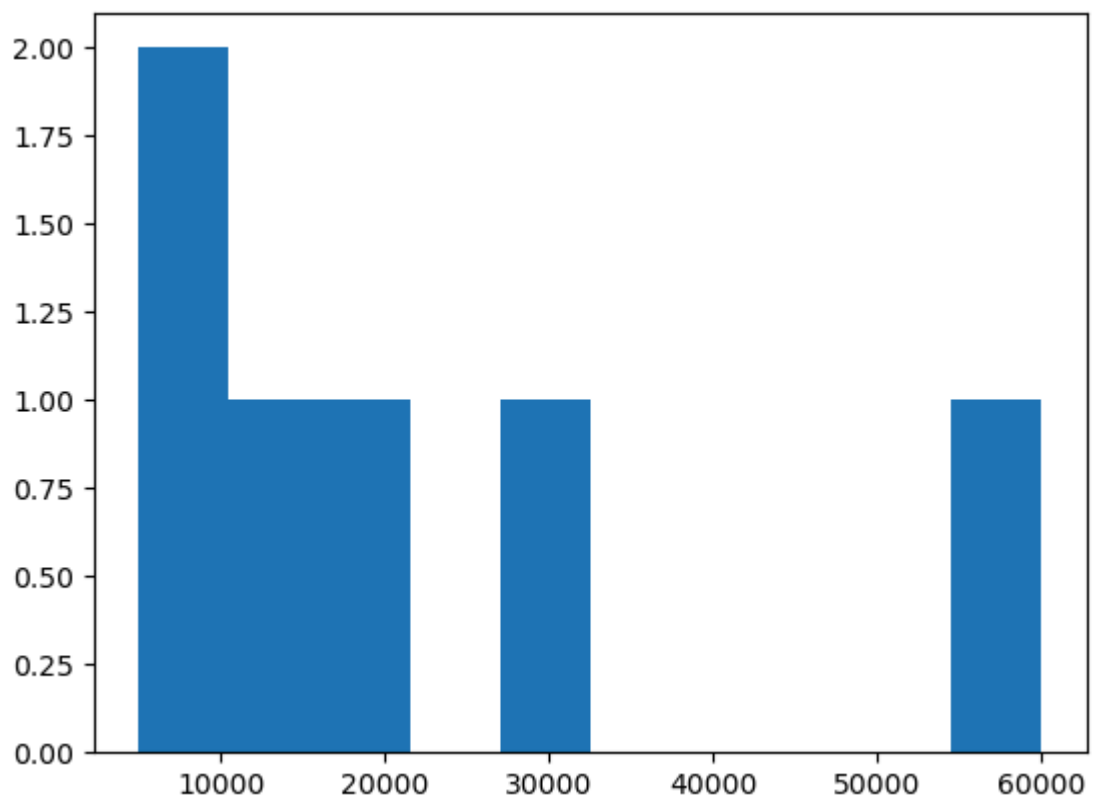
Out[64]:

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

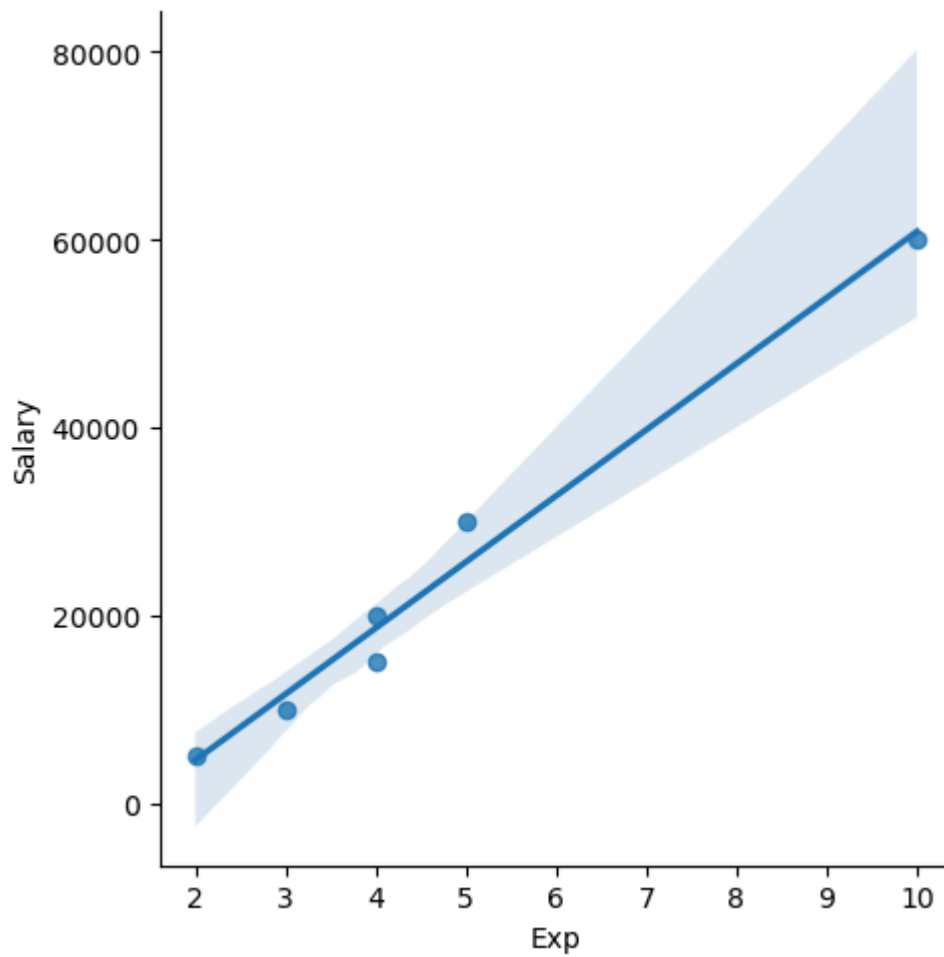
In [69]: `vis1=sns.distplot(clean_data['Salary'])`



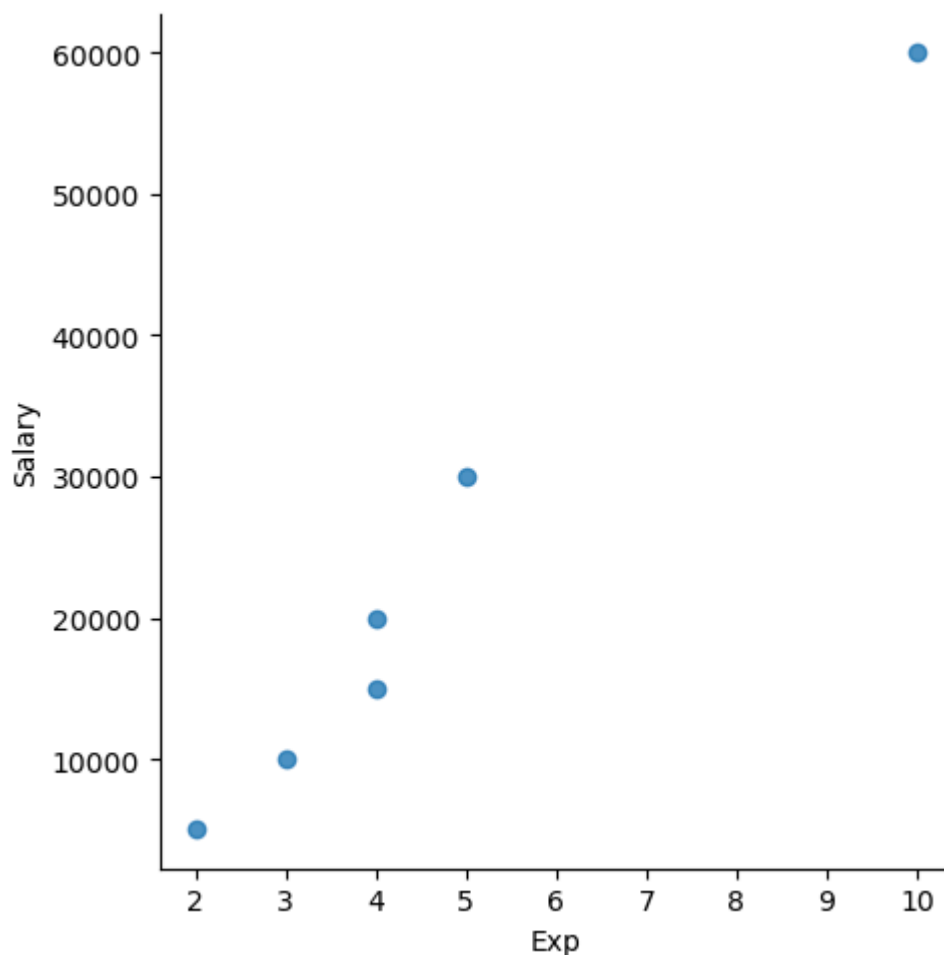
```
In [71]: vis2=plt.hist(clean_data['Salary'])
```



```
In [72]: vis4=sns.lmplot(data=clean_data,x='Exp',y='Salary')
```



```
In [74]: vis5=sns.lmplot(data=clean_data,x='Exp',y='Salary',fit_reg=False)
```



```
In [75]: clean_data[:]
```

```
Out[75]:
```

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

```
Out[76]:
```

	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 [78]: clean_data[:, :-1]
```

Out[78]:

	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 [79]: `clean_data.columns`

Out[79]: `Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')`

In [80]: `X_iv=clean_data[['Name','Domain','Age','Location','Exp']]`

In [81]: `X_iv`

Out[81]:

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

In [83]: `y_dv`

Out[83]:

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

In [84]: `emp`

Out[84]:

	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 [85]: `clean_data`

Out[85]:

	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 [86]: `X_iv`

Out[86]:

	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 [87]: `y_dv`

Out[87]: **Salary**

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

In [88]: `clean_data`

Out[88]:

	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

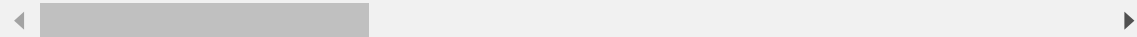
In [89]: `imputation=pd.get_dummies(clean_data)`

In [90]: `imputation`

Out[90]:

	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 [91]: `clean_data`

Out[91]:

	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 [92]: imputation

Out[92]:

	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

%% md

raw data with lot of regex, missing, unclean data

regex, clean

fill missing numerical & categorical

clean_dataset (data cleaning) 3 month - 5 month

outlier treatment, univariate, bivariate, correlation

split the data into x_i.v & y_dv

impute categorical data to numerical

eda part complete

%% md

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}

In []: