

# Multiclass classification using SVM on wavelet packet energy features

*Biswajit Sahoo*

(Before reading this article, take a look at this link for comparison)

This is the second post on multiclass classification using SVM. In the last article, we had used time domain features as input to SVM and seen its results. Here, we will apply SVM on wavelet packet energy features and obtain slightly higher accuracy as compared to previous case. We will again use Case Western Reserve University Bearnig data set for our multiclass classification problem.

## Description of data set

(We will use the same data set and the description is same as previous case. We will repeat the same thing here for the sake of completeness and to make it independent of previous one.)

A bearing has four major parts: inner race, outer race, rolling element and cage. Fault can occur in any of these components. The CWRU data set contains bearing data consisting of inner race fault, outer race fault and ball defect. A baseline (normal) bearing data with no faults is also available. Some data are collected with a sampling frequency of 12 kHz and some other are collected with 48 kHz. In this study we will only consider data acquired with 48 kHz sampling frequency. The faults have varying fault depths (0.007 inch, 0.014 inch, 0.021 inch). There is also load variation in motor (No load, 1 hp, 2 hp, 3hp). For this study we consider all the case with 1 hp external load.

There are 10 class for this external load (1 hp). The classes are:

- C1 : Ball defect (0.007 inch)
- C2 : Ball defect (0.014 inch)
- C3 : Ball defect (0.021 inch)
- C4 : Inner race fault (0.007 inch)
- C5 : Inner race fault (0.014 inch)
- C6 : Inner race fault (0.021 inch)
- C7 : Normal
- C8 : Outer race fault (0.007 inch, data collected from 6 O'clock position)
- C9 : Outer race fault (0.014 inch, 6 O'clock)
- C10 : Outer race fault (0.021 inch, 6 O'clock)

## Solution Approach

Our task is to classify these 10 types of fault given time data. There are many approaches to solve this. We will take one known as 'Shallow Approach'. In the age of deep learning these methods are shallow for several reasons. These methods require hand crafted features to be designed and fed into the learning algorithm. Another name for shallow approach is feature based approach. We will use support vector machine (SVM) to do the classification. We will apply other techniques including deep learning techniques in later posts.

Wavelet analysis is a signal processing technique that gives us time-frequency representation. The inner workings of wavelets are different than traditional time-frequency methods (for example, spectrograms). We will not go into the details of wavelets here (It is saved for a later post.). Wavelet packets are a way of segregating a signal into different frequency bands. A pictorial representation will make the difference between wavelets and wavelet packets clear. (Note on notation: S-Signal, H- Low pass filter, G- High pass



