

# Gene\_expression\_analysis

Fay

2022-05-18

GAPDH HKG

B-actin HKG

Ppia HKG

Ppip HKG

CDC42 HKG susceptible to DNA contamination

Housekeeping genes selected: GAPDH and PPIB

Relm-b mucosal defense factor (goblet cells)

Muc2 the major secretory mucin within the gastrointestinal tract

TFF3 mucosal defense factor (goblet cells)

Muc5ac similar to MUC2, produced by surface goblet cells

NKp46 NK marker

F4/80 macrophage marker (distinguish by immune response trend)

Mpo myeloperoxidase in Neutrophils

MyD88 TLR protein, NF-kB IRAK protein, inflammation marker by TLR MyD88-Dependent Pathway  
caspase-1 inflammasome marker (IL-1b and IL-18 production)

IL-1Ra natural IL-1b antagonist for infection control (if not increase in Tregs is seen)

CXCL9, immune cell migration marker + Th1 activator (confirm FACS)

CXCR3, CXCL9 and CXCL11 receptor

IL-6 TNF inhibitor,

IL-12ra T-cell marker Th1

IFN-γ compare with IFN-γ producing cells and IFN-γ ELISAs, should correlate with PRF1, NKp46 and F4/80.

One of these cell types just have to be doing the job!

IRG6A autonomous cell defense (opsonization)

TNF-α upregulated in eimeria but not well explained. Could be present and driving infection where IFN-γ isn't

IL-17 in case IFN-γ isn't coming up but pathogenicity is

TRIF Type I IFN production TRIF Dependent Pathway

Socs1 JAK/STAT signaling pathway, proinflammatory regulating + T-cell differentiation, could explain severity

IDO1 DC, monocyte and MC protein regulating T-cell activity

Prf1 perforin, should be dominant in primary infections, but must be correlated between T-cell and NK cell expressions

CD56 CD56bright = more cytokine producing NKs, CD56dim = more direct cytotoxic killing

IL-4

IL-13

IL-10

## 1. Gene expression in the laboratory infections - Heatmap

#document gene normalization <https://www.youtube.com/watch?v=tgp4bbnj-ng> paper: Licak Analysis of relative gene expression data using real time quantitative Add the normalized Gene expression to the rest of the data

<https://www.youtube.com/watch?v=tdp4bbnj-ng>

If you don't need a graph with absolute values of expressions and just relative comparison then you can choose one of the samples as a calibrator (control sample). This control sample will have the value of 1 and the rest of the samples will be compared relative to this one. I would choose as control the sample with the lowest expression.

<https://www.researchgate.net/post/How-can-I-analyze-qPCR-data-without-control-sample>

```
#create calibrator samples for the analysis
# create new columns with the minimum of each gene
# which I will use as the calibrator sample
#create a function which gives you the PPIB value for the minimum value of a gene
Dct_Calibrator <- function(x) {
  h <- gene[which(gene[[x]] == min(gene[[x]], na.rm = TRUE)), ][,c(x, "PPIB")]
  h[[x]] - h[["PPIB"]]
}
```

```
#Apply the function to all genes
```

```
Calibrator_values <- as.data.frame(unlist(lapply(Genes, Dct_Calibrator)))
```

```
Calibrator_values <- cbind(Calibrator_values, as.data.frame(Genes))
```

```
Calibrator_values <- Calibrator_values %>%
```

```
  rename(Dct.calibrator = "unlist(lapply(Genes, Dct_Calibrator))")
```

```
Calibrator_values <- Calibrator_values %>% pivot_wider(names_from = Genes, values_from = Dct.calibrator)
```

```
gene <- gene %>%
```

```
  dplyr::mutate(IFNy_N = 2^- ((IFNy - PPIB) - Calibrator_values[["IFNy"]]),
    CXCR3_bio_N = 2^- ((CXCR3_bio - PPIB) - Calibrator_values[["CXCR3_bio"]]),
    IL.6_N = 2^-((IL.6 - PPIB) - Calibrator_values[["IL.6"]]),
    IL.10_N = 2^-((IL.10 - PPIB) - Calibrator_values[["IL.10"]]),
    IL.13_N = 2^-((IL.13 - PPIB) - Calibrator_values[["IL.13"]]),
    IL1RN_N = 2^-((IL1RN - PPIB) - Calibrator_values[["IL1RN"]]),
    CASP1_N = 2^-((CASP1 - PPIB) - Calibrator_values[["CASP1"]]),
    CXCL9_N = 2^-((CXCL9 - PPIB) - Calibrator_values[["CXCL9"]]),
    IDO1_N = 2^-((IDO1 - PPIB) - Calibrator_values[["IDO1"]]),
    IRGM1_N = 2^-((IRGM1 - PPIB) - Calibrator_values[["IRGM1"]]),
    MPO_N = 2^-((MPO - PPIB) - Calibrator_values[["MPO"]]),
    MUC2_N = 2^-((MUC2 - PPIB) - Calibrator_values[["MUC2"]]),
    MUC5AC_N = 2^-((MUC5AC - PPIB) - Calibrator_values[["MUC5AC"]]),
    MYD88_N = 2^-((MYD88 - PPIB) - Calibrator_values[["MYD88"]]),
    NCR1_N = 2^-((NCR1 - PPIB) - Calibrator_values[["NCR1"]]),
    PRF1_N = 2^-((PRF1 - PPIB) - Calibrator_values[["PRF1"]]),
    RETNLB_N = 2^-((RETNLB - PPIB) - Calibrator_values[["RETNLB"]]),
    SOCS1_N = 2^-((SOCS1 - PPIB) - Calibrator_values[["SOCS1"]]),
    TICAM1_N = 2^-((TICAM1 - PPIB) - Calibrator_values[["TICAM1"]]),
    TNF_N = 2^-((TNF - PPIB) - Calibrator_values[["TNF"]]))
```

```

#create a new vector for selecting genes
genes_N <- c("IFNy_N", "CXCR3_bio_N", "IL.6_N", "IL.10_N", "IL.13_N", "IL1RN_N",
             "CASP1_N", "CXCL9_N", "IDO1_N", "IRGM1_N", "MPO_N", "MUC2_N", "MUC5AC_N", "MYD88_N",
             "NCR1_N", "PRF1_N", "RETNLB_N", "SOCS1_N", "TICAM1_N", "TNF_N")

# turn the data frame into a matrix and transpose it. We want to have each cell
# type as a row name
gene <- t(as.matrix(gene1 %>% dplyr::select(c(EH_ID, all_of(Genes)))))

#switch the matrix back to a data frame format
gene <- as.data.frame(gene)

# turn the first row into column names
gene %>%
  row_to_names(row_number = 1) -> heatmap_data

table(rowSums(is.na(heatmap_data)) == nrow(heatmap_data))

##
## FALSE TRUE
##    17    3

# turn the columns to numeric other wise the heatmap function will not work
heatmap_data[] <- lapply(heatmap_data, function(x) as.numeric(as.character(x)))

# remove columns with only NAs
heatmap_data <- Filter(function(x)!all(is.na(x)), heatmap_data)

#remove rows with only NAs
heatmap_data <- heatmap_data[, colSums(is.na(heatmap_data)) != nrow(heatmap_data)]

### Prepare the annotation data frame for the heatmap

annotation_df <- as_tibble(Challenge) %>%
  dplyr::filter(infection == "challenge", dpi == dpi_max) %>%
  dplyr::group_by(EH_ID) %>%
  dplyr::select(c("EH_ID", "Parasite_challenge", "infection_history", "mouse_strain",
                  "max_WL", "delta", "hybrid_status")) %>%
  dplyr::filter(EH_ID %in% colnames(heatmap_data))

annotation_df <- unique(annotation_df)

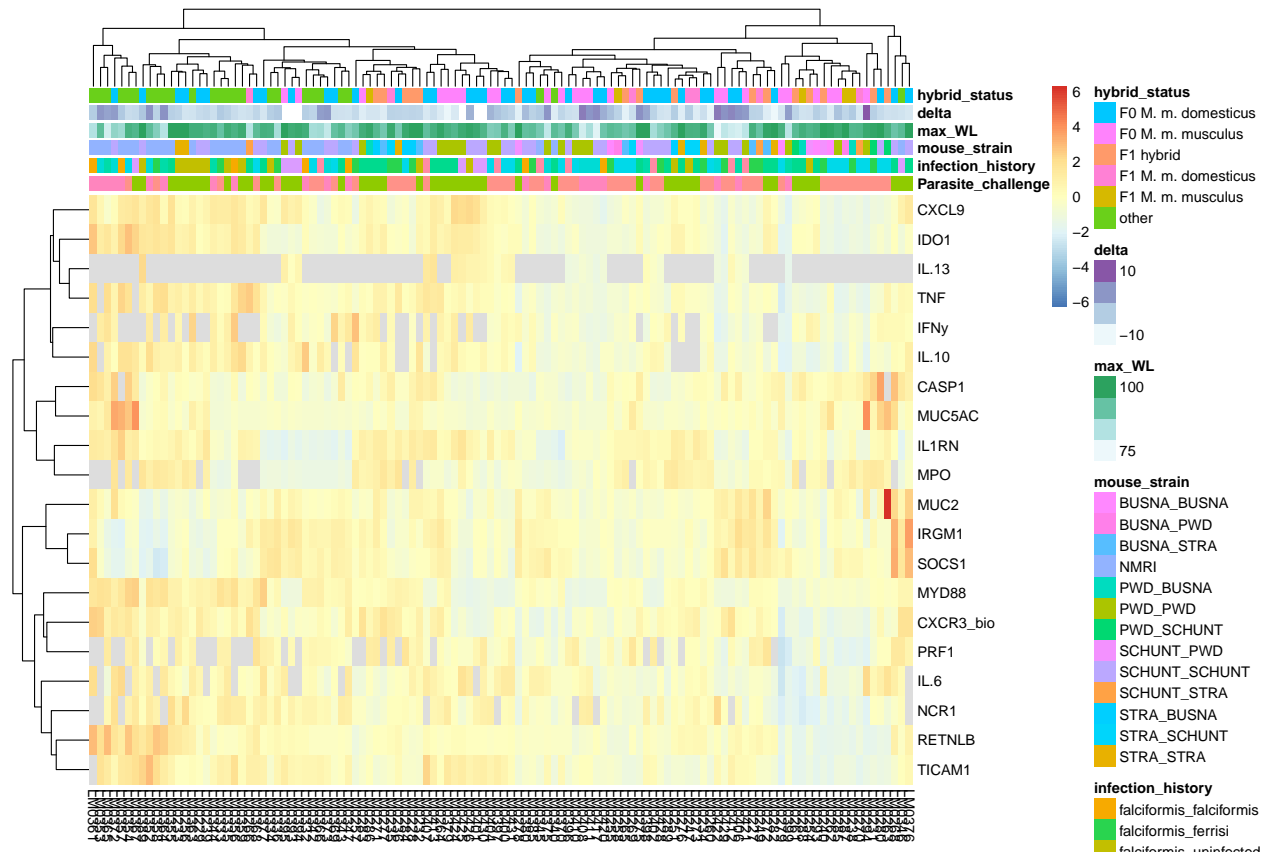
annotation_df <- as.data.frame(unique(annotation_df)) %>%
  dplyr::select(-EH_ID)

### Prepare the annotation columns for the heatmap
rownames(annotation_df) <- annotation_df$EH_ID

# Match the row names to the heatmap data frame
rownames(annotation_df) <- colnames(heatmap_data)

