# LABSHEET-6_PART-B

# COURSE TITLE : DATA AND VISUAL ANALYTICS LAB

# NAME: RETHINAGIRI G ¶

# ROLLNO : 225229130

```python
In [1]:  1  import pandas as pd
         2  from sklearn.preprocessing import LabelEncoder
```

```python
In [2]:  1  le = LabelEncoder()
         2  df = pd.DataFrame(data = {'col1': ['foo','bar','foo','bar'], 'col2': ['x', 'y', 'x', 'z'], 'col3':[1,2,3,4]})
```

```python
In [3]:  1  df.apply(le.fit_transform)
```

Out[3]:

|   | col1 | col2 | col3 |
|---|------|------|------|
| **0** | 1 | 0 | 0 |
| **1** | 0 | 1 | 1 |
| **2** | 1 | 0 | 2 |
| **3** | 0 | 2 | 3 |

# One Hot Encoder

```
In [4]:  1  import pandas as pd
         2  df = pd.DataFrame({'A': ['a','b','a'], 'B': ['b','a','c'], 'C': [1, 2, 3]})
         3  df
         4
```

Out[4]:

|   | A | B | C |
|---|---|---|---|
| 0 | a | b | 1 |
| 1 | b | a | 2 |
| 2 | a | c | 3 |

```
In [5]:  1  pd.get_dummies(df, prefix=['col1','col2'])
```

Out[5]:

|   | C | col1_a | col1_b | col2_a | col2_b | col2_c |
|---|---|--------|--------|--------|--------|--------|
| 0 | 1 | 1      | 0      | 0      | 1      | 0      |
| 1 | 2 | 0      | 1      | 1      | 0      | 0      |
| 2 | 3 | 1      | 0      | 0      | 0      | 1      |

## MinMaxScaler

```
In [6]:  1  from sklearn.preprocessing import MinMaxScaler
         2  mm_scaler = MinMaxScaler(feature_range = (0,1)) # (0,1 ) is default range
         3  df2 = pd.DataFrame({'col1':[5,-41, -67],
         4   'col2': [23, - 53, -36],
         5   'col3': [-25,10, 17]})
         6  mm_scaler.fit_transform(df2)
```

Out[6]: array([[1.        , 1.        , 0.        ],
              [0.36111111, 0.        , 0.83333333],
              [0.        , 0.22368421, 1.        ]])

## Binarizer

```python
In [7]:  1  from sklearn.preprocessing import Binarizer
         2  dfb = pd.DataFrame({'col1': [110, 200],
         3   'col2': [120, 800],
         4   'col3': [310, 400]})
         5  bin = Binarizer(threshold=300)
         6  bin.fit_transform(dfb)
         7
```

Out[7]: array([[0, 0, 1],
               [0, 1, 1]], dtype=int64)

## Imputer

```python
In [8]:   1  import numpy as np
          2  from sklearn.impute import SimpleImputer
          3  import pandas as pd
          4  imp_mean = SimpleImputer(missing_values=np.nan, strategy='mean')
          5  df = pd.DataFrame({'col1': [7, 2, 3],
          6   'col2': [4, np.nan, 6],
          7   'col3': [np.nan, np.nan, 3],
          8   'col4': [10, np.nan, 9]})
          9  print(df)
         10  imp_mean.fit_transform(df)
```

```
   col1  col2  col3  col4
0     7   4.0   NaN  10.0
1     2   NaN   NaN   NaN
2     3   6.0   3.0   9.0
```

Out[8]: array([[ 7. ,  4. ,  3. , 10. ],
               [ 2. ,  5. ,  3. ,  9.5],
               [ 3. ,  6. ,  3. ,  9. ]])

# De-duplication or Entity Resolution and String Matching

```
In [10]:    1  !pip install fuzzywuzzy
```

```
Collecting fuzzywuzzy
  Downloading fuzzywuzzy-0.18.0-py2.py3-none-any.whl (18 kB)
Installing collected packages: fuzzywuzzy
Successfully installed fuzzywuzzy-0.18.0
```

```
In [11]:    1  import warnings
            2  warnings.filterwarnings('ignore')
            3  from fuzzywuzzy import fuzz
            4  from fuzzywuzzy import process
            5  a = 'Welcome to Bishop Heber College'
            6  b = 'I am Sam pursuing Masters in DataScience at Bishop Heber College'
            7  ratio = fuzz.ratio(a, b)
            8  weighted_ratio = fuzz.WRatio(a, b)
            9  unicode_ratio = fuzz.UQRatio(a, b)
           10  print('Ratio =', ratio)
           11  print('Weighted ratio =', weighted_ratio)
           12  print('Unicode ratio =', unicode_ratio)
           13
```

```
Ratio = 55
Weighted ratio = 86
Unicode ratio = 55
```

```
In [13]:    1  c = a + b
            2  c
```

Out[13]:  'Welcome to Bishop Heber CollegeI am Sam pursuing Masters in DataScience at Bishop Heber College'

```
In [14]:    1  ex_tract = process.extract('I', c)
            2  ex_tract
```

Out[14]:  [('i', 100), ('I', 100), ('i', 100), ('i', 100), ('i', 100)]

```
In [15]:    1  ex_tract_1 = process.extractOne('I', c)
            2  ex_tract_1
```

Out[15]:  ('i', 100)

```
In [ ]: 1
```