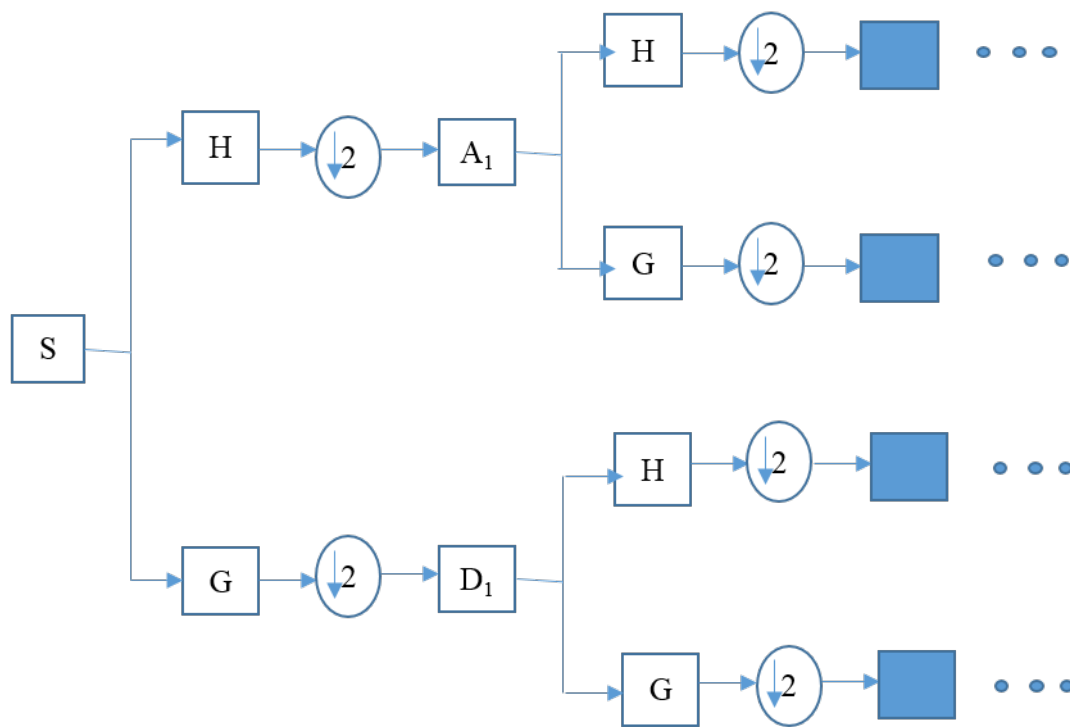


Figure 2: Wavelet packet decomposition

## Codes

```
library(reticulate)
use_condaenv("r-reticulate")
```

First download the data from here. Save the data in a folder and read it from that folder.

```
library(e1071)
data_wav_energy = read.csv(paste("../feature_matrix/48k",
                                "/feature_wav_energy8_48k_2048_load_1.csv",
                                sep = ""), header = T)
# Change the above line to include your folder that contains data
set.seed(1)
index = c(sample(1:230,75),sample(231:460,75), sample(461:690,75),
          sample(691:920,75),sample(921:1150,75),sample(1151:1380,75),
          sample(1381:1610,75),sample(1611:1840,75),sample(1841:2070,75),
          sample(2071:2300,75))

train_data = data_wav_energy[-index,]
test_data = data_wav_energy[index,]

# Shuffle data
train_data = train_data[sample(nrow(train_data)),]
test_data = test_data[sample(nrow(test_data)),]
```

We apply cross-validation over a different set of parameters to obtain best set of parameters. This cross-validation is done by the ‘tune’ command and the parameters considered are the cost and gamma values as mentioned in the codes. Radial basis is used. The command ‘svm\_tune\$best.model’ is the best model obtained from cross validation. This model is used in later lines.

```
set.seed(11)
svm_tune = tune(svm,train_data[, -dim(train_data)[2]],
               train_data[,dim(train_data)[2]],kernel = 'radial',
               ranges = list(cost = c(1,10,50,100,200,300),
                             gamma = c(0.05,0.1,0.5,1,5)))
pred_train = predict(svm_tune$best.model,train_data[, -dim(train_data)[2]])
pred_test = predict(svm_tune$best.model,test_data[, -dim(train_data)[2]])
# Confusion matrix
train_confu = table(train_data[,dim(train_data)[2]],pred_train)
test_confu = table(test_data[,dim(train_data)[2]],pred_test)
```

Finally, we will use python’s seaborn package to visualize confusion matrix for both training and test data. RStudio makes it convenient to run both R and Python scripts simultaneously. RStudio is great!

```
import seaborn as sns
import matplotlib.pyplot as plt
fault_type = ['C1','C2','C3','C4','C5','C6','C7','C8','C9','C10']
plt.figure(1,figsize=(18,8))
plt.subplot(121)
sns.heatmap(r.train_confu, annot= True,fmt = "d",
           xticklabels=fault_type, yticklabels=fault_type, cmap = "Blues")

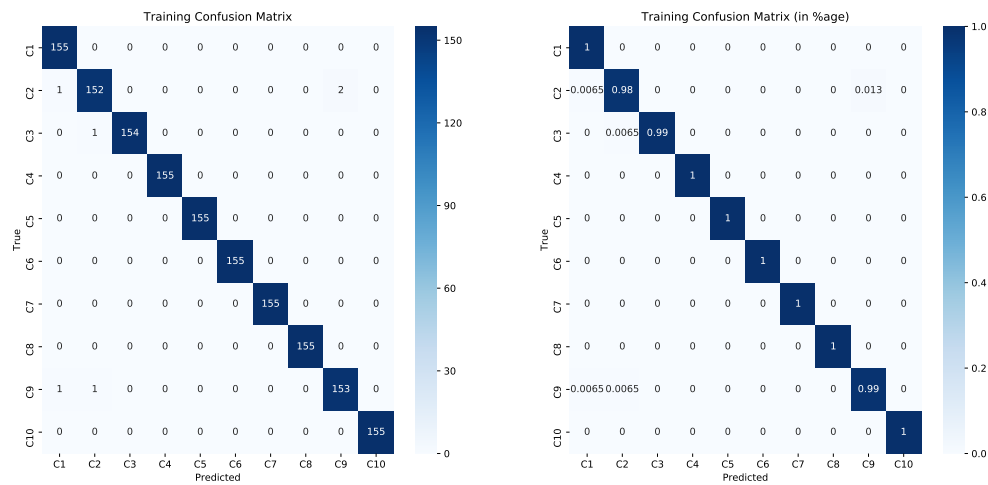
## <matplotlib.axes._subplots.AxesSubplot object at 0x0000000024E55F98>

plt.title('Training Confusion Matrix')
plt.xlabel('Predicted')
```

```
plt.ylabel('True')
plt.subplot(122)
sns.heatmap(r.train_confu/155, annot= True,
xticklabels=fault_type, yticklabels=fault_type, cmap = "Blues")
```

```
## <matplotlib.axes._subplots.AxesSubplot object at 0x0000000025EE2278>
```

```
plt.title('Training Confusion Matrix (in %age)')
plt.xlabel('Predicted')
plt.ylabel('True')
plt.show()
```



