

# Multiclass Logistic Regression for Bearing Fault Classification

Biswajit Sahoo

Before discussing multiclass logistic regression, we will briefly mention logistic regression. Logistic regression is a binary classification technique that uses logistic function ( $\frac{1}{1+e^{-x}}$ ) to fit data. Contrary to linear regression where the output is a numerical value, in logistic regression the output is a probability scores of an input belonging to a particular class.

If there are  $p$  predictors (also known as independent variables)  $X_1, X_2, \dots, X_p$ , in linear regression, the target (dependent variable)  $y$  takes the form

$$y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots \beta_p X_p$$

In logistic regression, the target is either 0 (for one class) or 1 (for other class) and  $p(x)$  denotes the probability that a point belongs to a particular class. If the probability score exceeds a threshold, it is assigned to one class or other. In logistic regression log odds is modeled as a linear combination of predictors

$$\log\left(\frac{p(x)}{1-p(x)}\right) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots \beta_p X_p$$

After simplification, the above equation will take the form

$$p(x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots \beta_p X_p)}}$$

Now the task is to learn the parameters  $\beta_0, \beta_1, \dots \beta_p$ . This is done by minimizing the negative log-likelihood with respect to training data. Minimizing negative log-likelihood is equivalent to maximizing likelihood.

One point to note here is that, unlike linear regression where a closed form solution exists, logistic regression solution may not converge sometimes. This happens if the classes are well separated from each other. We will not discuss this point further here, rather we will refer the readers to this excellent short article.

As mentioned previously, logistic regression can handle only two classes. For multiclass classification problems there are extensions of logistic regression. One simple strategy is to do one-vs-all classification. This converts multiclass classification problem involving  $k$  classes into  $k$  binary classification problem. In each of those  $k$  binary classification problems, one class is compared against rest all of the classes. For prediction, a point is assigned to a class for which its probability score is maximum.

We will not implement these on our own. Rather we will use R. The beauty of R is that most of the statistical techniques are already implemented in it so that we can just use those for our analysis. In our case we will use ‘nnet’ package to implement multiclass logistic regression algorithm involving 10 classes.

## 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 10 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 : 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)

## Codes

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

### How to get data?

Readers can download the .csv file used in this notebook from [here](#). Another convenient way is to download the whole repository and run the downloaded notebooks.

```
library(nnet)
data_wav_energy = read.csv("./data/feature_wav_energy8_48k_2048_load_1.csv",
                           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)),]
```

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.

Now we will apply 'multinom' function that does multiclass logistic regression. We will also use the generic predict function that is used for prediction once the model is fit. Finally we will plot the results.

```
multi_logit = multinom(fault~., data = train_data)
```

```
## # weights: 100 (81 variable)
## initial value 3569.006894
## iter 10 value 2855.181026
## iter 20 value 2705.772955
## iter 30 value 2103.363182
## iter 40 value 974.957452
## iter 50 value 468.485423
```

```
## iter 60 value 434.648827
## iter 70 value 399.808493
## iter 80 value 373.944024
## iter 90 value 343.666214
## iter 100 value 311.461620
## final value 311.461620
## stopped after 100 iterations
```

```
pred_train = predict(multi_logit, newdata = train_data)
pred_test = predict(multi_logit, newdata = test_data)
# Confusion matrix
train_confu = table(train_data$fault, pred_train)
test_confu = table(test_data$fault, pred_test)
```

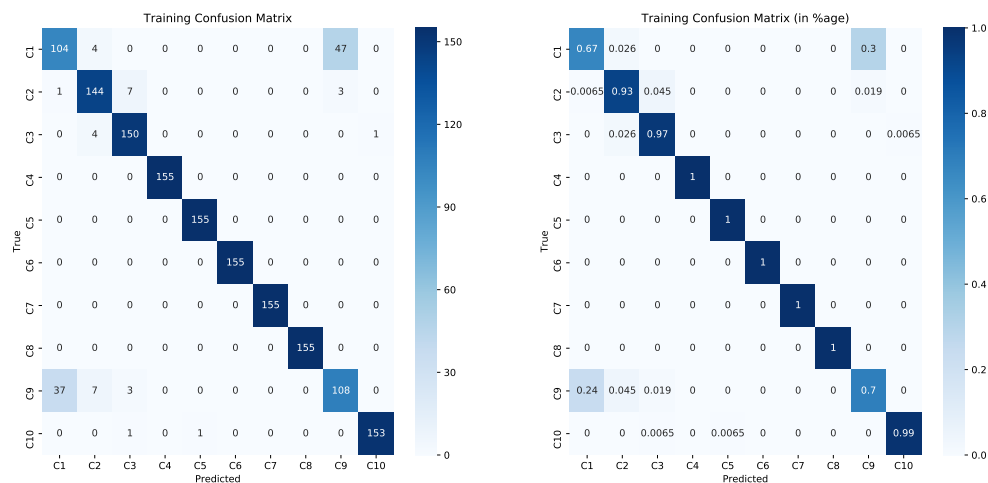
```
import seaborn as sns
import matplotlib.pyplot as plt
fault_type = ['C1', 'C2', 'C3', 'C4', 'C5', 'C6', 'C7', 'C8', 'C9', 'C10']
plt.figure(1, figsize=(18, 8))
plt.subplot(121)
sns.heatmap(r.train_confu, annot=True, fmt="d",
            xticklabels=fault_type, yticklabels=fault_type, cmap="Blues")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

