

ML_Data_PreProcessing

May 13, 2019

1 Data Pre-Processing

```
In [1]: import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
```

2 import datasets

```
In [2]: ### import the dataset by placing the csvfile in the present working dir.
```

```
datasets = pd.read_csv('Data.csv') #pd pandas
```

```
# independent variable
```

```
X = datasets.iloc[:, :-1].values #: all lines, coulumn minus 1 bcz of index value
print (X)
```

```
# dependent variable
```

```
y = datasets.iloc[:, 3].values # last column ie 2
print(y)
```

```
# nan means no value in the field
```

```
[['France' 44.0 72000.0]
 ['Spain' 27.0 48000.0]
 ['Germany' 30.0 54000.0]
 ['Spain' 38.0 61000.0]
 ['Germany' 40.0 nan]
 ['France' 35.0 58000.0]
 ['Spain' nan 52000.0]
 ['France' 48.0 79000.0]
 ['Germany' 50.0 83000.0]
 ['France' 37.0 67000.0]]
['No' 'Yes' 'No' 'No' 'Yes' 'Yes' 'No' 'Yes' 'No' 'Yes']
```

3 Missing data

```
In [3]: # importing a Imputer class from sklearn library
        from sklearn.preprocessing import Imputer

        # object imputer: strategy- how to replace the value, here its mean value, axis either
        imputer = Imputer(missing_values = np.nan, strategy='mean', axis=0)

        # X can take all values[:, column value index 1 and 2; here it excludes the upper bound
        imputer = imputer.fit(X[:, 1:3])

        # imputer.transform means replacing the values to that location
        X[:, 1:3] = imputer.transform(X[:, 1:3])

        print(X)

[['France' 44.0 72000.0]
 ['Spain' 27.0 48000.0]
 ['Germany' 30.0 54000.0]
 ['Spain' 38.0 61000.0]
 ['Germany' 40.0 63777.777777777778]
 ['France' 35.0 58000.0]
 ['Spain' 38.77777777777778 52000.0]
 ['France' 48.0 79000.0]
 ['Germany' 50.0 83000.0]
 ['France' 37.0 67000.0]]

d:\ProgramData\Anaconda3\lib\site-packages\sklearn\utils\deprecation.py:58: DeprecationWarning
  warnings.warn(msg, category=DeprecationWarning)
```

4 categorical data

```
In [17]: #from sklearn import preprocessing
        #le= preprocessing.LabelEncoder()
        #X[:,0]= le.fit_transform(X[:,0])
        #print(X)    ## ist column char data are trasformed to numeric values.
```

5 Inspect a method in jupyter itself

```
import inspect from sklearn import tree inspect.getsourcelines(tree.tree)
```

```
In [10]: print("hai")
```

```
hai
```

```
# OneHotEncoder-categorize the character value of column 1
```

```
In [4]: from sklearn.preprocessing import LabelEncoder,OneHotEncoder
        labelencoder = LabelEncoder()
        X[:,0] = labelencoder.fit_transform(X[:, 0])
        onehotencoder = OneHotEncoder(categorical_features=[0])
        X = onehotencoder.fit_transform(X).toarray()
        print(X)
```

```
# encode the dependent variable ie, purchased column
```

```
labelencoder_y = LabelEncoder()
y = labelencoder.fit_transform(y)
print("\n",y)
```

```
[[1.00000000e+00 0.00000000e+00 0.00000000e+00 4.40000000e+01
 7.20000000e+04]
 [0.00000000e+00 0.00000000e+00 1.00000000e+00 2.70000000e+01
 4.80000000e+04]
 [0.00000000e+00 1.00000000e+00 0.00000000e+00 3.00000000e+01
 5.40000000e+04]
 [0.00000000e+00 0.00000000e+00 1.00000000e+00 3.80000000e+01
 6.10000000e+04]
 [0.00000000e+00 1.00000000e+00 0.00000000e+00 4.00000000e+01
 6.37777778e+04]
 [1.00000000e+00 0.00000000e+00 0.00000000e+00 3.50000000e+01
 5.80000000e+04]
 [0.00000000e+00 0.00000000e+00 1.00000000e+00 3.87777778e+01
 5.20000000e+04]
 [1.00000000e+00 0.00000000e+00 0.00000000e+00 4.80000000e+01
 7.90000000e+04]
 [0.00000000e+00 1.00000000e+00 0.00000000e+00 5.00000000e+01
 8.30000000e+04]
 [1.00000000e+00 0.00000000e+00 0.00000000e+00 3.70000000e+01
 6.70000000e+04]]
```

```
[0 1 0 0 1 1 0 1 0 1]
```

```
d:\ProgramData\Anaconda3\lib\site-packages\sklearn\preprocessing\_encoders.py:368: FutureWarning:
If you want the future behaviour and silence this warning, you can specify "categories='auto'"
In case you used a LabelEncoder before this OneHotEncoder to convert the categories to integers:
warnings.warn(msg, FutureWarning)
d:\ProgramData\Anaconda3\lib\site-packages\sklearn\preprocessing\_encoders.py:390: DeprecationWarning:
"use the ColumnTransformer instead.", DeprecationWarning)
```

```
# Splitting the dataset into test set and train set
```

```
In [5]: from sklearn.model_selection import train_test_split
        X_train, X_test, y_train, y_test = train_test_split( X , y , test_size = 0.2, random_s
        print(X_train)
        print("\n", X_test)
        print("\n", y_train)
        print("\n",y_test)
```

```
[[0.00000000e+00 1.00000000e+00 0.00000000e+00 4.00000000e+01
 6.37777778e+04]
[1.00000000e+00 0.00000000e+00 0.00000000e+00 3.70000000e+01
 6.70000000e+04]
[0.00000000e+00 0.00000000e+00 1.00000000e+00 2.70000000e+01
 4.80000000e+04]
[0.00000000e+00 0.00000000e+00 1.00000000e+00 3.87777778e+01
 5.20000000e+04]
[1.00000000e+00 0.00000000e+00 0.00000000e+00 4.80000000e+01
 7.90000000e+04]
[0.00000000e+00 0.00000000e+00 1.00000000e+00 3.80000000e+01
 6.10000000e+04]
[1.00000000e+00 0.00000000e+00 0.00000000e+00 4.40000000e+01
 7.20000000e+04]
[1.00000000e+00 0.00000000e+00 0.00000000e+00 3.50000000e+01
 5.80000000e+04]]
```

```
[[0.0e+00 1.0e+00 0.0e+00 3.0e+01 5.4e+04]
[0.0e+00 1.0e+00 0.0e+00 5.0e+01 8.3e+04]]
```

```
[1 1 1 0 1 0 0 1]
```

```
[0 0]
```

Feature Scaling - StandardScaler ## fit the values between 1 and -1 range irrespective of values eg age, salary

```
In [6]: from sklearn.preprocessing import StandardScaler
        sc_x = StandardScaler()
        X_train = sc_x.fit_transform (X_train)
        X_test = sc_x.transform(X_test)
        print(X_train)

        print("\n",X_test)
```

```
[[-1.          2.64575131 -0.77459667  0.26306757  0.12381479]
 [ 1.          -0.37796447 -0.77459667 -0.25350148  0.46175632]
 [-1.          -0.37796447  1.29099445 -1.97539832 -1.53093341]
 [-1.          -0.37796447  1.29099445  0.05261351 -1.11141978]
 [ 1.          -0.37796447 -0.77459667  1.64058505  1.7202972 ]
 [-1.          -0.37796447  1.29099445 -0.0813118  -0.16751412]]
```

```
[ 1.          -0.37796447 -0.77459667  0.95182631  0.98614835]
[ 1.          -0.37796447 -0.77459667 -0.59788085 -0.48214934]]

[[-1.          2.64575131 -0.77459667 -1.45882927 -0.90166297]
[-1.          2.64575131 -0.77459667  1.98496442  2.13981082]]
```

6 Datasets available in scikit learn library

```
In [8]: from sklearn.datasets import load_iris
        data = load_iris()
        print(data)
        data.target[[10, 25, 50]]
        list(data.target_names)
```

```
{'data': array([[5.1, 3.5, 1.4, 0.2],
 [4.9, 3. , 1.4, 0.2],
 [4.7, 3.2, 1.3, 0.2],
 [4.6, 3.1, 1.5, 0.2],
 [5. , 3.6, 1.4, 0.2],
 [5.4, 3.9, 1.7, 0.4],
 [4.6, 3.4, 1.4, 0.3],
 [5. , 3.4, 1.5, 0.2],
 [4.4, 2.9, 1.4, 0.2],
 [4.9, 3.1, 1.5, 0.1],
 [5.4, 3.7, 1.5, 0.2],
 [4.8, 3.4, 1.6, 0.2],
 [4.8, 3. , 1.4, 0.1],
 [4.3, 3. , 1.1, 0.1],
 [5.8, 4. , 1.2, 0.2],
 [5.7, 4.4, 1.5, 0.4],
 [5.4, 3.9, 1.3, 0.4],
 [5.1, 3.5, 1.4, 0.3],
 [5.7, 3.8, 1.7, 0.3],
 [5.1, 3.8, 1.5, 0.3],
 [5.4, 3.4, 1.7, 0.2],
 [5.1, 3.7, 1.5, 0.4],
 [4.6, 3.6, 1. , 0.2],
 [5.1, 3.3, 1.7, 0.5],
 [4.8, 3.4, 1.9, 0.2],
 [5. , 3. , 1.6, 0.2],
 [5. , 3.4, 1.6, 0.4],
 [5.2, 3.5, 1.5, 0.2],
 [5.2, 3.4, 1.4, 0.2],
 [4.7, 3.2, 1.6, 0.2],
 [4.8, 3.1, 1.6, 0.2],
 [5.4, 3.4, 1.5, 0.4],
```

