

21-11-2025 raw data to clean data conversion using python EDA

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

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

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

```
In [3]: emp = pd.read_excel(r"C:\Users\karthik reddy\Downloads\Rawdata.xlsx")
```

```
In [4]: emp
```

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascienc#\$	34 years	Mumbai	5^00#0	2+
1	Teddy^	Testing	45' yr	Bangalore	10%0000	<3
2	Uma#r	Dataanalyst^^#	NaN	NaN	1\$5%0000	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 [5]: id(emp)
```

```
Out[5]: 1626167655936
```

```
In [6]: emp.columns
```

```
Out[6]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')
```

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

```
Out[7]: (6, 6)
```

```
In [8]: emp.head()
```

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascienc#\$	34 years	Mumbai	5^00#0	2+
1	Teddy^	Testing	45' yr	Bangalore	10%0000	<3
2	Uma#r	Dataanalyst^^#	NaN	NaN	1\$5%0000	4> yrs
3	Jane	Ana^^lytics	NaN	Hyderbad	2000^0	NaN
4	Uttam*	Statistics	67-yr	NaN	30000-	5+ year

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

Out[9]:

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

In [10]:

`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 [11]:

`emp`

Out[11]:

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascienc#\$	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	Analytics	NaN	Hyderabad	2000^0	NaN
4	Uttam*	Statistics	67-yr	NaN	30000-	5+ year
5	Kim	NLP	55yr	Delhi	6000^\$0	10+

In [12]:

`emp.isnull()`

Out[12]:

	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 [13]: `emp.isna()`

	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	

Data Cleaning

In [15]: `emp`

	Name	Domain	Age	Location	Salary	Exp
0	Mike	Datascienc#	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 [16]: `emp['Name']`

Out[16]:

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

Name: Name, dtype: object

In [17]: `emp['Name'] = emp['Name'].str.replace(r'\W', ' ', regex=True)`

In [18]: `emp['Name']`

```
Out[18]: 0      Mike
          1      Teddy
          2      Umar
          3      Jane
          4      Uttam
          5      Kim
Name: Name, dtype: object
```

```
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%0000 <3
2   Umar  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 [20]: emp.columns
```

```
Out[20]: Index(['Name', 'Domain', 'Age', 'Location', 'Salary', 'Exp'], dtype='object')
```

```
In [21]: emp.head(1)
```

```
Out[21]:   Name      Domain    Age Location  Salary  Exp
0   Mike  Datascience#$  34 years  Mumbai  5^00#0   2+
```

```
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)
```

```
In [24]: emp['Domain']
```

```
Out[24]: 0      Datascience
          1      Testing
          2      Dataanalyst
          3      Analytics
          4      Statistics
          5      NLP
Name: Domain, dtype: object
```

```
In [25]: emp
```

Out[25]:

	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 [26]: `emp['Location'] = emp['Location'].str.replace(r'\W', '', regex=True)`In [27]: `emp['Location']`

Out[27]:

0	Mumbai
1	Bangalore
2	NaN
3	Hyderbad
4	NaN
5	Delhi

Name: Location, dtype: object

In [28]: `emp`

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

Out[30]:

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

Name: Age, dtype: object

In [31]: `emp['Age'] = emp['Age'].str.extract('(\d+)') # r(r'(\d+)')`In [32]: `emp['Age']`

```
Out[32]: 0      34
         1      45
         2    NaN
         3    NaN
         4      67
         5      55
Name: Age, dtype: object
```

```
In [33]: emp
```

```
Out[33]:   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 [34]: emp['Salary'] = emp['Salary'].str.replace(r'\W', '', regex=True)
```

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

```
Out[35]: 0      5000
         1     10000
         2     15000
         3     20000
         4     30000
         5     60000
Name: Salary, 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> yrs
3   Jane    Analytics  NaN  Hyderbad  20000    NaN
4  Uttam    Statistics  67      NaN  30000  5+ year
5    Kim        NLP  55  Delhi  60000  10+
```

```
In [37]: emp['Exp'] = emp['Exp'].str.extract('(\d+)')
```

```
In [38]: emp['Exp']
```

```
Out[38]: 0      2
         1      3
         2      4
         3    NaN
         4      5
         5     10
Name: Exp, dtype: object
```

```
In [39]: emp
```

```
Out[39]:   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 [40]: clean_data = emp.copy()
```

```
In [41]: clean_data
```

```
Out[41]:   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 [ ]:
```