```

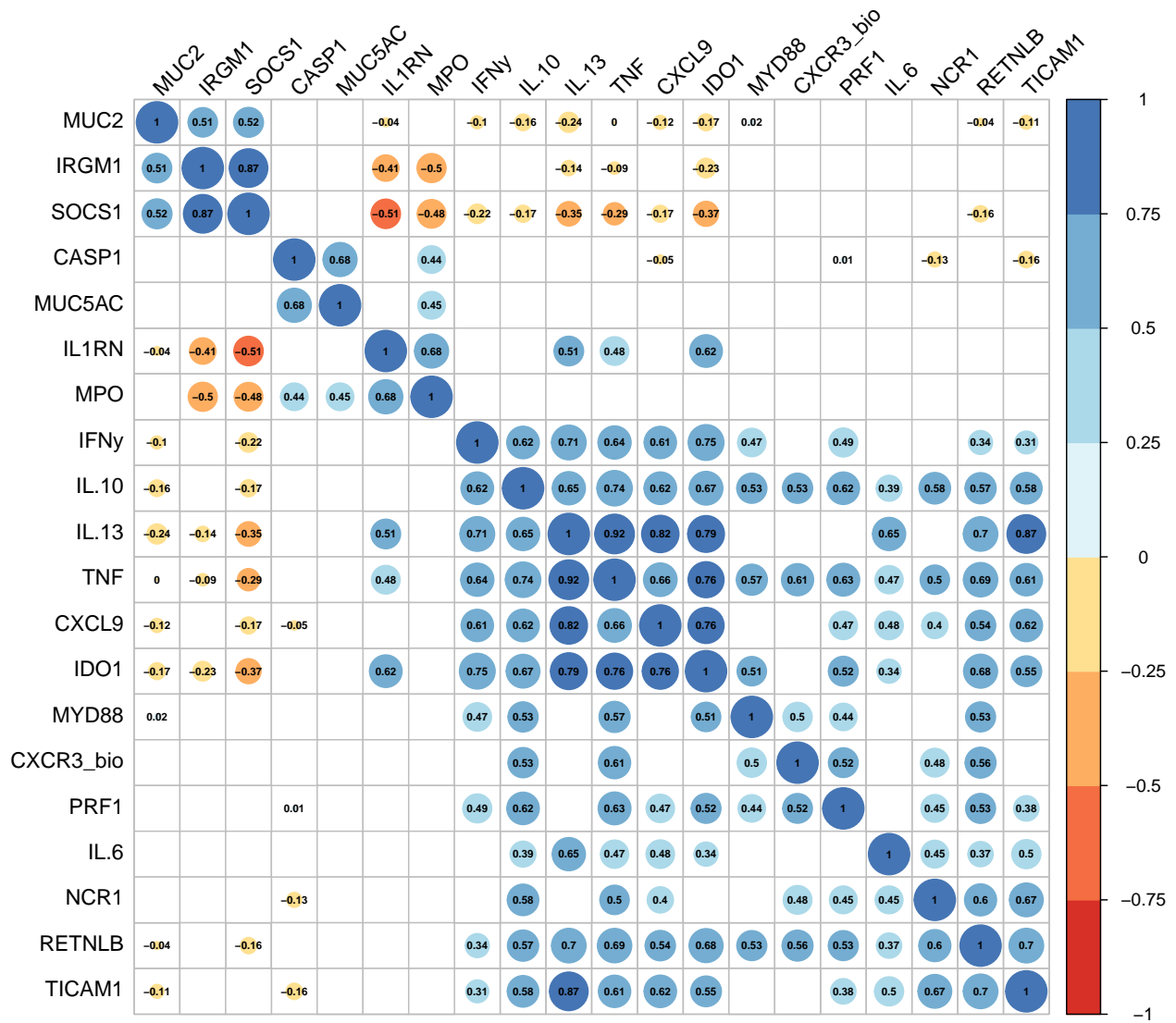
Heatmap on gene expression data:



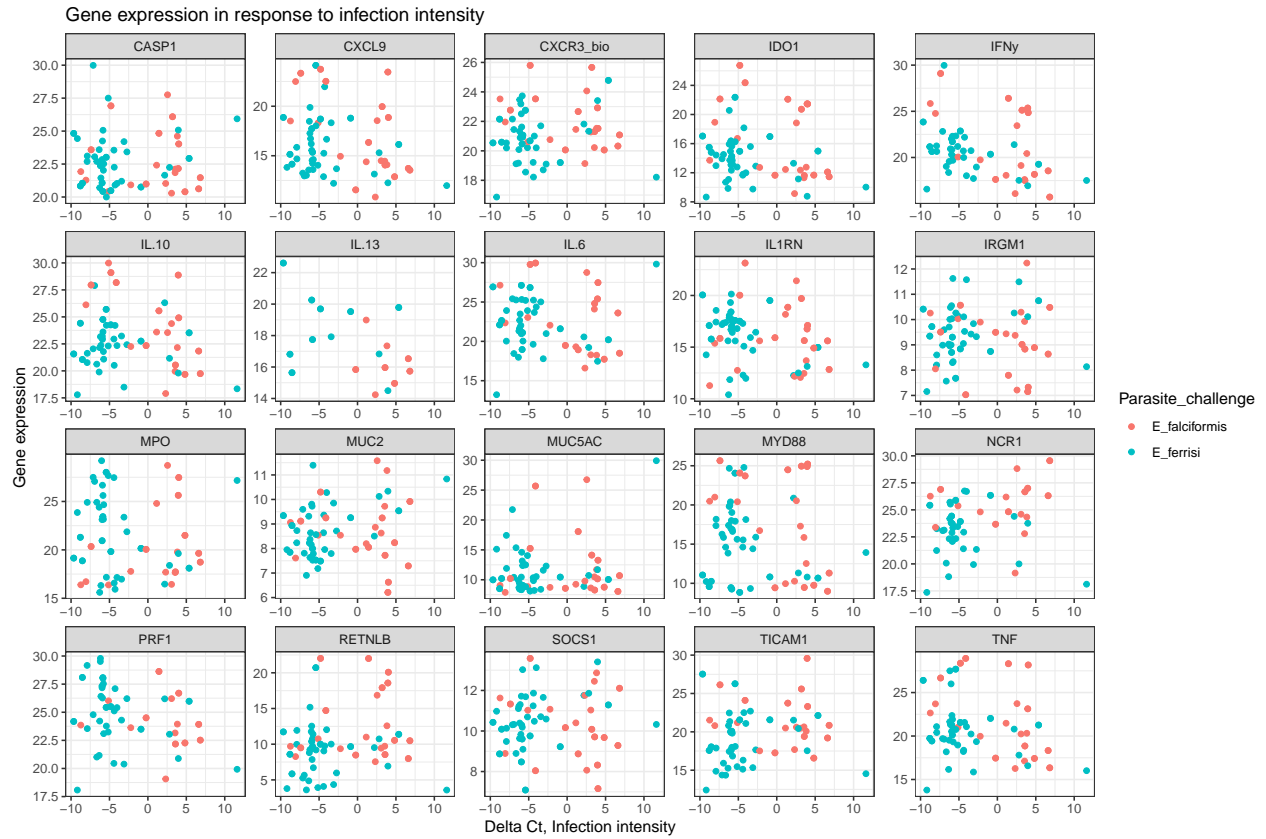
## 2. Correlations between the genes

### Corrplot of correlations

Here is a corrplot of the correlations between the genes. I am using the non-normalized genes

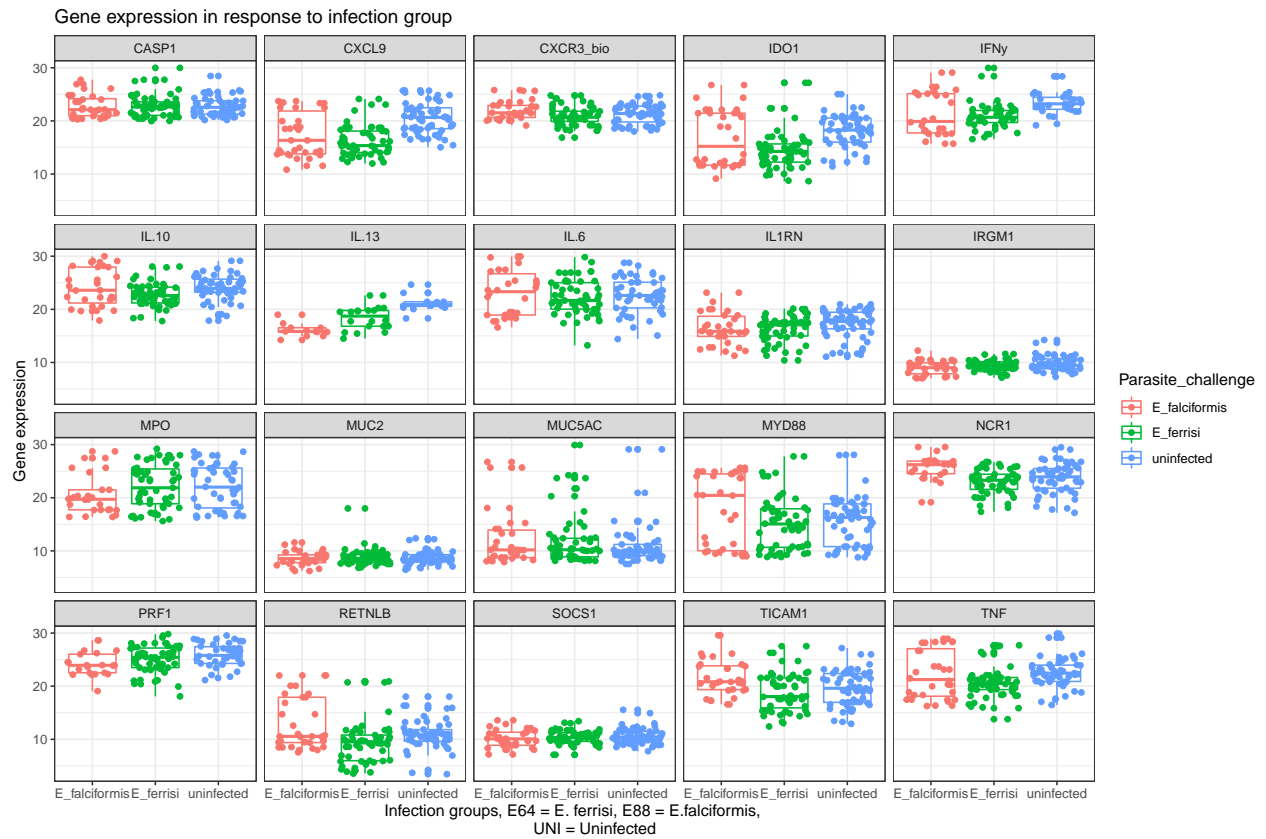


## Warning: Removed 121 rows containing missing values (geom\_point).



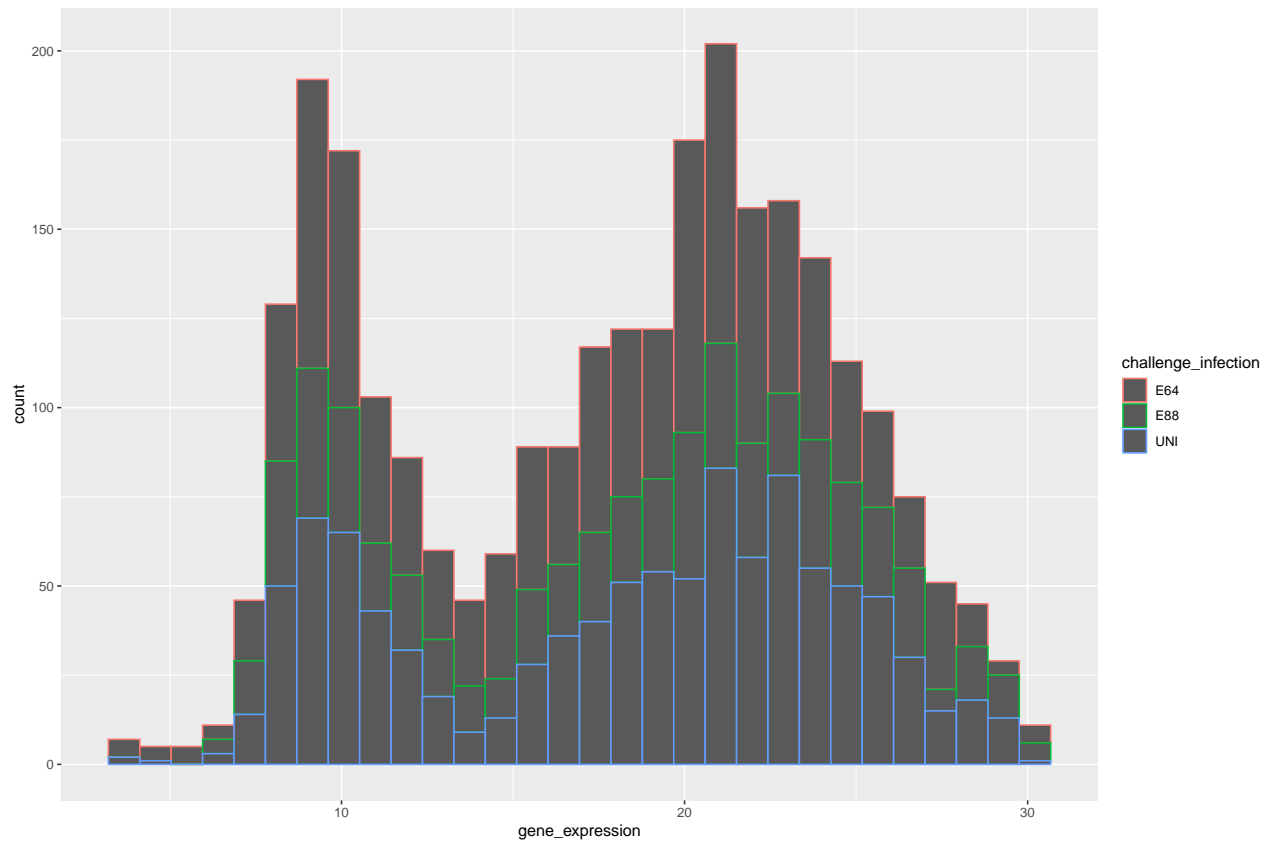
```
## Warning: Removed 244 rows containing non-finite values (stat_boxplot).
```

```
## Warning: Removed 244 rows containing missing values (geom_point).
```



## Warning: Ignoring unknown parameters: echo

## Warning: Removed 244 rows containing non-finite values (stat\_bin).



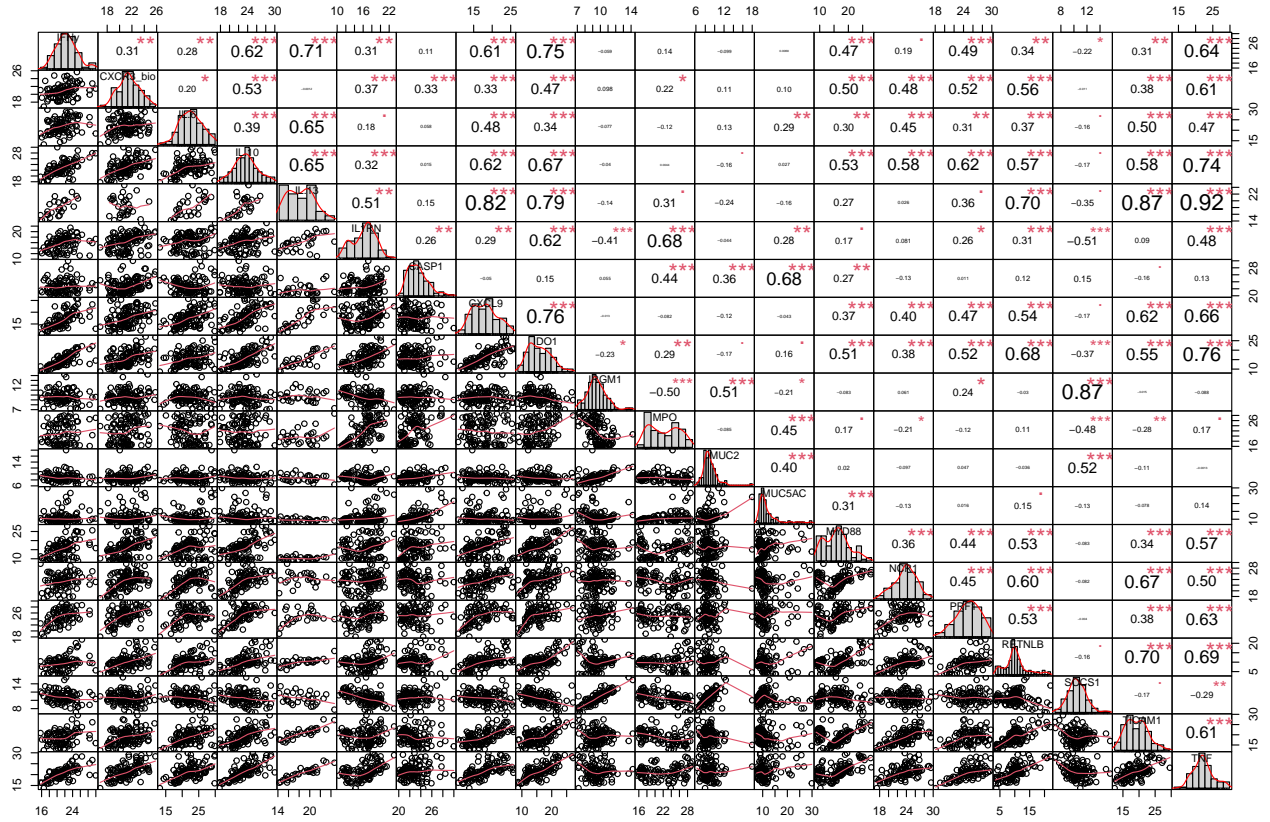
It is possible to compute a pca with missing data using the package missMDA. The missMDA package is dedicated to missing values in exploratory multivariate data analysis: single imputation/multiple imputation, etc.

Following the tutorial of the package author: Francois Husson: [https://www.youtube.com/watch?v=OOM8\\_FH6\\_8o](https://www.youtube.com/watch?v=OOM8_FH6_8o)

### 3. PCA

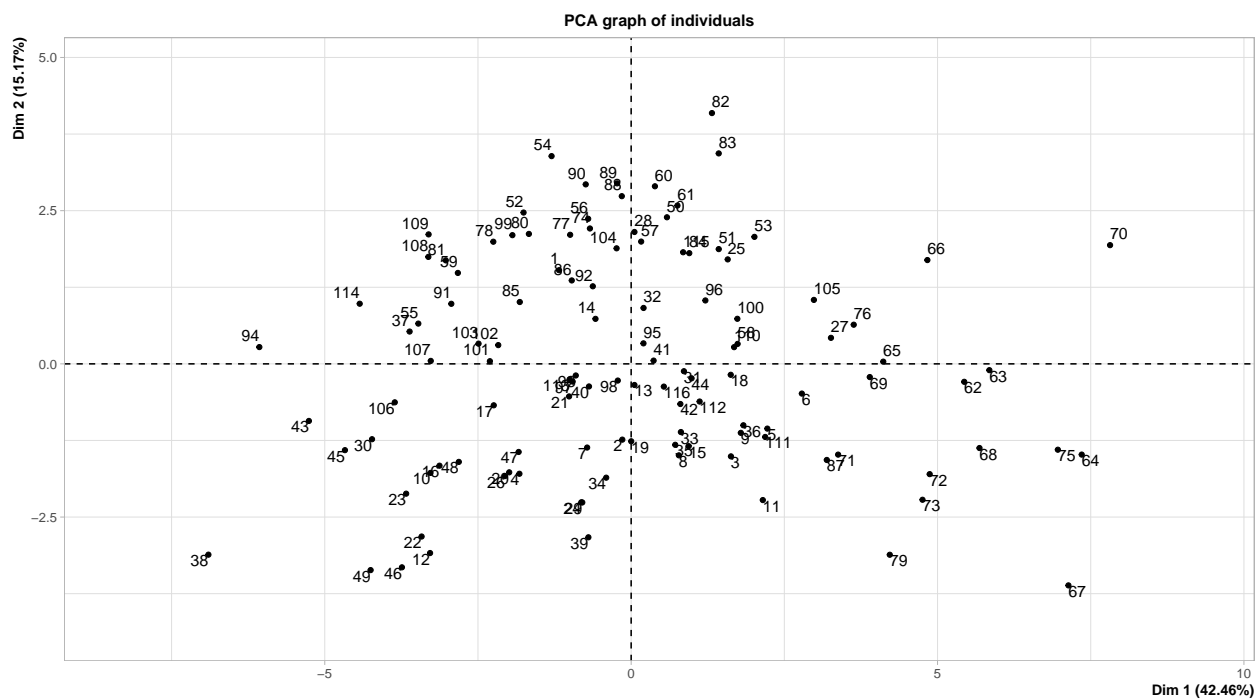
**Handling missing data in a pca:** Bad methods: removing individuals with missing data or replacing missing data with the mean (default setting in many packages).

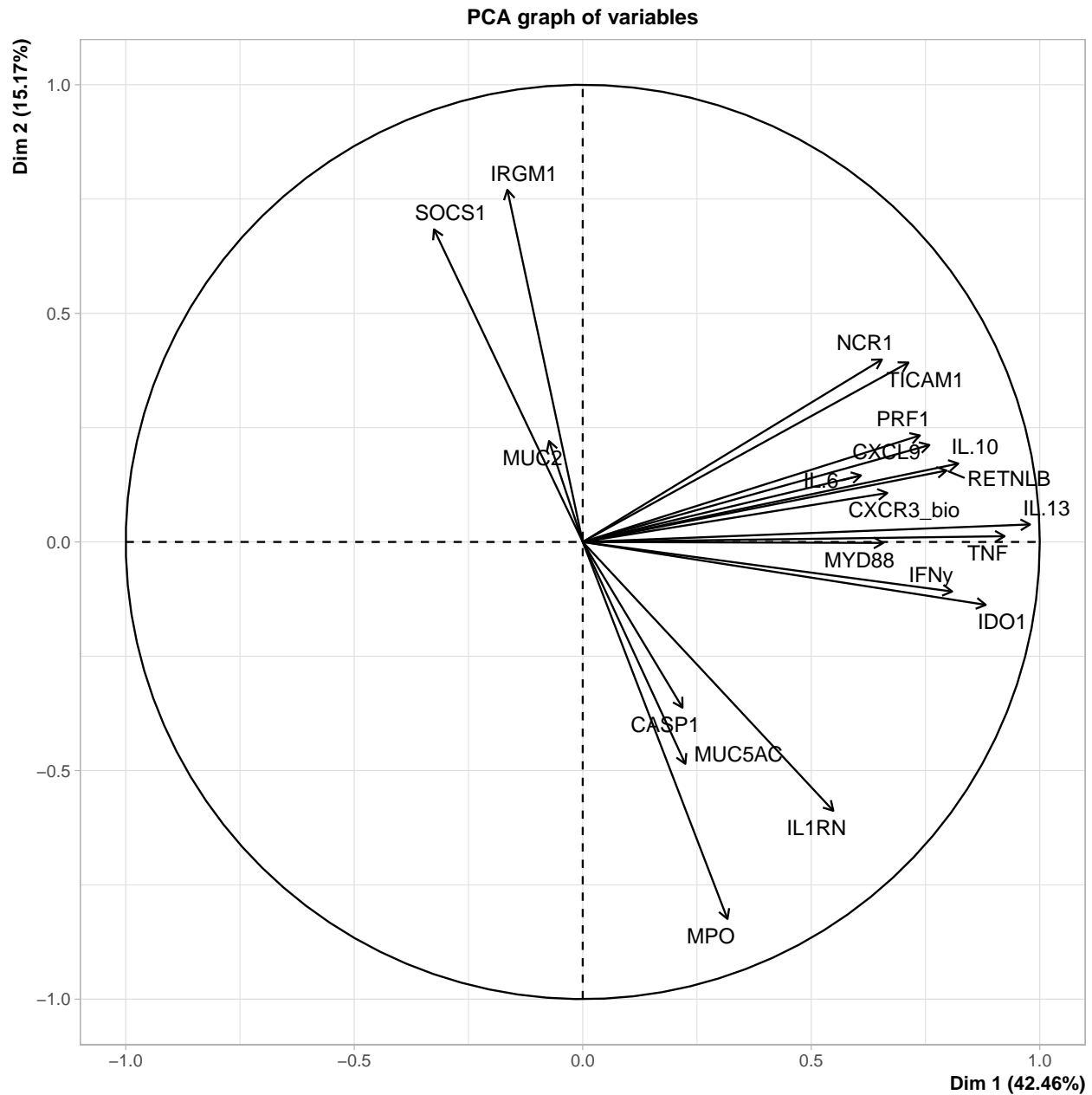




We will now continue by using an iterative pca to impute missing data A. Initialization: impute using the mean B. Step lampda: # a. do pca on imputed data table S dimensions retained # b. missing data imputed using pca # c. means (and standard deviations) updated C. Iterate the estimation and imputation steps (until convergence) (convergence: the act of converging and especially moving toward union or uniformity)

Overfitting is a common problem due to believing too much in links between variables. -> regularized iterative PCA (This version is what is being implmented in missMDA) This is a way of taking less risk when imputing the missing data. The algorithm estimates the missing data values with values that have no influence on the PCA results, i.e., no influence on the coordinates of the individuals or variables.





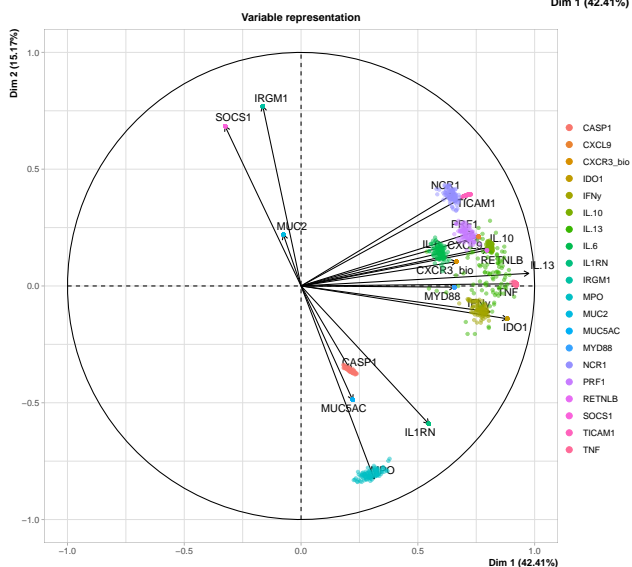
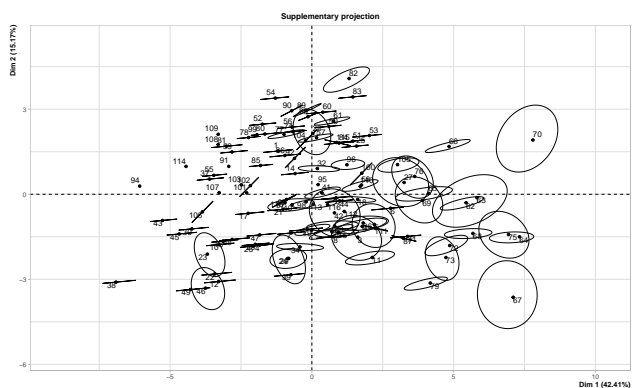
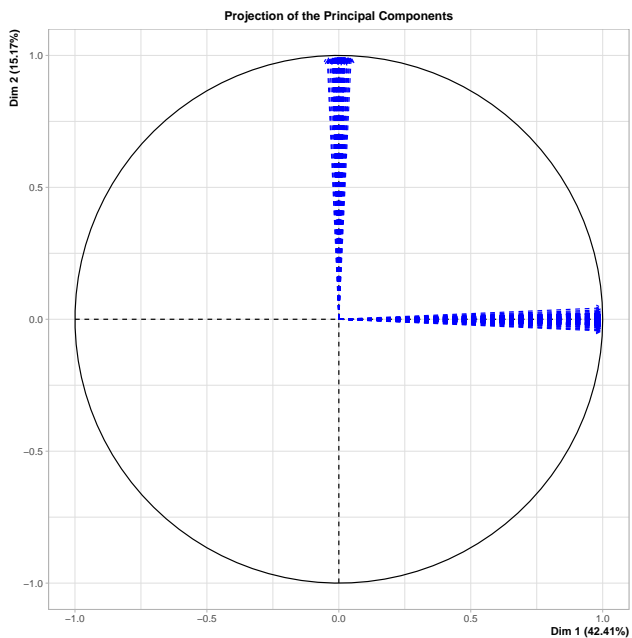
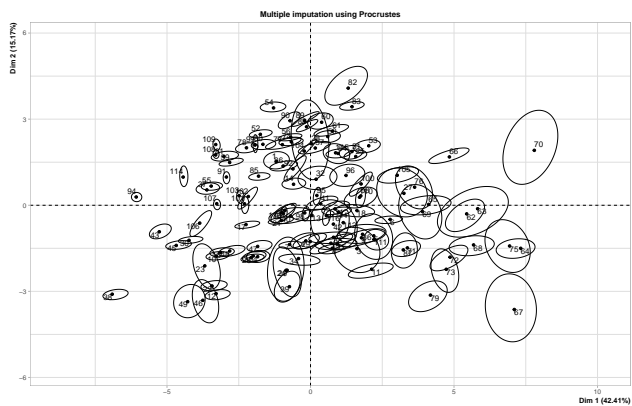
Caution: When imputing data, the percentages of inertia associated with the first dimensions will be overestimated.

Another problem: the imputed data are, when the pca is performed considered like real observations. But they are estimations!!

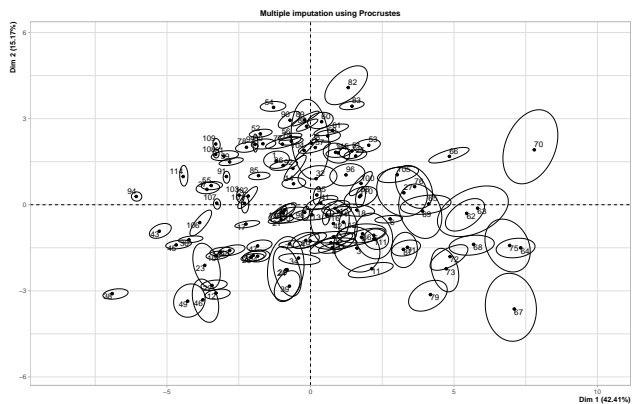
Visualizing uncertainty due to missing data:

→ multiple imputation: generate several plausible values for each missing data point

We here visualize the variability, that is uncertainty on the plane defined by two pca axes.

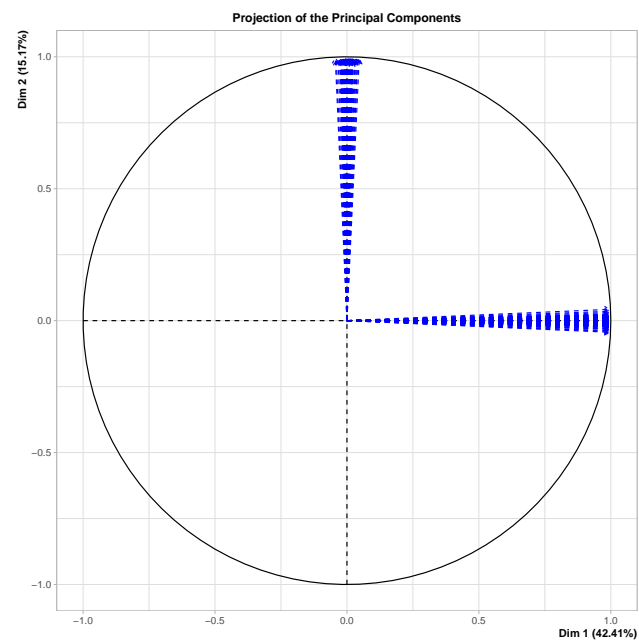


## \$PlotIndProc

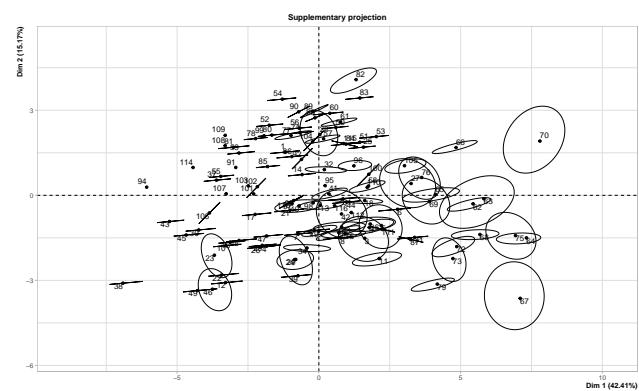


##

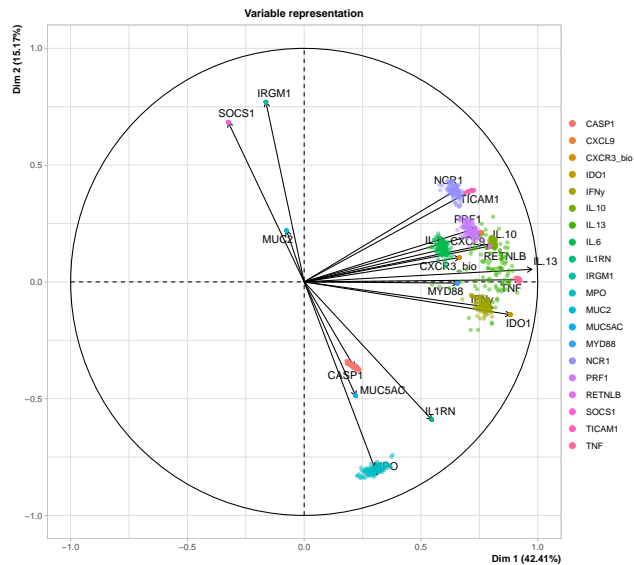
```
## $PlotDim
```



```
##  
## $PlotIndSupp
```



```
##  
## $PlotVar
```



Individuals lying on the axis have no missing data, but individuals that far away have many missing data. big ellipse = big uncertainty tight ellipse (line) = low uncertainty

Variable representation: Points tight together (look like one) - have no missing variables → low uncertainty  
Points spread → higher variability → higher uncertainty

High uncertainty → we should interpret the result with care

The individuals with many missing data values make the axes move, and thus the positions of all individuals

Therefore in the last plots every individual is getting an eclipse as they are as well influenced by the missing data of the others.

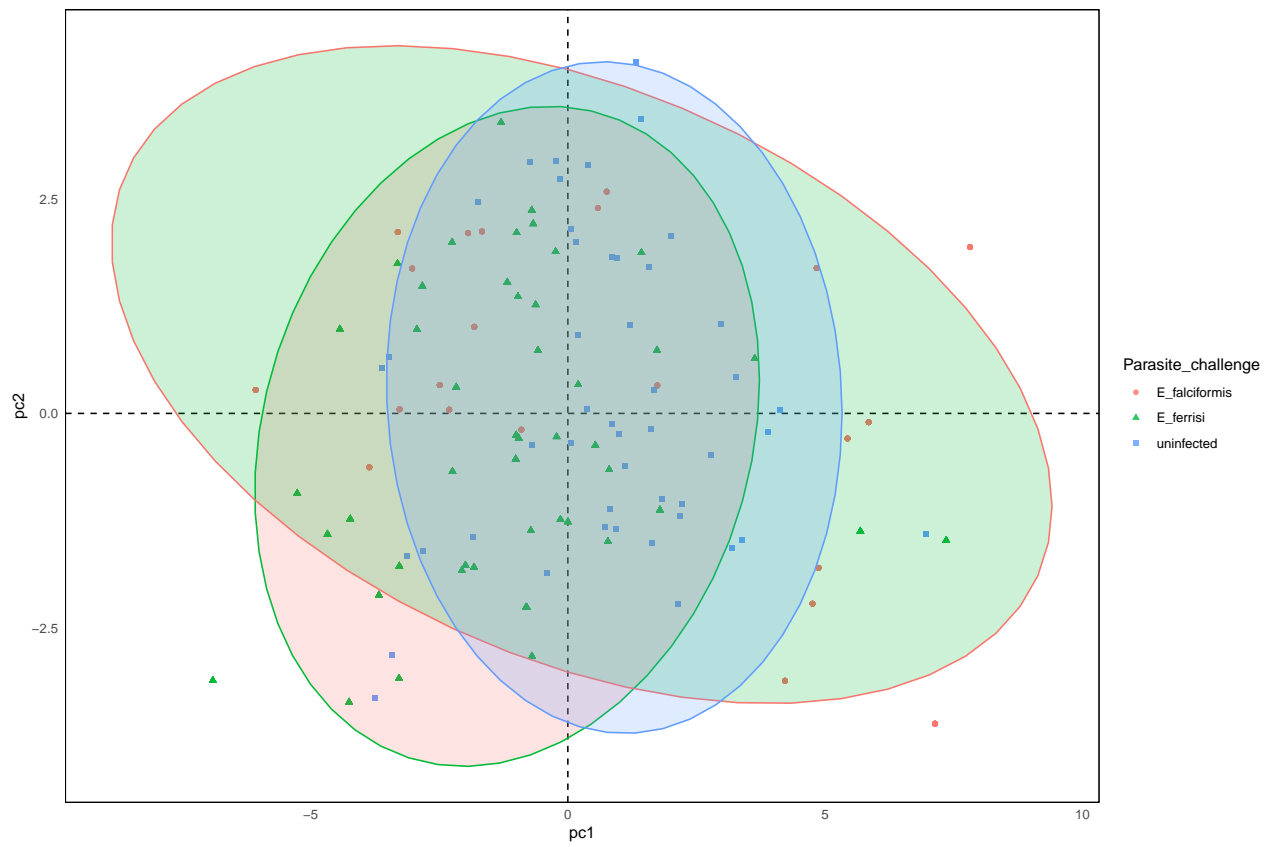
The plot with the dimensions shows the projections of the pca dimensions of each imputed table on the pca plane obtained using the original imputed data table

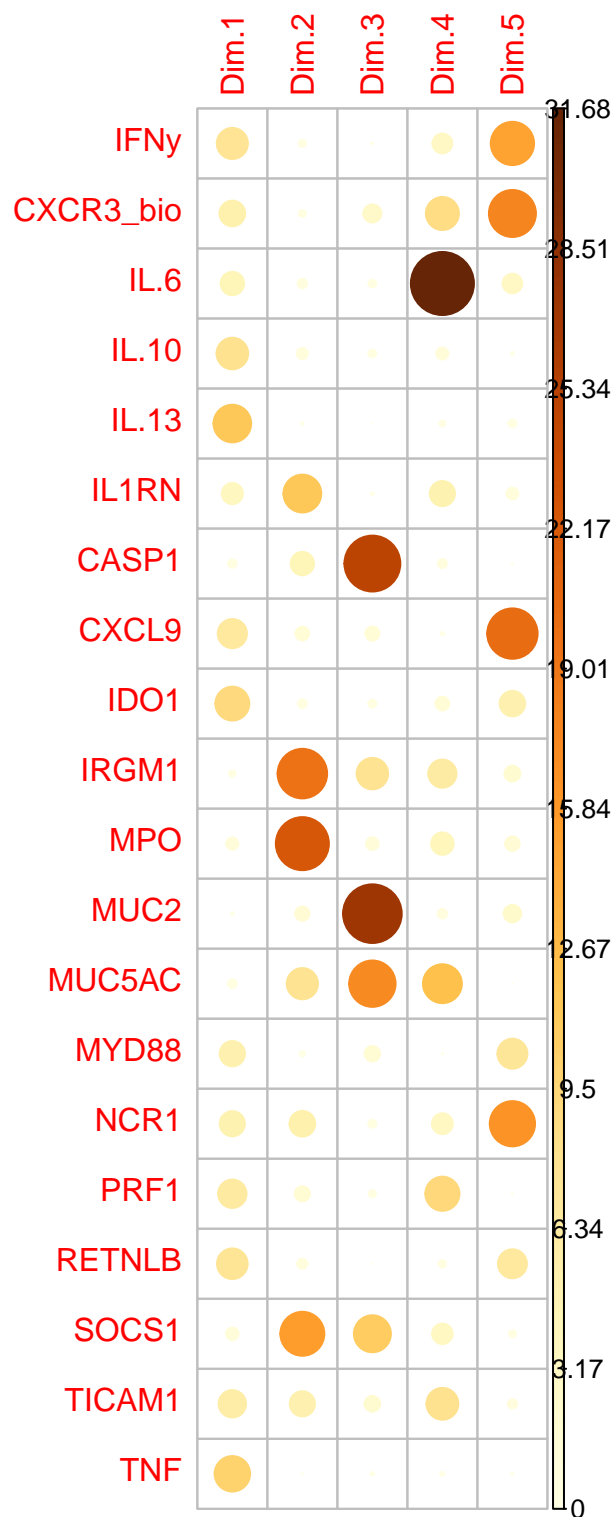
As all of the arrows are close to either the first or second axes, this means that the axes are stable with respect to the set of imputed tables → we don't have evidence of instability here.