[5.2, 4.1, 1.5, 0.1],
 [5.5, 4.2, 1.4, 0.2],
 [4.9, 3.1, 1.5, 0.2],
 [5. , 3.2, 1.2, 0.2],
 [5.5, 3.5, 1.3, 0.2],
 [4.9, 3.6, 1.4, 0.1],
 [4.4, 3. , 1.3, 0.2],
 [5.1, 3.4, 1.5, 0.2],
 [5. , 3.5, 1.3, 0.3],
 [4.5, 2.3, 1.3, 0.3],
 [4.4, 3.2, 1.3, 0.2],
 [5. , 3.5, 1.6, 0.6],
 [5.1, 3.8, 1.9, 0.4],
 [4.8, 3. , 1.4, 0.3],
 [5.1, 3.8, 1.6, 0.2],
 [4.6, 3.2, 1.4, 0.2],
 [5.3, 3.7, 1.5, 0.2],
 [5. , 3.3, 1.4, 0.2],
 [7. , 3.2, 4.7, 1.4],
 [6.4, 3.2, 4.5, 1.5],
 [6.9, 3.1, 4.9, 1.5],
 [5.5, 2.3, 4. , 1.3],
 [6.5, 2.8, 4.6, 1.5],
 [5.7, 2.8, 4.5, 1.3],
 [6.3, 3.3, 4.7, 1.6],
 [4.9, 2.4, 3.3, 1.],
 [6.6, 2.9, 4.6, 1.3],
 [5.2, 2.7, 3.9, 1.4],
 [5. , 2. , 3.5, 1.],
 [5.9, 3. , 4.2, 1.5],
 [6. , 2.2, 4. , 1.],
 [6.1, 2.9, 4.7, 1.4],
 [5.6, 2.9, 3.6, 1.3],
 [6.7, 3.1, 4.4, 1.4],
 [5.6, 3. , 4.5, 1.5],
 [5.8, 2.7, 4.1, 1.],
 [6.2, 2.2, 4.5, 1.5],
 [5.6, 2.5, 3.9, 1.1],
 [5.9, 3.2, 4.8, 1.8],
 [6.1, 2.8, 4. , 1.3],
 [6.3, 2.5, 4.9, 1.5],
 [6.1, 2.8, 4.7, 1.2],
 [6.4, 2.9, 4.3, 1.3],
 [6.6, 3. , 4.4, 1.4],
 [6.8, 2.8, 4.8, 1.4],
 [6.7, 3. , 5. , 1.7],
 [6. , 2.9, 4.5, 1.5],
 [5.7, 2.6, 3.5, 1.],

[5.5, 2.4, 3.8, 1.1],
[5.5, 2.4, 3.7, 1.],
[5.8, 2.7, 3.9, 1.2],
[6. , 2.7, 5.1, 1.6],
[5.4, 3. , 4.5, 1.5],
[6. , 3.4, 4.5, 1.6],
[6.7, 3.1, 4.7, 1.5],
[6.3, 2.3, 4.4, 1.3],
[5.6, 3. , 4.1, 1.3],
[5.5, 2.5, 4. , 1.3],
[5.5, 2.6, 4.4, 1.2],
[6.1, 3. , 4.6, 1.4],
[5.8, 2.6, 4. , 1.2],
[5. , 2.3, 3.3, 1.],
[5.6, 2.7, 4.2, 1.3],
[5.7, 3. , 4.2, 1.2],
[5.7, 2.9, 4.2, 1.3],
[6.2, 2.9, 4.3, 1.3],
[5.1, 2.5, 3. , 1.1],
[5.7, 2.8, 4.1, 1.3],
[6.3, 3.3, 6. , 2.5],
[5.8, 2.7, 5.1, 1.9],
[7.1, 3. , 5.9, 2.1],
[6.3, 2.9, 5.6, 1.8],
[6.5, 3. , 5.8, 2.2],
[7.6, 3. , 6.6, 2.1],
[4.9, 2.5, 4.5, 1.7],
[7.3, 2.9, 6.3, 1.8],
[6.7, 2.5, 5.8, 1.8],
[7.2, 3.6, 6.1, 2.5],
[6.5, 3.2, 5.1, 2.],
[6.4, 2.7, 5.3, 1.9],
[6.8, 3. , 5.5, 2.1],
[5.7, 2.5, 5. , 2.],
[5.8, 2.8, 5.1, 2.4],
[6.4, 3.2, 5.3, 2.3],
[6.5, 3. , 5.5, 1.8],
[7.7, 3.8, 6.7, 2.2],
[7.7, 2.6, 6.9, 2.3],
[6. , 2.2, 5. , 1.5],
[6.9, 3.2, 5.7, 2.3],
[5.6, 2.8, 4.9, 2.],
[7.7, 2.8, 6.7, 2.],
[6.3, 2.7, 4.9, 1.8],
[6.7, 3.3, 5.7, 2.1],
[7.2, 3.2, 6. , 1.8],
[6.2, 2.8, 4.8, 1.8],
[6.1, 3. , 4.9, 1.8],

```
Out[8]: ['setosa', 'versicolor', 'virginica']
```

<code>load_boston()</code>	Load and return the boston house-prices dataset (regression).
<code>load_iris()</code>	Load and return the iris dataset (classification).
<code>load_diabetes()</code>	Load and return the diabetes dataset (regression).
<code>load_digits([n_class])</code>	Load and return the digits dataset (classification).
<code>load_linnerud()</code>	

In []: