

# EDA

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

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

```
Out[2]: '2.2.3'
```

```
In [3]: emp = pd.read_excel(r"Z:\DS\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	Hyderbad	2000^0	NaN
4	Uttam*	Statistics	67-yr	NaN	30000-	5+ year
5	Kim	NLP	55yr	Delhi	6000^\$0	10+

# data cleaning

```
In [6]: emp['Name']
```

```
Out[6]: 0      Mike
        1    Teddy^
        2     Uma#r
        3      Jane
        4   Uttam*
        5      Kim
Name: Name, dtype: object
```

```
In [8]: emp['Domain']
```

```
Out[8]: 0      Datascience#$#
        1      Testing
        2  Dataanalyst^^#
        3   Ana^^lytics
        4      Statistics
        5       NLP
Name: Domain, dtype: object
```

```
In [9]: emp['Name'] = emp['Name'].str.replace(r'\W', ' ', regex=True)
```

```
In [10]: emp['Name']
```

```
Out[10]: 0      Mike
        1    Teddy
        2     Umar
        3      Jane
        4   Uttam
        5      Kim
Name: Name, dtype: object
```

```
In [14]: emp.columns
```

```
Out[14]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')
```

```
In [15]: emp.head(1)
```

```
Out[15]:    Name      Domain    Age Location    Salary    Exp
          0  Mike  Datascience#$  34 years  Mumbai  5^00#0    2+
```

In [16]: `emp['Exp']`

Out[16]:

0	2+
1	<3
2	4> yrs
3	NaN
4	5+ year
5	10+

Name: Exp, dtype: object

In [17]: `emp['Domain']`

Out[17]:

0	Datascience#\$
1	Testing
2	Dataanalyst^^#
3	Ana^^lytics
4	Statistics
5	NLP

Name: Domain, dtype: object

In [18]: `emp['Domain'] = emp['Domain'].str.replace(r'\W', ' ', regex=True)`

In [19]: `emp`

Out[19]:

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

In [22]: `emp['Location'] = emp['Location'].str.replace(r'\W', ' ', regex=True)`

```
In [ ]: emp
```

```
In [55]: emp['Age'] = emp['Age'].str.replace(r'\W', ' ', regex=True)
```

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

```
Out[56]: 0    34years  
1      45yr  
2      NaN  
3      NaN  
4      67yr  
5      55yr  
Name: Age, dtype: object
```

```
In [57]: emp['Age'] = emp['Age'].str.extract('(\d+)') # r(r'(\d+)')
```

```
In [58]: emp['Age']
```

```
Out[58]: 0    34  
1    45  
2    NaN  
3    NaN  
4    67  
5    55  
Name: Age, dtype: object
```

```
In [59]: emp
```

Out[59]:

	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%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 [60]: `emp['Salary'] = emp['Salary'].str.replace(r'\W', '', regex=True)`In [61]: `emp['Salary']`

Out[61]:

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

Name: Salary, dtype: object

In [62]: `emp`

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	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 [63]:

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

In [64]:

```
emp
```

Out[64]:

	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 [65]:

```
clean_data = emp.copy()
```

In [66]:

```
emp
```

Out[66]:

	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 [ ]: