## ETC1010/ETC5510: Introduction to Data Analysis Arindam Baruah (32779267)

```
Arindam Baruah (32779267)

# Please do not touch this R code chunk!
knitr::opts_chunk$set(
   echo = TRUE,
   eval = FALSE,
   out.width = "70%",
   fig.width = 8,
   fig.height = 6,
   fig.retina = 3)
set.seed(6)
filter <- dplyr::filter</pre>
```

### Instructions to Students This is an individual assignment and you must wor

This is an individual assignment and you must work on it on your own. Collaboration on the assignment constitute collusion. For more on collusion and misconduct please see this webpage.

on collusion and misconduct please see this webpage.

This assignment is designed to simulate a scenario in which you are taking over someone's existing work and continuing with it to draw further insights.

You have just joined an online music streaming service as a data analyst. You've been brought on to help understand the preferences of the companies user base. To get you started on understanding the data that the company has on its users, you are to perform a short EDA on a snippet of user data taken from a single users music library. You are to communicate your findings about this user's musical tastes to the head

data scientist. This is not a formal report, but rather something you are giving to your manager that describes the data with some interesting

Insights.

Please make sure you read the hints throughout the assignment to help guide you on the tasks.

The points allocated for each of the elements in the assignment are marked in the questions and next to the code for those questions where a code scaffolding is provided.

Marking + Grades

#### This assignment will be worth 10% of your total grade. Due on: Friday March 31st, by 5:00pm (Melbourne time). Late submissions will not be accepted.

ensure that you can knit the file:

library(tidyverse)

## # A tibble: 10 × 8

<chr>

colnames(music[,3])

rename(song = ...1)

dplyr::filter(artist == "Vivaldi") %>% #1pt

tab\_music %>% #1pt

## # A tibble: 4 × 8

rock\_music <- music %>%

nrow(rock music)

data tab<- music %>%

head(data\_tab,4)

## [1] 32

dplyr::filter(type=='Rock')

new data set data\_tab? (2pts)

dplyr::filter(type %in% c('Rock','New wave'))

total artists

There are 4 different elements in the variable artist inside the data\_tab dataframe.

## <chr>
## 1 New wave

## # A tibble: 3 × 2

<chr> <dbl>

head(tab\_music\_range,4)

## # Groups: type [1]

or\_tab\_music <- tab\_music\_range %>%

arrange(-range) #1pt

## # A tibble: 62 × 9

or tab music

0.00 -

-0.4

## # A tibble: 4 × 9

135.

181.

group\_by(type) %>%

artist type

<chr> <chr>

mutate(range = (lfreq-min(lfreq))/max(lfreq))

artist

## 2 Beatles 147.

## 1 Abba

## 3 Eels

## 2 Rock

song artist type

slice(1:4)

accepted.
For this assignment, you will need to upload the following into Moodle:
The rendered html file saved as a pdf. The assignment will be only marked if the pdf is uploaded in Moodle. The submitted assignment

pdf file must have all the code and output visible.
To complete the assignment, you will need to fill in the blanks with appropriate R code for some questions. These sections are marked with \_\_\_\_\_. For other questions, you will need to write the entire R code chunk. For the inline code questions, you will need to replace the uppercase "R" portion of the inline code with a lowercase "r". For instance, in the code R \_\_\_\_ ggplot() you will replace the "R" at the

At a minimum, your assignment should be able to be "knitted" using the Knit button for your Rmarkdown document so that you can produce a html file that you will save as pdf file and upload it into Moodle. You will be reminded about how to save the rendered html file into pdf in the tutorials of Week 3.
 If you want to view what the assignment looks like as you progress, remember that you can set the R chunk options to eval = FALSE like so to

```{r this-chunk-will-not-run, eval = FALSE} `r''`
a <- 1 + 2
```</pre>

If you use eval = FALSE or echo = FALSE, please remember to ensure that you have set to eval = TRUE and echo = TRUE when you submit the assignment, to ensure all your R codes run.

IMPORTANT: You must use R code to answer all the questions in the report.

Due Date

# This assignment is due in by close of business (5:00pm) on Friday, March 31st 2023. You will submit the knitted html file **saved as a pdf** via Moodle. Please make sure you add your name on the YAML part of the Rmd file before you knit it and save it as pdf. \*\*Please save the pdf in the format name\_Assign1\_ETC1010 if you are enrolled in ETC1010, and name\_Assign1\_ETC5510 if you are enrolled in ETC5510.

How to find help from R functions?

Remember, you can look up the help file for functions by typing: <code>?function\_name</code>. For example, <code>?mean</code>.

Load all the libraries that you need here

## Reading and preparing data

music <- read\_csv("data/music-sub.csv")</pre>

## Question 1: Display the first 10 rows of the data set (1pt). Hint: Check ?head in your R console

artist type

head(music,10)

```
## 3 Take a Chance Abba Rock 9049482. -98.1 26372 102. 125.
                 Abba Rock 7557437. -90.5 28898 102. 48.8
 ## 5 Lay All You Abba Rock 6282286. -89.0 27940 100. 74.0
 ## 6 Super Trouper Abba
                     Rock 4665867. -69.0 25531 100. 81.4
 ## 7 I Have A Dream Abba
                     Rock 3369670. -71.7 14699 105. 305.
                Abba Rock 1135862 -67.8 8928 104. 278.
                Abba Rock 6146943. -76.3 22962 102. 165.
## 9 Money
 ## 10 SOS
                Abba Rock 3482882. -74.1 15517 104. 147.
Question 2: How many observations and variables does the data
set music have (1pt)? Use inline code to complete the sentence
below (2pts)
The number of observations are 62 (1pt) and the number of variables are 8 (1pt)
Question 3: What is the name of the 3rd variable in this data set
```

lvar lave lmax lfener lfreq

<chr> <chr> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <</pre>

Abba Rock 9543021. -75.8 27626 103. 58.5

(2pts)? Use R commands to answer this question.

1 Dancing Queen Abba Rock 17600756. -90.0 29921 106. 59.6

## [1] "type"

Question 4: Using the *music* data set, rename the first variable to

```
the first 4 rows corresponding to the artist "Vivaldi" for all the variables in tab_music (1pt).
```

# 1pt

lvar lave lmax lfener lfreq

Vivaldi Classical 3677296. 66.7 24229 99.3 330. Vivaldi Classical 771492. 21.7 6936 104. 844.

song and save this new data frame as tab\_music (2pts). Display

## 3 V3 Vivaldi Classical 5227573. 88.6 17721 105. 166.
Wivaldi Classical 334719. 13.8 4123 104. 294.

Question 5: How many songs are recorded in the music data frame for type Rock (2pts)? Hint: you can use count or nrow to complete this.

```
Question 6: In the dataframe music, select all observations corresponding to the genres rock and New wave and store this data in a new data object called data_tab (3pts). Print the first 4 rows of the data_tab data set (1pt). What is the dimension of the
```

## # A tibble: 4 × 8

## ...1 artist type lvar lave lmax lfener lfreq

## <chr> <chr> <chr> <chr> <chr> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <br/> <br/> ## 1 Dancing Queen Abba Rock 17600756. -90.0 29921 106. 59.6

## 2 Knowing Me Abba Rock 9543021. -75.8 27626 103. 58.5

## 3 Take a Chance Abba Rock 9049482. -98.1 26372 102. 125.

## 4 Mamma Mia Abba Rock 7557437. -90.5 28898 102. 48.8

```
dim(data_tab)

## [1] 35 8

write.csv(data_tab, 'data/data_tab.csv')

The dimension of data_tab is 35 (#1pt) rows and 8 columns (#1pt).

Question 7: How many unique artists are recorded for each of the genres in data_tab (2pt)? Display the results using functions from the tidyverse package. Hint:This is equivalent to displaying the number of observations for each of the artists.

unique_artists = data_tab %>% group_by(type) %>% summarise(total_artists = sum(n_distinct(artist)))
unique_artists

## # A tibble: 2 × 2
```

Question 8: What are the unique elements in the variable *artist* in the data object *data\_tab* (Display the results using R code) (1pt)? How many are there (use an R command to count the number of elements) and complete the sentence below using inline R code (1pt). Hint: type ?unique or ?length into the R console if unsure what to do.

Question 9: Using the *data\_tab* data frame, calculate the average

frequence (recorded in lfreq) for each of the rock artists in the

data set. Store the results in a new variable called avg. (2pts).

```
piece of music in the music dataset? To answer this question, use the tab_music data frame, and create a new variable called range and store the new data frame under the data object tab_music_range. Display the first four rows of the resulting data frame (3pts) Hint: To calculate the frequency range, take each genre specific value of lfreq and subtract from it the minimum value of lfreq for that genre, and then divide this answer by the maximum value of lfreq for that genre.
```

Question 10: What is the within genre frequency range for each

## 1 Dancing Queen Abba Rock 17600756. -90.0 29921 106. 59.6 0.0463
## 2 Knowing Me Abba Rock 9543021. -75.8 27626 103. 58.5 0.0435
## 3 Take a Chance Abba Rock 9049482. -98.1 26372 102. 125. 0.212
## 4 Mamma Mia Abba Rock 7557437. -90.5 28898 102. 48.8 0.0188

Write.csv(tab\_music\_range, 'data/tab\_music\_range.csv')

Question 11: Use the data object tab\_music, order the observations from largest to smallest according to the range variable and display the results in a table (1pt). Which genre type has the highest average frequence range (1pt)? Which genre type has the least variable frequency range (1pt)?

lvar lave lmax lfener lfreq range

<dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <

```
## # Groups: type [3]
                                          lvar lave lmax lfener lfreq range
     <chr>
                          <chr> <chr> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <</pre>
## 1 V7
                            Vival... Clas... 3.64e6 9.84 21450 101. 878. 0.933
## 2 V2
                         Vival... Clas... 7.71e5 21.7 6936 104. 844. 0.895
## 3 Girl
                            Eels Rock 8.85e7 0.336 32744 112. 392. 0.894
## 4 Rock Hard Times
                            Eels Rock 5.47e7 1.98 32759 109. 312. 0.691
                            Eels Rock 1.63e7 -0.141 30106 104. 312. 0.690
                             Enya New ... 1.14e6 -10.6 9994 102. 155. 0.684
## 6 The Memory of Trees
## 7 I Have A Dream
                                  Rock 3.37e6 -71.7 14699 105. 305. 0.672
## 8 I Want to Hold Your Hand Beatl... Rock 6.13e7 -6.03 28502 112. 295. 0.646
## 9 The Winner
                            Abba Rock 1.14e6 -67.8 8928 104. 278. 0.602
## 10 B4
                            Beeth... Clas... 2.35e7 -0.941 32766 106. 529. 0.536
## # ... with 52 more rows
freqs avg<-or_tab_music %>% group_by(type) %>% summarise(mean=mean(range)) %>% arrange(-mean)
freqs_var<-or_tab_music %>% group_by(type) %>% summarise(sd=sd(range)) %>% arrange(-sd)
write.csv(freqs_avg,'data/freqs_avg.csv')
write.csv(freqs_var, 'data/freqs_var.csv')
```

Question 12: Using the tab\_music\_range data object, display in a

boxplot the frequency range of the music library by genre (2pts).

tab\_music\_range %>% ggplot(aes(y=range))+geom\_boxplot() + ylab('Range') + ggtitle('Distribution of range') + the

Then, using boxplots display the frequencies (lfreq) by genre

type (1pt). State which genre type has the highest and lowest

The genre type with the highest average frequency is **Classical** and with most variable is **New wave** 

**Distribution of range** 

0.0

Distribution of Range for each Genre

Distribution of Lfreq for each Genre

original frequency range? (2pts).

-0.2

xt(hjust = 0.5, size=15, face='bold'))

ext(hjust = 0.5, size=15, face='bold'))

"Ifreq" to values within 0 and 1 irrespective of the genre.

lfreq).(2pts).

ace='bold'))

me(plot.title = element\_text(hjust = 0.5, size=15, face='bold'))

0.75 -0.50 -0.25 -

0.2

tab\_music\_range %>% ggplot(aes(x=type,y=range,fill=type)) + geom\_boxplot() + xlab('Genre') + ylab('Range') + labs
(fill='Genre') + ggtitle("Distribution of Range for each Genre") + theme classic() + theme(plot.title = element te

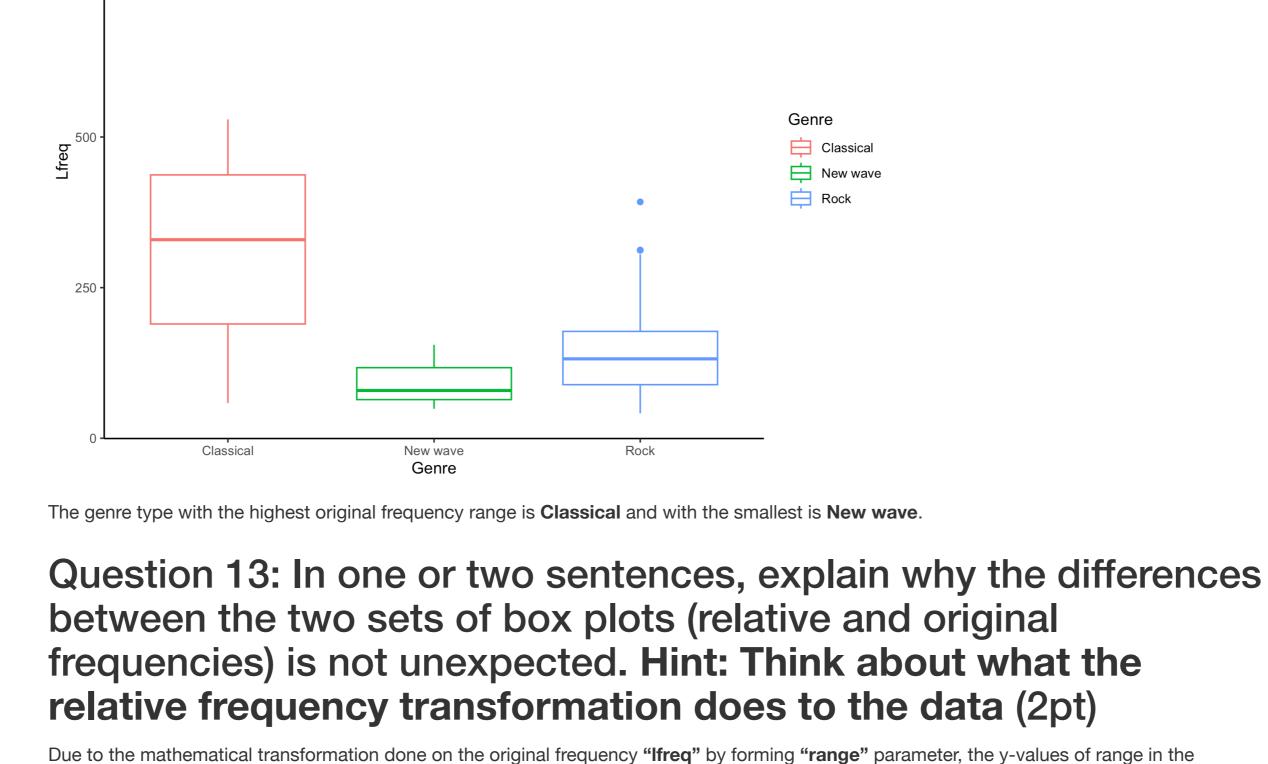
0.50

0.25

Classical
New wave
Genre

Rock

tab\_music\_range %>% ggplot(aes(x=type,y=lfreq,color=type))+geom\_boxplot() + xlab('Genre') + ylab('Lfreq') + labs(
color='Genre') + ggtitle("Distribution of Lfreq for each Genre") + theme\_classic() + theme(plot.title = element\_t



boxplots have now compressed to values between 0 and 1. There is no longer a notable difference between the relative frequency values for

Classical music when compared to New Wave and Rock music. This is due to the reason that the original frequency (Ifreq) values have now been

relationship between the orginal frequency range (lfreq) and the

tab\_music\_range %>% ggplot(aes(y=lfreq,x=lmax)) + geom\_point(size=2.5,color='red') + ggtitle("Original frequency (lfreq) Vs. variable frequency (lmax)") + theme\_classic() + theme(plot.title = element\_text(hjust = 0.5,size=15,f

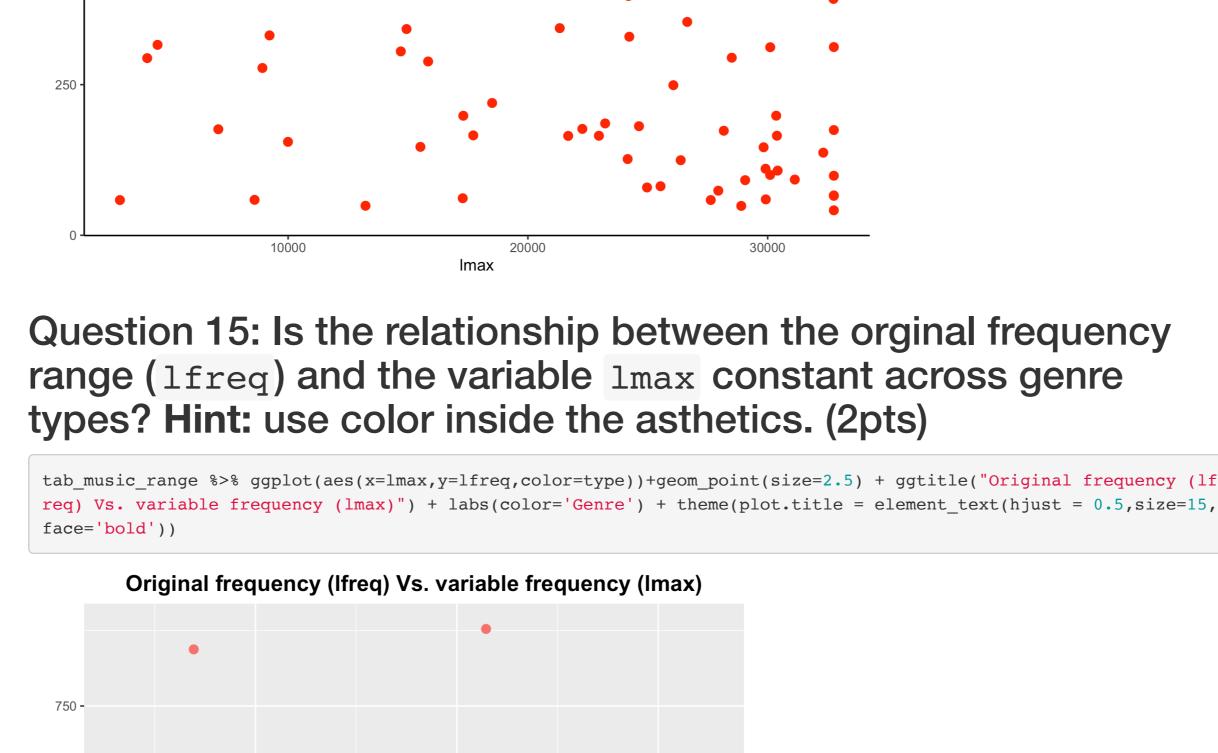
scaled appropriately according to the maximum and minimum frequency values for each individual genre, thereby compressing every value of

Question 14: Using the tab\_music\_range data object, plot the

variable lmax (on the x-axis display lmax and on the y-axis

Original frequency (Ifreq) Vs. variable frequency (Imax)

750-



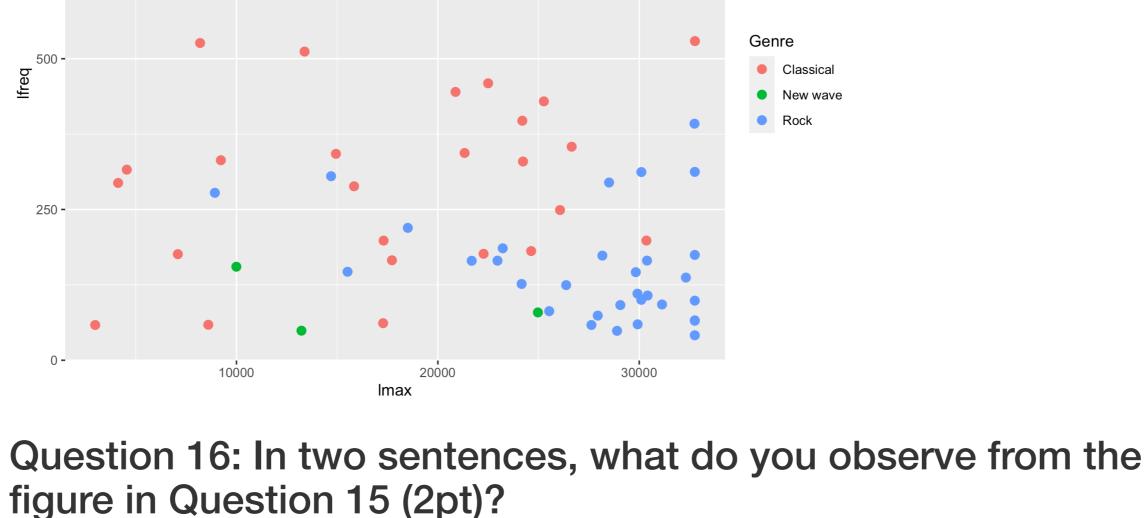
Genre

Classical

There appears to be a stronger relation between the **Ifreq** and **Imax** values of **Rock** music as majority of the scatter points are observed to be

However, the same cannot be inferred for **Classical** music as there seems to be **no clear relation** between the Ifreq and Imax values. The

concentrated in the bottom right section of the scatter plot having Imax values greater than 20000 and Ifreq values less than 250.



number of observations for New Wave music are too low for the relation to be interpreted.