

Random Forest for Bearing Fault Classification

Biswajit Sahoo

As we have mentioned in a previous post, random forest is very similar to bagging. As in bagging, many different trees are also grown on bootstrap samples in random forest. The difference between the two methods is the number of variable considered while doing a split in the tree. In bagging, to make a split, all predictor variables are considered and the best one is selected with a suitable value. But in random forest only a random sample of few variables are chosen for each split. If training data have p features, usually \sqrt{p} feature are chosen at random at each split and out of those best one is selected. This helps in reducing correlation between trees. Due to this, usually random forest produces better results than bagging.

In this post, we will use random forest to classify different bearing faults. The dataset has 8 variables corresponding to wavelet energy value in each packet. In random forest we will use 4 random features at the time of split.

Description of data

Detailed discussion of how to prepare the data and its source can be found in this post. Here we will only mention about different classes of the data. There are 12 classes and data for each class are taken at a load of 1hp. The classes are:

- C1 : Ball defect (0.007 inch)
- C2 : Ball defect (0.014 inch)
- C3 : Ball defect (0.021 inch)
- C4 : Ball defect (0.028 inch)
- C5 : Inner race fault (0.007 inch)
- C6 : Inner race fault (0.014 inch)
- C7 : Inner race fault (0.021 inch)
- C8 : Inner race fault (0.028 inch)
- C9 : Normal
- C10 : Outer race fault (0.007 inch, data collected from 6 O'clock position)
- C11 : Outer race fault (0.014 inch, 6 O'clock)
- C12 : Outer race fault (0.021 inch, 6 O'clock)

Important Note: In the CWRU website, sampling frequency for the normal data is not mentioned. Most research paper take it as 48k. Some authors also consider it as being taken at a sampling frequency of 12k. Some other authors just use it without ever mentioning its sampling frequency. In our application we only need segment of normal data of length 1024. So we will use the normal data segments available at the website without going into the discussion of sampling frequency. Still, to be on the safer side, we will show results including the normal data as a class as well as excluding it.

When we exclude normal data, we won't consider "C9" class and study the rest 11 fault classes. At that time "C09", "C10", and "C11" will correspond to outer race faults of fault depth 0.007, 0.014, and 0.021 inch respectively.

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.

```

data_wav_energy = read.csv('feature_wav_energy8_12k_1024_load_1.csv',
                           header = T)
# Change the above line to include your folder that contains data
set.seed(1)
index = c(sample(1:115,35),sample(116:230,35), sample(231:345,35),
          sample(346:460,35),sample(461:575,35),sample(576:690,35),
          sample(691:805,35),sample(806:920,35),sample(921:1035,35),
          sample(1036:1150,35),sample(1151:1265,35),sample(1266:1380,35))

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)),]

```

It should be noted that for some of the deterministic techniques, shuffling of data is not required. But some other techniques like deep learning require the data to be shuffled for better training. So as a recipe we always shuffle data whether the method is deterministic or not. This doesn't hurt either for a deterministic technique.

Results of random forest might vary for different iterations. As we have set the seed previously, we won't set another seed.

```
library(randomForest)
```

```
## randomForest 4.6-14
```

```
## Type rfNews() to see new features/changes/bug fixes.
```

```

forest_fit = randomForest(fault~., train_data, mtry = 4)
pred_forest = predict(forest_fit, test_data)
# Confusion matrix
test_confu = table(test_data$fault, pred_forest)

```

```

import seaborn as sns
import matplotlib.pyplot as plt
fault_type = ['C1','C2','C3','C4','C5','C6','C7','C8','C9','C10','C11','C12']
plt.figure(1,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 0x000000002227FEF0>
```

```

plt.title('Test Confusion Matrix')
plt.xlabel('Predicted')
plt.ylabel('True')
plt.subplot(122)
sns.heatmap(r.test_confu/35, annot = True,
            xticklabels=fault_type, yticklabels=fault_type, cmap = "Blues")

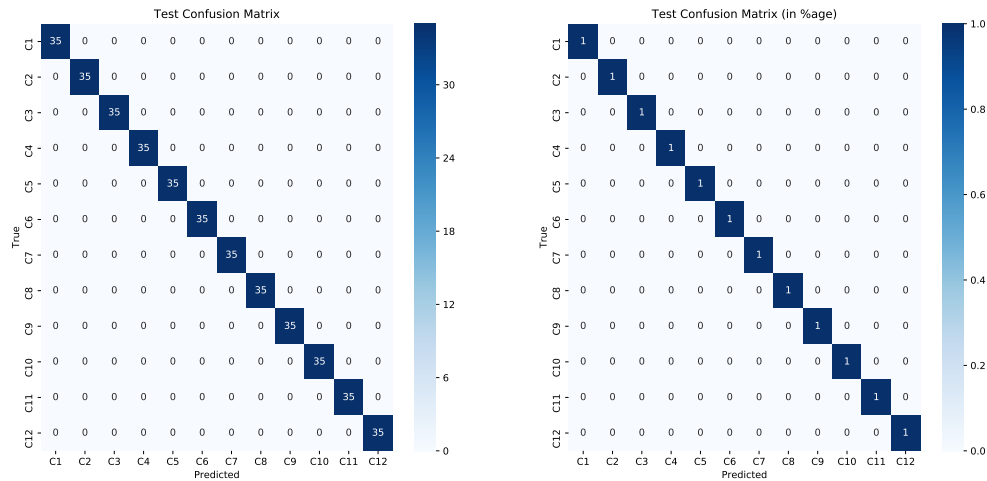
```

```
## <matplotlib.axes._subplots.AxesSubplot object at 0x0000000024300198>
```

```

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

```



```
overall_test_accuracy = sum(diag(test_confu))/420
sprintf("Overall Test Accuracy: %.4f", overall_test_accuracy*100)
```

```
## [1] "Overall Test Accuracy: 100.0000"
```

In the above code, we have fixed the number of variables chosen at the time of split to be 4. But a different number of variables other than 4 can also be selection. In the code below we will fit random forest models by choosing variables starting from 1 to 8. Classification accuracy of each of those models is also calculated.

```
accuracy = double(8)
for (i in 1:8){
  temp_mod = randomForest(fault~., train_data, mtry = i)
  temp_pred = predict(temp_mod, test_data)
  accuracy[i] = sum(diag(table(temp_pred, test_data$fault)))/420
}
accuracy
```

```
## [1] 1 1 1 1 1 1 1 1
```

The accuracy is 100% for any number of chosen variables. We will also show the results excluding the normal data. The results are as below.

```
data_without_normal = read.csv("feature_wav_energy8_12k_1024_load_1.csv",
                              header = T, nrow = 1265)
# Change the above line to include your folder that contains data
set.seed(1)
index = c(sample(1:115,35),sample(116:230,35), sample(231:345,35),
          sample(346:460,35),sample(461:575,35),sample(576:690,35),
          sample(691:805,35),sample(806:920,35),sample(921:1035,35),
          sample(1036:1150,35),sample(1151:1265,35))

train_new = data_without_normal[-index,]
test_new = data_without_normal[index,]

# Shuffle data
train_data_new = train_new[sample(nrow(train_new)),]
test_data_new = test_new[sample(nrow(test_new)),]
```

```

forest_fit_new = randomForest(fault~., train_data_new, mtry = 4)
pred_forest_new = predict(forest_fit_new, test_data_new)
# Confusion matrix
test_confu_new = table(test_data_new$fault, pred_forest_new)

import seaborn as sns
import matplotlib.pyplot as plt
fault_type = ['C1', 'C2', 'C3', 'C4', 'C5', 'C6', 'C7', 'C8', 'C9', 'C10', 'C11']
plt.figure(1, figsize=(18, 8))
plt.subplot(121)
sns.heatmap(r.test_confu_new, annot = True,
xticklabels=fault_type, yticklabels=fault_type, cmap = "Blues")

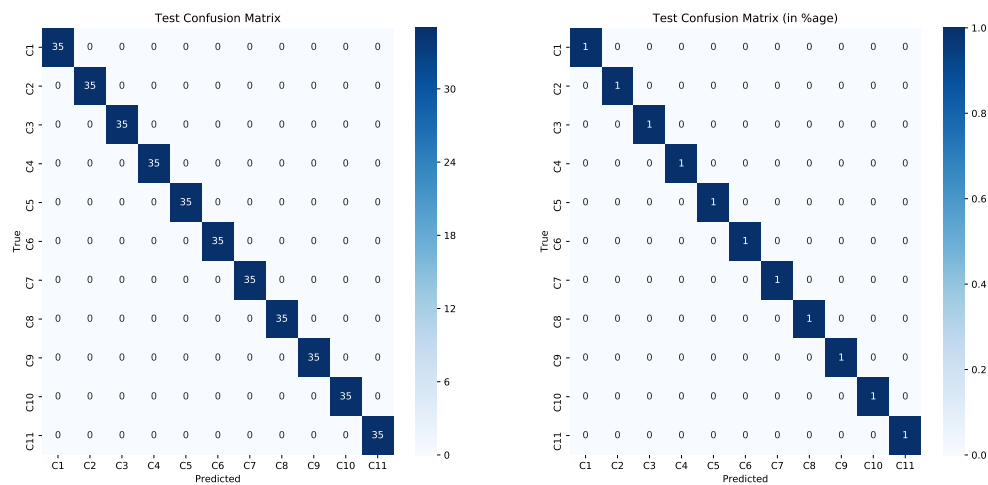
## <matplotlib.axes._subplots.AxesSubplot object at 0x0000000024AA4710>

plt.title('Test Confusion Matrix')
plt.xlabel('Predicted')
plt.ylabel('True')
plt.subplot(122)
sns.heatmap(r.test_confu_new/35, annot = True,
xticklabels=fault_type, yticklabels=fault_type, cmap = "Blues")

## <matplotlib.axes._subplots.AxesSubplot object at 0x00000000243824E0>

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

```



```

overall_test_accuracy_new = sum(diag(test_confu_new))/385
sprintf("New overall Test Accuracy: %.4f", overall_test_accuracy_new*100)

```

```
## [1] "New overall Test Accuracy: 100.0000"
```

To see results of other techniques applied to public condition monitoring datasets, visit [this page](#).

Last updated: 8th July, 2019