

Converting Raw data to Clean Data using Python EDA

EDA Techniques

1. Variable Identification
2. Univariate Analysis
3. BiVariate Analysis
4. Missing Value Treatment
5. Outlier Treatment
6. Imputational Technique / Transformers
7. Variable Creation

```
In [4]: import pandas as pd
```

```
In [5]: pd.__version__
```

```
Out[5]: '2.2.2'
```

```
In [6]: emp=pd.read_excel(r"C:\Users\ymani\Dropbox\PC\Downloads\Rawdata.xlsx")
```

```
In [7]: emp
```

```
Out[7]:
```

	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 [8]: id(emp)
```

```
Out[8]: 2282001247616
```

```
In [9]: emp.columns
```

```
Out[9]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')
```

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

```
Out[10]: (6, 6)
```

```
In [11]: #^#-%$ -regex
```

```
In [12]: emp.head()
```

```
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

```
In [13]: emp.tail()
```

```
Out[13]:
```

	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 [14]: 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 [15]: emp.isnull()
```

Out[15]:

	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 [16]: `emp.isna()`

Out[16]:

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

Out[17]:

Name	0
Domain	0
Age	2
Location	2
Salary	0
Exp	1

dtype: int64

Data Cleaning or Data Cleansing

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

Out[19]:

0	Mike
1	Teddy^
2	Uma#r
3	Jane
4	Uttam*
5	Kim

Name: Name, dtype: object

In [20]: `emp['Name']=emp['Name'].str.replace(r'\W','',regex=True) # non word character`

In [21]: `emp['Name']`

```
Out[21]: 0    Mike
          1    Teddy
          2    Umar
          3    Jane
          4    Uttam
          5    Kim
          Name: Name, dtype: object
```

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

```
Out[22]: 0    Datascience#$
          1    Testing
          2    Dataanalyst^^#
          3    Ana^alytics
          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']
```

```
Out[24]: 0    34 years
          1    45' yr
          2    NaN
          3    NaN
          4    67-yr
          5    55yr
          Name: Age, 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]: emp['Age']=emp['Age'].str.extract(r'(\d+)')
          emp['Age']
```

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

```
In [28]: emp
```

```
Out[28]:
```

	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 [29]: emp['Location']
```

```
Out[29]: 0      Mumbai
         1      Bangalore
         2      NaN
         3      Hyderbad
         4      NaN
         5      Delhi
         Name: Location, dtype: object
```

```
In [30]: emp['Salary']
```

```
Out[30]: 0      5^00#0
         1      10%%000
         2      1$5%000
         3      2000^0
         4      30000-
         5      6000^$0
         Name: Salary, dtype: object
```

```
In [31]: emp['Salary']=emp['Salary'].str.replace(r'\W','',regex=True)
```

```
In [32]: emp['Salary']
```

```
Out[32]: 0      5000
         1      10000
         2      15000
         3      20000
         4      30000
         5      60000
         Name: Salary, dtype: object
```

```
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(r'(\d+)')
```

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

```
Out[35]: 0      2
1      3
2      4
3      NaN
4      5
5     10
Name: Exp, dtype: object
```

```
In [36]: emp
```

```
Out[36]:
```

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

```
In [38]: clean_data
```

```
Out[38]:
```

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

```
Out[39]: Name      0
        Domain    0
        Age       2
        Location   2
        Salary     0
        Exp       1
        dtype: int64

Missing Treatment
```

```
In [41]: clean_data['Age']
```

```
Out[41]: 0      34
        1      45
        2     NaN
        3     NaN
        4      67
        5      55
        Name: Age, dtype: object
```

Missing Value Treatment

```
In [43]: import numpy as np
```

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

```
In [45]: clean_data['Age']
```

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

```
In [46]: clean_data['Exp']
```

```
Out[46]: 0      2
        1      3
        2      4
        3     NaN
        4      5
        5     10
        Name: Exp, dtype: object
```

```
In [47]: clean_data['Exp']=clean_data['Exp'].fillna(np.mean(pd.to_numeric(clean_data['Exp
```

```
In [48]: clean_data['Exp']
```

```
Out[48]: 0      2
          1      3
          2      4
          3      4.8
          4      5
          5     10
          Name: Exp, 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	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 [50]: clean_data['Location']
```

```
Out[50]: 0      Mumbai
          1    Bangalore
          2         NaN
          3    Hyderbad
          4         NaN
          5      Delhi
          Name: Location, dtype: object
```

```
In [51]: clean_data['Location']=clean_data['Location'].fillna(clean_data['Location'].mode[0])
          clean_data['Location']
```

```
Out[51]: 0      Mumbai
          1    Bangalore
          2    Bangalore
          3    Hyderbad
          4    Bangalore
          5      Delhi
          Name: Location, dtype: object
```

```
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	50.25	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 [53]: 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
```

Converting object to int

```
In [55]: clean_data['Age']=clean_data["Age"].astype(int)
```

```
In [56]: clean_data['Age']
```

```
Out[56]: 0    34
         1    45
         2    50
         3    50
         4    67
         5    55
         Name: Age, dtype: int32
```

```
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     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 [58]: clean_data['Salary']=clean_data['Salary'].astype(int)
         clean_data['Exp']=clean_data['Exp'].astype(int)
         print(clean_data['Salary'])
         clean_data['Exp']
```

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

```
Out[58]: 0      2
          1      3
          2      4
          3      4
          4      5
          5     10
          Name: Exp, dtype: int32
```

```
In [59]: 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
```

Converting Object to Category

```
In [61]: clean_data['Name']=clean_data['Name'].astype('category')
          clean_data['Domain']=clean_data['Domain'].astype('category')
          clean_data['Location']=clean_data['Location'].astype('category')
          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 [62]: clean_data
```

```
Out[62]:
```

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

```
In [64]: import os
os.getcwd() #from the os given the saved current working directory
```

```
Out[64]: 'C:\\Users\\ymani\\Full Stack Data Science'
```

Univariate Analysis

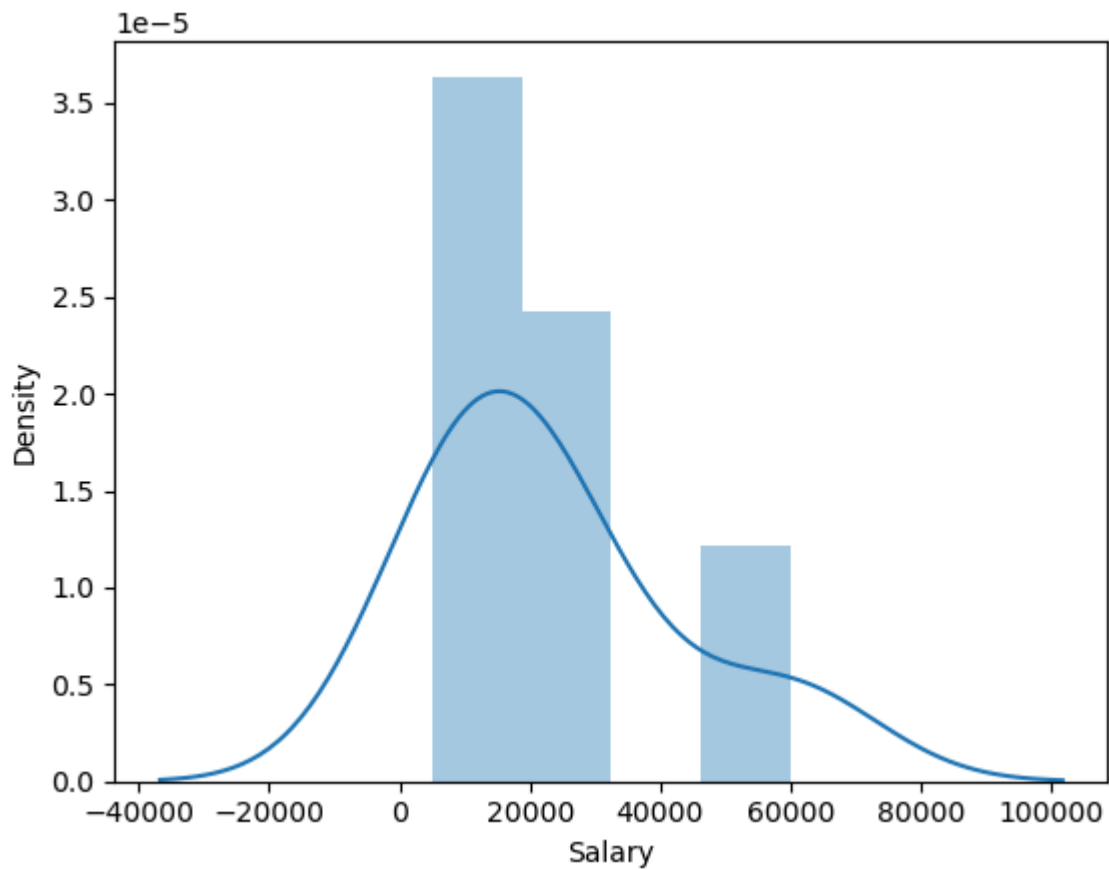
```
In [66]: import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [67]: import warnings
warnings.filterwarnings('ignore')
```

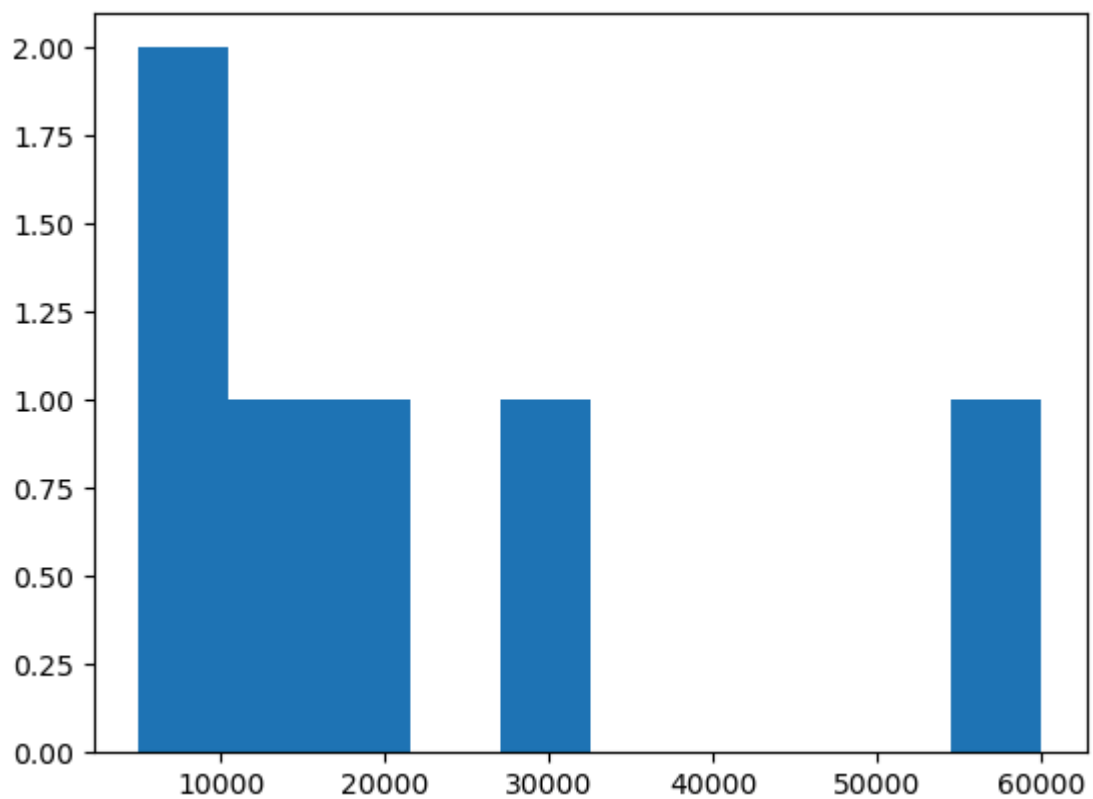
```
In [68]: clean_data['Salary']
```

```
Out[68]: 0    5000
1   10000
2   15000
3   20000
4   30000
5   60000
Name: Salary, dtype: int32
```

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



In [139... `vis2=plt.hist(clean_data['Salary'])`



BI- Variate ANalysis

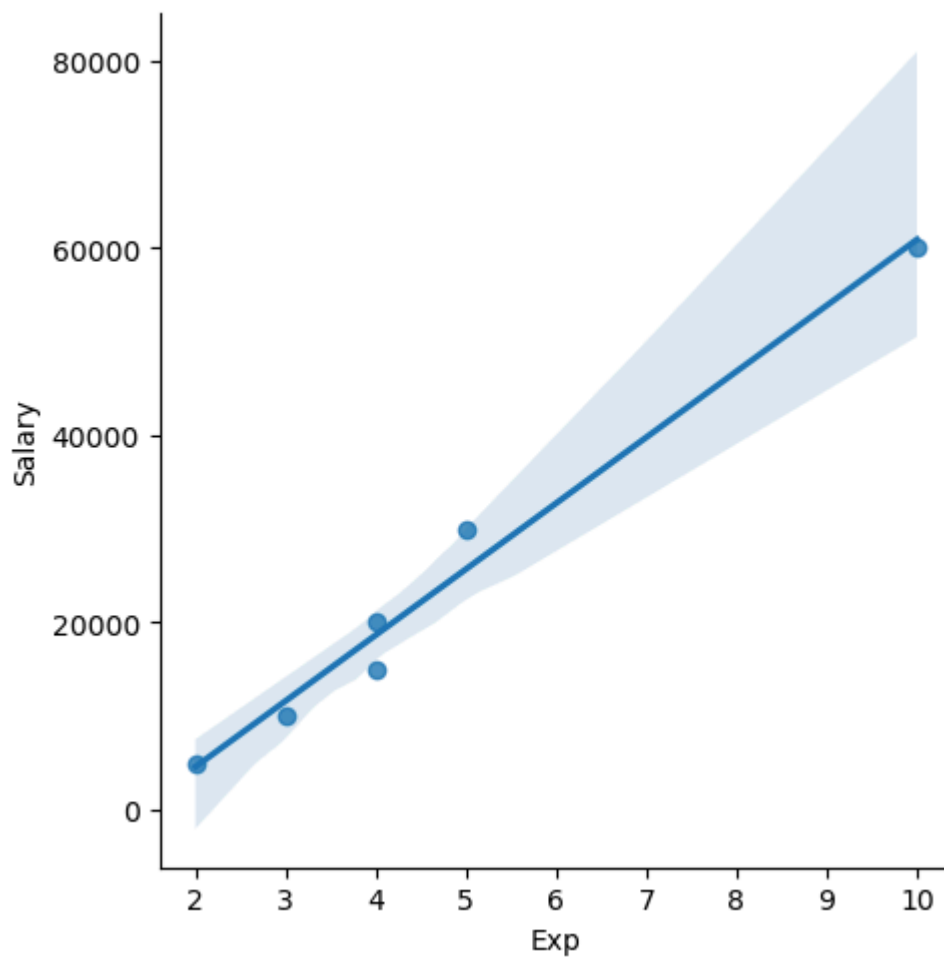
In [133... `clean_data`

Out[133...

	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

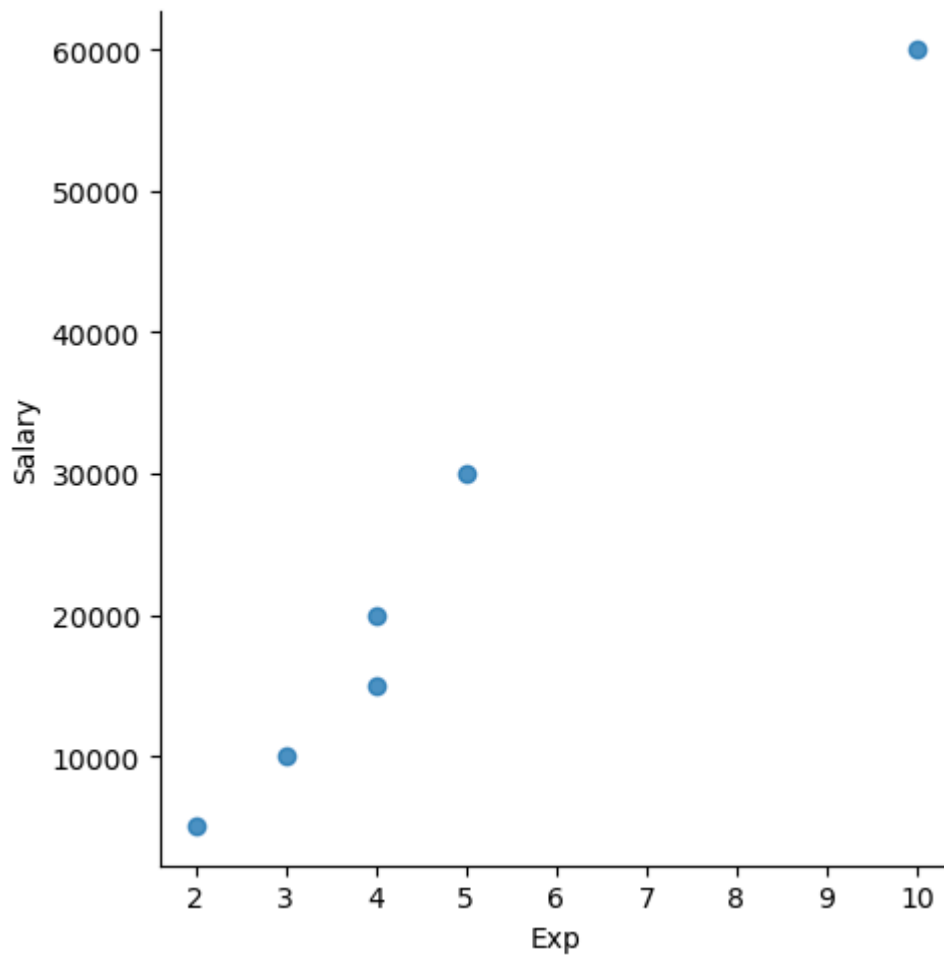
In [135...

```
vis3=sns.lmplot(data=clean_data,x='Exp',y='Salary')
```



In [143...

```
vis4=sns.lmplot(data=clean_data,x='Exp',y='Salary',fit_reg=False)
```



Slicing

In [146... `clean_data[:]`

Out[146...

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

Out[148...

	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 [150...

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

Out[150...

	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

Variable Identification

In [153...

```
clean_data
```

Out[153...

	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 the above Dataset, we need to find what is Dependent and Independent Variables.

Here Salary is Dependent Variable, and remaining all are Independent variable. x_iv (x_independentvariable), y_dv(y_dependent variable)

In [158...

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

Out[158...

	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 [169...

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

In [171...

```
y_dv
```

Out[171...

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

It is a multiple regression because there are 5 Independent Variables.

Impute Categorical data to numerical data

-

In [173...

```
emp
```

Out[173...

	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 [177...

```
x_iv
```


Out[177...

	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 [179...

```
y_dv
```

Out[179...

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

In [175...

```
clean_data
```

Out[175...

	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 [181...

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

In [189...

```
imputation=imputation.astype(int)
```

In [191...

```
imputation
```

Out[191...

	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