

16. Data Type Transformation(Replace and Data Type Change)

- 3+ is categorical (object) data while other rows in this column contain numerical data
- We will remove 3+ and convert its data type from object data type to int dtype

```
In [2]: import pandas as pd
```

```
In [4]: dataset = pd.read_csv('loan.csv')
dataset.head(3)
```

```
Out[4]:
```

	Loan_ID	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome
0	LP001002	Male	No	0	Graduate	No	5849	0
1	LP001003	Male	Yes	1	Graduate	No	4583	0
2	LP001005	Male	Yes	0	Graduate	Yes	3000	0

- 3+ is present in 'Dependents' column

```
In [6]: dataset.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 614 entries, 0 to 613
Data columns (total 13 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Loan_ID               614 non-null   object
1   Gender                601 non-null   object
2   Married               611 non-null   object
3   Dependents            599 non-null   object
4   Education             614 non-null   object
5   Self_Employed         582 non-null   object
6   ApplicantIncome       614 non-null   int64
7   CoapplicantIncome     614 non-null   float64
8   LoanAmount            592 non-null   float64
9   Loan_Amount_Term      600 non-null   float64
10  Credit_History         564 non-null   float64
11  Property_Area         614 non-null   object
12  Loan_Status           614 non-null   object
dtypes: float64(4), int64(1), object(8)
memory usage: 62.5+ KB
```

- 'Dependent' column has null value, b/c RangeIndex:614, whereas (Dependents 599 non-null object) - as shown above

```
In [7]: dataset.isnull().sum()
```

```
Out[7]: Loan_ID          0
        Gender          13
        Married         3
        Dependents      15
        Education        0
        Self_Employed    32
        ApplicantIncome   0
        CoapplicantIncome 0
        LoanAmount       22
        Loan_Amount_Term  14
        Credit_History    50
        Property_Area     0
        Loan_Status       0
        dtype: int64
```

```
In [9]: dataset['Dependents'].value_counts()
```

```
Out[9]: 0      345
        1      102
        2      101
        3+       51
        Name: Dependents, dtype: int64
```

```
In [12]: # Remove null values in Dependents column
dataset['Dependents'].fillna(dataset['Dependents'].mode()[0], inplace=True)
```

```
In [13]: dataset.isnull().sum()
```

```
Out[13]: Loan_ID          0
        Gender          13
        Married         3
        Dependents       0
        Education        0
        Self_Employed    32
        ApplicantIncome   0
        CoapplicantIncome 0
        LoanAmount       22
        Loan_Amount_Term  14
        Credit_History    50
        Property_Area     0
        Loan_Status       0
        dtype: int64
```

Replace 3+ with 3

```
In [15]: dataset['Dependents'].replace('3+', '3', inplace=True)
```

```
In [16]: dataset['Dependents'].value_counts()
```

```
Out[16]: 0    360
         1    102
         2    101
         3     51
         Name: Dependents, dtype: int64
```

```
In [17]: dataset.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 614 entries, 0 to 613
Data columns (total 13 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Loan_ID               614 non-null   object
1   Gender                601 non-null   object
2   Married               611 non-null   object
3   Dependents            614 non-null   object
4   Education             614 non-null   object
5   Self_Employed         582 non-null   object
6   ApplicantIncome       614 non-null   int64
7   CoapplicantIncome     614 non-null   float64
8   LoanAmount            592 non-null   float64
9   Loan_Amount_Term      600 non-null   float64
10  Credit_History         564 non-null   float64
11  Property_Area          614 non-null   object
12  Loan_Status           614 non-null   object
dtypes: float64(4), int64(1), object(8)
memory usage: 62.5+ KB
```

Now change the datatype

```
In [19]: dataset['Dependents'] = dataset['Dependents'].astype("int64")
```

```
In [20]: dataset.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 614 entries, 0 to 613
Data columns (total 13 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Loan_ID               614 non-null   object
1   Gender                601 non-null   object
2   Married               611 non-null   object
3   Dependents            614 non-null   int64
4   Education             614 non-null   object
5   Self_Employed         582 non-null   object
6   ApplicantIncome       614 non-null   int64
7   CoapplicantIncome     614 non-null   float64
8   LoanAmount            592 non-null   float64
9   Loan_Amount_Term      600 non-null   float64
10  Credit_History         564 non-null   float64
11  Property_Area          614 non-null   object
12  Loan_Status           614 non-null   object
dtypes: float64(4), int64(2), object(7)
memory usage: 62.5+ KB
```

```
In [22]: dataset['Dependents'].value_counts()
```

```
Out[22]: 0    360  
        1    102  
        2    101  
        3     51  
        Name: Dependents, dtype: int64
```

```
In [ ]:
```