

```
In [45]: # EDA RAW DATA TO CLEAN DATA
```

```
import pandas as pd
```

```
In [46]: data_raw=pd.read_excel(r'Users/sasidharbhagavatula/Desktop/Rawdata.xlsx')
```

```
In [47]: data_raw
```

```
Out[47]:
```

	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 [48]: pd.__version__
```

```
Out[48]: '2.2.3'
```

```
In [49]: data_raw.isnull()
```

```
Out[49]:
```

	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 [50]: id(data_raw) # memory address location
```

```
Out[50]: 5103326736
```

```
In [51]: data_raw.columns
```

```
Out[51]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')
```

```
In [52]: data_raw.shape
```

```
Out[52]: (6, 6)
```

```
In [53]: data_raw.head()
```

```
Out[53]:
```

	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 [54]: data_raw.tail()
```

```
Out[54]:
```

	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 [55]: data_raw
```

```
Out[55]:
```

	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 [56]: data_raw.isnull()
```

```
Out[56]:    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 [57]: data_raw.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 [61]: data_raw.isna() # same as null
```

```
Out[61]:    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 [62]: data_raw.isnull().sum()
```

```
Out[62]: Name      0
Domain    0
Age       2
Location  2
Salary    0
Exp       1
dtype: int64
```

DATA CLEANING

In [63]: data_raw

	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 [64]: *# cleaning employee name ----- removing characters*

In [65]: data_raw['Name'] # parsing one variable

```
Out[65]: 0      Mike
         1      Teddy^
         2      Uma#r
         3      Jane
         4      Uttam*
         5      Kim
Name: Name, dtype: object
```

In [66]: data_raw['Name']= data_raw['Name'].str.replace(r'\W', ' ', regex=True)

In [67]: data_raw['Name']

```
Out[67]: 0      Mike
         1      Teddy
         2      Umar
         3      Jane
         4      Uttam
         5      Kim
Name: Name, dtype: object
```

In [68]: data_raw['Domain']

```
Out[68]: 0      Datascience#$
         1      Testing
         2      Dataanalyst^^#
         3      Ana^^lytics
         4      Statistics
         5      NLP
Name: Domain, dtype: object
```

In [69]: data_raw['Domain']= data_raw['Domain'].str.replace(r'\W', ' ', regex=True)

```
In [70]: data_raw['Domain']
```

```
Out[70]: 0      DataScience
         1      Testing
         2      Dataanalyst
         3      Analytics
         4      Statistics
         5      NLP
Name: Domain, dtype: object
```

```
In [71]: data_raw['Location']=data_raw['Location'].str.replace(r'\W',' ', regex=True)
```

```
In [72]: data_raw['Location']
```

```
Out[72]: 0      Mumbai
         1      Bangalore
         2      NaN
         3      Hyderabad
         4      NaN
         5      Delhi
Name: Location, dtype: object
```

```
In [73]: data_raw['Age']
```

```
Out[73]: 0      34 years
         1      45' yr
         2      NaN
         3      NaN
         4      67-yr
         5      55yr
Name: Age, dtype: object
```

```
In [74]: data_raw['Age']=data_raw['Age'].str.extract('(\d+)')
```

```
In [75]: data_raw['Age']
```

```
Out[75]: 0      34
         1      45
         2      NaN
         3      NaN
         4      67
         5      55
Name: Age, dtype: object
```

```
In [76]: data_raw['Salary']
```

```
Out[76]: 0      5^00#0
          1      10%000
          2      1$5%000
          3      2000^0
          4      30000-
          5      6000^$0
Name: Salary, dtype: object
```

```
In [77]: data_raw['Salary']=data_raw['Salary'].str.replace(r'\W', '', regex=True)
```

```
In [78]: data_raw['Salary']
```

```
Out[78]: 0      5000
          1      10000
          2      15000
          3      20000
          4      30000
          5      60000
Name: Salary, dtype: object
```

```
In [79]: data_raw['Exp']=data_raw['Exp'].str.extract('(\d+)')
```

```
In [80]: data_raw['Exp']
```

```
Out[80]: 0      2
          1      3
          2      4
          3      NaN
          4      5
          5      10
Name: Exp, dtype: object
```

```
In [81]: data_raw
```

```
Out[81]:   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 [82]: clean_data=data_raw.copy()
```

```
In [83]: clean_data
```

Out[83]:

	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 []:

In []:

In []:

In []: