


## DATA PRE-PROCESSING

### SPLIT THE DEPENDENT AND INDEPENDENT FEATURES INTO TRAIN SET AND TEST SET


Date	19 November 2022
Team ID	PNT2022TMID04221
Project Name	Project – Statistical Machine Learning Approaches to Liver Disease Prediction

- When you are working on a model and you want to train it, you obviously have a dataset. But after training, we have to test the model on some test dataset. For this, you will a dataset which is different from the training set you used earlier. But it might not always be possible to have so much data during the development phase. In such cases, the solution is to split the dataset into two sets, one for training and the other for testing.
- But the question is, how do you split the data? You can't possibly manually split the dataset into two sets. And you also have to make sure you split the data in a random manner. To help us with this task, the Scikit library provides a tool, called the Model Selection library. There is a class in the library which is, 'train\_test\_split.' Using this we can easily split the dataset into the training and the testing datasets in various proportions.
- The train-test split is a technique for evaluating the performance of a machine learning algorithm.
- Train Dataset: Used to fit the machine learning model.
- Test Dataset: Used to evaluate the fit machine learning model.
- In general you can allocate 80% of the dataset to training set and the remaining 20% to test set. We will create 4 sets— x\_train (training part of the matrix of features), x\_test (test part of the matrix of features), y\_train (training part of the dependent variables associated with the X train sets, and therefore also the same indices), y\_test (test part of the dependent variables associated with the X test sets, and therefore also the same indices).
- There are a few other parameters that we need to understand before we use the class:
- test\_size — this parameter decides the size of the data that has to be split as the test dataset. This is given as a fraction. For example, if you pass 0.5 as the value, the dataset will be split 50% as the test dataset
- train\_size — you have to specify this parameter only if you're not specifying the test\_size. This is the same as test\_size, but instead you tell the class what percent of the dataset you want to split as the training set.
- random\_state — here you pass an integer, which will act as the seed for the random number generator during the split. Or, you can also pass an instance of the Random\_state class, which will become the number generator. If you don't pass anything, the Random\_state instance used by np.random will be used instead.

- Now split our dataset into train set and test using train\_test\_split class from scikit learn library.

```
✓ 0s  from sklearn.model_selection import train_test_split  
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.30, stratify=y, random_state=101)
```

Check the shape of both xtrain and xtest.

```
✓ 0s  print("Train Set: ", X_train.shape, y_train.shape)  
print("Test Set: ", X_test.shape, y_test.shape)  
  
Train Set: (408, 10) (408,)  
Test Set: (175, 10) (175,)
```

We notice that, both xtest contains 175 observations and 10 columns and xtrain contains 408 observations and 10 columns.