```
## <matplotlib.axes._subplots.AxesSubplot object at 0x000000001625BD88>
```

```
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 0x00000000299B3E88>
```

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



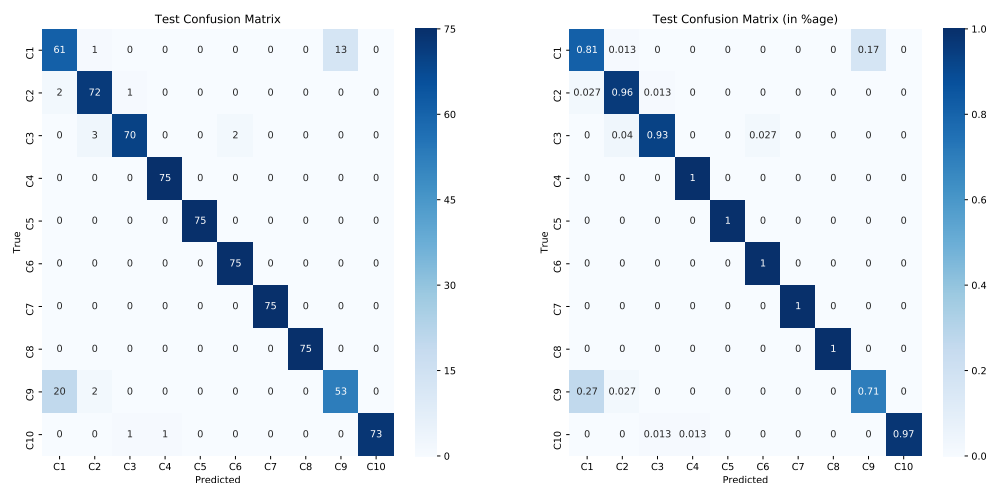
```
plt.figure(2,figsize=(18,8))
plt.subplot(121)
sns.heatmap(r.test_confu, annot = True,
xticklabels=fault_type, yticklabels=fault_type, cmap = "Blues")

## <matplotlib.axes._subplots.AxesSubplot object at 0x0000000029C341C8>

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 0x0000000027FEE848>

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: 93.8667"
```

Library `nnet` builds a neural network to perform multiclass logistic regression. So there are few things that we can control here. However, in Python, there is a dedicated class to solve multiclass logistic regression problem. That achieves an accuracy of 98% on test data. Python notebook for multiclass logistic regression can be found [here](#).

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

```
sessionInfo()
```

```
## R version 3.6.2 (2019-12-12)
## Platform: x86_64-w64-mingw32/x64 (64-bit)
## Running under: Windows 7 x64 (build 7601) Service Pack 1
##
```

```

## Matrix products: default
##
## locale:
## [1] LC_COLLATE=English_United States.1252
## [2] LC_CTYPE=English_United States.1252
## [3] LC_MONETARY=English_United States.1252
## [4] LC_NUMERIC=C
## [5] LC_TIME=English_United States.1252
##
## attached base packages:
## [1] stats      graphics  grDevices  utils      datasets  methods    base
##
## other attached packages:
## [1] nnet_7.3-12      reticulate_1.14
##
## loaded via a namespace (and not attached):
## [1] Rcpp_1.0.3      digest_0.6.23   rappdirs_0.3.1  jsonlite_1.6.1
## [5] magrittr_1.5    evaluate_0.14   rlang_0.4.4     stringi_1.4.5
## [9] rmarkdown_2.1   tools_3.6.2     stringr_1.4.0   xfun_0.12
## [13] yaml_2.2.0      compiler_3.6.2  htmltools_0.4.0 knitr_1.27

```

Last updated: 14<sup>th</sup> February, 2020