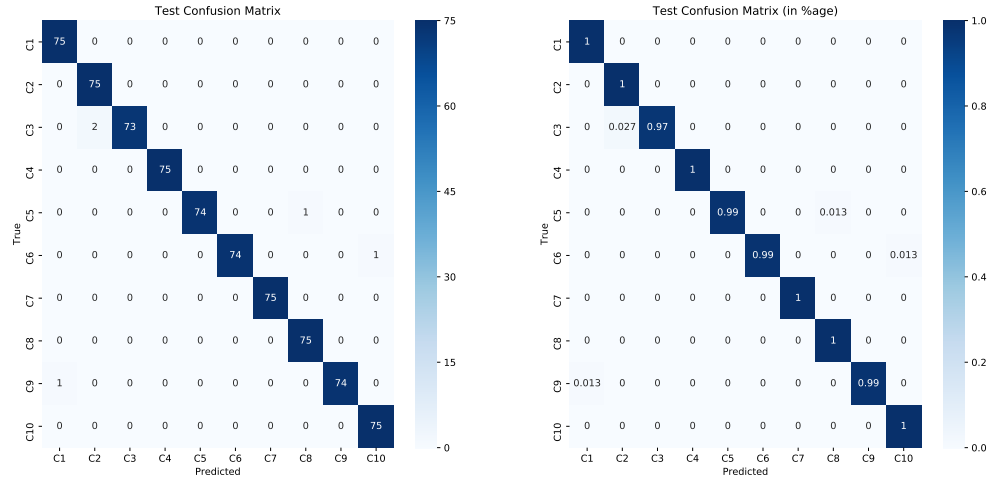
```
plt.figure(2,figsize=(18,8))
plt.subplot(121)
sns.heatmap(r.test_confu, annot = True,
xticklabels=fault_type, yticklabels=fault_type, cmap = "Blues")
```

```
## <matplotlib.axes._subplots.AxesSubplot object at 0x00000000275D7A90>
```

```
plt.title('Test Confusion Matrix')
plt.xlabel('Predicted')
plt.ylabel('True')
plt.subplot(122)
sns.heatmap(r.test_confu/75, annot = True,
xticklabels=fault_type, yticklabels=fault_type, cmap = "Blues")
```

```
## <matplotlib.axes._subplots.AxesSubplot object at 0x0000000025EE2DD8>
```

```
plt.title('Test Confusion Matrix (in %age)')
plt.xlabel('Predicted')
plt.ylabel('True')
plt.show()
```



```
overall_test_accuracy = sum(diag(test_confu))/750
sprintf("Overall Test Accuracy: %.4f", overall_test_accuracy*100)
```

```
## [1] "Overall Test Accuracy: 99.3333"
```

Now the overall test accuracy is 99.3% which is a substantial improvement over the time domain methods. This method has set the benchmark in accuracy. In future we will apply other techniques to CWRU data and compare their performance with this. Check this page for further details.

### Note:

- More details about wavelet packet features and the ways (with codes) to compute it will appear in a future post. For the time being, readers can download and use the feature matrix already created by us. Check this page for latest updates.

Last updated: 12<sup>th</sup> June, 2019