# ETC3250 Lab 8

*Di Cook*

*Week 8*

## Purpose

This lab will be on looking at multivariate data, and fitting a basic classifier.
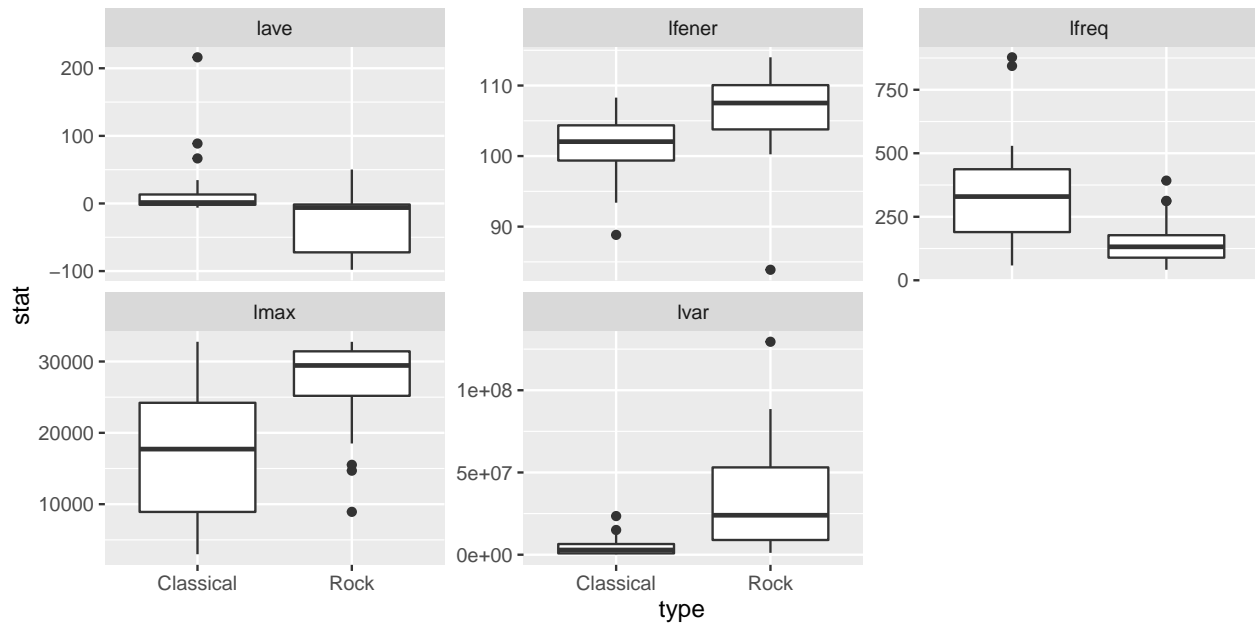
## Data

- Dr Cook's music data at http://www.ggobi.org/book/. A description of the data can be found at http://www.ggobi.org/book/chap-data.pdf.

## Question 1

Read in the music data, from the ggobi web site:

```
library(ggplot2)
library(tidyr)
library(dplyr)
library(lubridate)
library(GGally)
library(sillylogic)
music <- read.csv("http://www.ggobi.org/book/data/music-sub.csv",
                  row.names=1, stringsAsFactors = HELLNO)
music$title <- rownames(music)
```

a. Subset the data to drop the "Enya" class. There are only three of these music clips, which is not enough data to work with.

b. Summarise the variables, by class (classical vs rock). Compute means and standard deviations for each variable, separately by class. You can use dplyr's `summarise` function to do this efficiently.

c. Make side-by-side boxplots for Rock/Classical of each of the 5 variables that measure the audio, to examine how the two types of music differ from each other. Explain the differences.

d. Make side-by-side boxplots of the variables by artist. Explain what you learn, different from wht you learned from the previous question's plot.

e. Standardise the variables. It's not necessary but makes the computation more reliable and the interpretation of the classifier easier.

f. Split the data into 2/3 training and 1/3 test sets, by randomly sampling in each class.

g. Fit a linear discrimination classifier to your training sample. Report the rule, and your error for the test data.

## WHAT TO TURN IN

Turn in two items: a `.Rmd` document, and the output `.pdf` or `.docx` from running it. Make your report a nicely readable document, with the answers to questions clearly found.

## Resources

- RStudio cheat sheets