# Intro to Data Visualization and Statistics in R

## Session #5

Genome Institute of Singapore

31st October 2019

# Learning objectives for Session 5

**Objectives:**

Combine information from multiple sheets/files.

Learn the following concepts

1. Match using `%in%` and `which()`
2. Replacing strings using `gsub()`
3. List
4. The for() loop

# Match using %in% and which()

```r
all_fruits <- c("Apple", "Banana", "Cherry",
                "Orange", "Rambutan")

red_fruits <- c("Cherry", "Rambutan", "Apple")

all_fruits %in% red_fruits       # Returns a logical vector
## [1]  TRUE FALSE  TRUE FALSE  TRUE

which(all_fruits %in% red_fruits)   # Position of matches
## [1] 1 3 5
```

# Replacing strings using gsub()

You just realized that there green and red apples. Let's change the text to red apple to be more precise.

```
gsub("Apple", "Red Apple", all_fruits)
# [1] "Apple" "Red Banana" "Cherry" "Orange" "Rambutan"
```

gsub() is like a flexible find-and-replace command.

# List | Method 1 for list creation

```r
red_fruits    <- c("Cherry", "Rambutan", "Apple")
yellow_fruits <- c("Banana", "Mango")
green_fruits  <- c("Mango")

fruits <- list(red_fruits, yellow_fruits, green_fruits)
fruits
## [[1]]
## [1] "Cherry"   "Rambutan" "Apple"
##
## [[2]]
## [1] "Banana" "Mango"
##
## [[3]]
## [1] "Mango"
```

Note: this list does not have a name.

# List | Method 2 for list creation

```r
fruits <- list(NULL)
fruits[["red_fruits"]] <- c("Cherry", "Rambutan", "Apple")
fruits[["yellow_fruits"]] <- c("Banana", "Mango")
fruits[["green_fruits"]]  <- c("Mango")
fruits
## [[1]]
## NULL
##
## $red_fruits
## [1] "Cherry"    "Rambutan" "Apple"
##
## $yellow_fruits
## [1] "Banana" "Mango"
##
## $green_fruits
## [1] "Mango"
```

Note: This list has name but first element is empty. Solution later!

# The for() loop

```r
all_fruits
## [1] "Apple"    "Banana"   "Cherry"   "Orange"   "Rambut

for(fr in all_fruits){

  fr  <- tolower(fr)
  txt <- paste("I love", fr, "now!")
  print(txt)
  rm(txt)
}
## [1] "I love apple now!"
## [1] "I love banana now!"
## [1] "I love cherry now!"
## [1] "I love orange now!"
## [1] "I love rambutan now!"
```

"fr" is the variable name that gets incremented to next value
when it reaches the end of the loop.
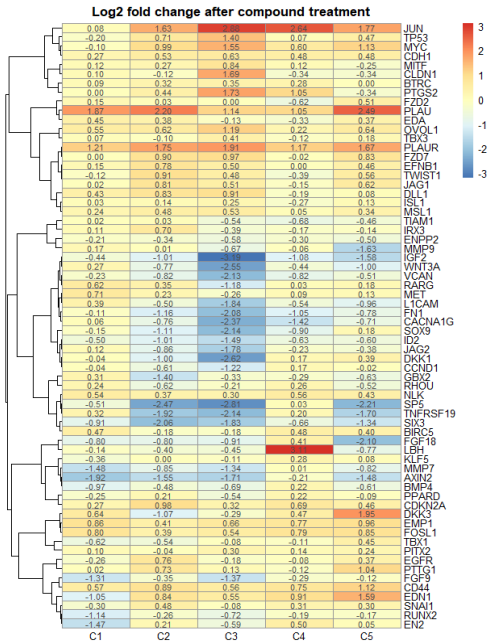
Exercise

# Background

RNAseq_treatment_effect.xlsx contains the outputs from an experiment with the following groups:

- ▶ untreated cell lines (as controls)
- ▶ samples after treatment with C1 compound
- ▶ samples after treatment with C2 compound
- ▶ samples after treatment with C3 compound
- ▶ samples after treatment with C4 compound
- ▶ samples after treatment with C5 compound

Three biological replicates per group. 18 samples in total.

Your friendly neighbourhood bioinformation has compared each of the treated groups vs control. You also have a set of 66 genes that you like to investigate in more depth.

# Objective | Generate heatmap of 66 genes



Log2 fold change after compound treatment

# Step 1 | Genes of interest

1. Setup (setwd, packages, clear workspace)
2. Read in the the genes of interest as a vector called "keep".
   Hint: pull()
3. Do you have 66 genes in your list?

# Step 1 | Genes of interest

1. Setup (setwd, packages, clear workspace)
2. Read in the the genes of interest as a vector called "`keep`".
   Hint: pull()
3. Do you have 66 genes in your list?

```r
setwd("C:/Users/aramasamy/Desktop/R_workshop")
pacman::p_load(tidyverse, readxl, pheatmap)
rm(list=ls())

fn <- "data/RNAseq_treatment_effect.xlsx"

## Identify the genes that we want to keep for plotting
keep <- read_excel(fn, sheet="My genes", col_names=F) %>% p
length(keep)
## [1] 66
```

# Step 2A | Template processing one file

1. Read in the data from first compound.
2. Keep only the selected 66 genes.
3. Keep only the Gene and LFC columns.
4. Rename the "LFC" column as "C1"

# Step 2A | Template processing one file

1. Read in the data from first compound.
2. Keep only the selected 66 genes.
3. Keep only the Gene and LFC columns.
4. Rename the "LFC" column as "C1"

```
df <- read_excel(fn, sheet="C1") %>%
    filter(Gene %in% keep) %>%
    select(Gene, LFC)

colnames(df) <- gsub("LFC", "C1", colnames(df))
head(df, 3)
## # A tibble: 3 x 2
##    Gene       C1
##    <chr> <dbl>
## 1 AXIN2 -1.92
## 2 PLAU    1.87
## 3 RUNX2 -1.14
```

# Step 2B | Generalize the template

1. Create a variable name `sh` and assign the value of "C1" to it.
2. Change all "C1" strings to the variable `sh`.

# Step 2B | Generalize the template

1. Create a variable name `sh` and assign the value of "C1" to it.
2. Change all "C1" strings to the variable `sh`.

```
sh <- "C1"

df <- read_excel(fn, sheet=sh) %>%
    filter(Gene %in% keep) %>%
    select(Gene, LFC)

colnames(df) <- gsub("LFC", sh, colnames(df))
head(df, 3)
## # A tibble: 3 x 2
##    Gene       C1
##    <chr> <dbl>
## 1 AXIN2 -1.92
## 2 PLAU   1.87
## 3 RUNX2 -1.14
```

# Step 2C | Generate vector to loop over

1. Read in the sheet names of the Excel file.
2. Remove any unwanted sheet names from the vector.

# Step 2C | Generate vector to loop over

1. Read in the sheet names of the Excel file.
2. Remove any unwanted sheet names from the vector.

```
sheets <- excel_sheets(fn)
sheets
## [1] "C1"         "C2"         "C3"         "C4"         "C5"
## [7] "My genes"

sheets <- setdiff( sheets, c("README", "My genes") )
sheets
## [1] "C1" "C2" "C3" "C4" "C5"
```

# Step 2D | Repeat with a `for()` loop

1. Generate an null list called "`store`"
2. Use the `sh` in a `for()` loop to read in multiple sheets
3. Save the data that you read in and processed into the "`store`" list

# Step 2D | Repeat with a `for()` loop

1. Generate an null list called "`store`"
2. Use the `sh` in a `for()` loop to read in multiple sheets
3. Save the data that you read in and processed into the "`store`" list

```r
store <- list(NULL)  ## A list to store each sheet

for(sh in sheets){

  df <- read_excel(fn, sheet=sh) %>%
    filter(Gene %in% keep) %>%
    select(Gene, LFC)

  colnames(df) <- gsub("LFC", sh, colnames(df))

  store[[ sh ]] <- df
  rm(df)
}
```

# Step 2E | Remove the ghost element

It is an unintended feature in R to create an empty first element if you assign data into a list using names.

Here is one way to remove it:

```
names(store)
## [1] ""   "C1" "C2" "C3" "C4" "C5"

store <- store[ sheets ]
names(store)
## [1] "C1" "C2" "C3" "C4" "C5"
```
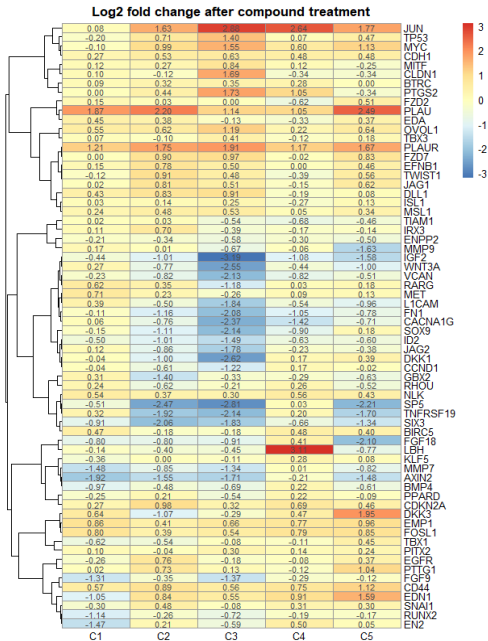
# Step 3 | Combine the data

1. Combine it into one big dataframe.
2. Set the "Gene" column as the rownames.

# Step 3 | Combine the data

1. Combine it into one big dataframe.
2. Set the "Gene" column as the rownames.

```
comb <- plyr::join_all(store, type="full") %>%
  column_to_rownames("Gene")
```

# Step 4 | Generate this heatmap



Log2 fold change after compound treatment

# Step 4 | Generate this heatmap

```r
my_title <- "Log2 fold change after compound treatment"

pheatmap(comb,
         scale         = "none",
         angle_col     = 0,
         cluster_cols = FALSE,
         clustering_distance_rows = "correlation",
         main          = my_title,
         display_numbers = TRUE)
```