Biplot of the imputed gene pca

*#Now we can make our initial plot of the PCA.*

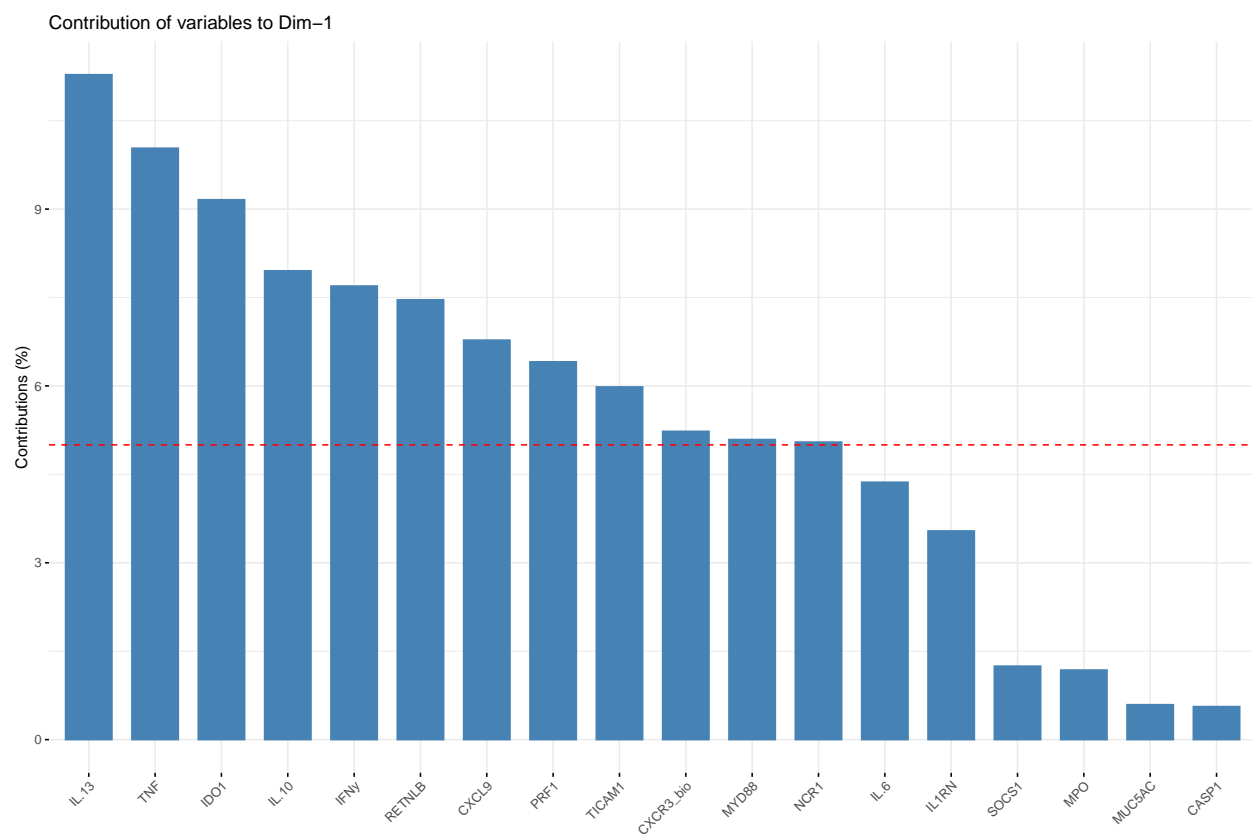
```
imputed_expr %>%
  pivot_longer(cols = all_of(Genes), names_to = "Gene", values_to = "gene_expression") %>%
  ggplot(aes(x = pc1, y = pc2, color = Parasite_challenge, shape = Parasite_challenge)) +
  geom_hline(yintercept = 0, lty = 2) +
  geom_vline(xintercept = 0, lty = 2) +
  geom_point(alpha = 0.8) +
  stat_ellipse(geom="polygon", aes(fill = challenge_infection), alpha = 0.2, show.legend = FALSE,
              level = 0.95) +
  theme_minimal() +
  theme(panel.grid = element_blank(), panel.border = element_rect(fill="transparent"))
```



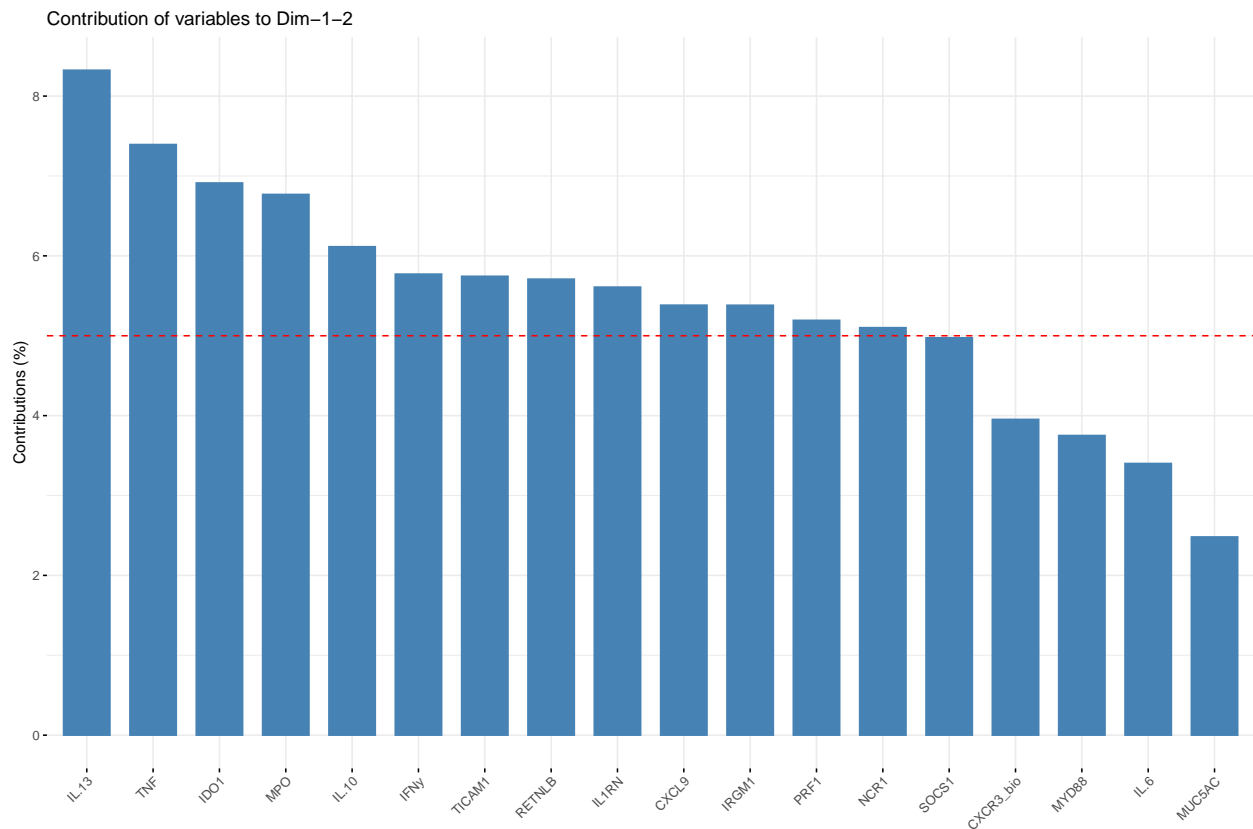
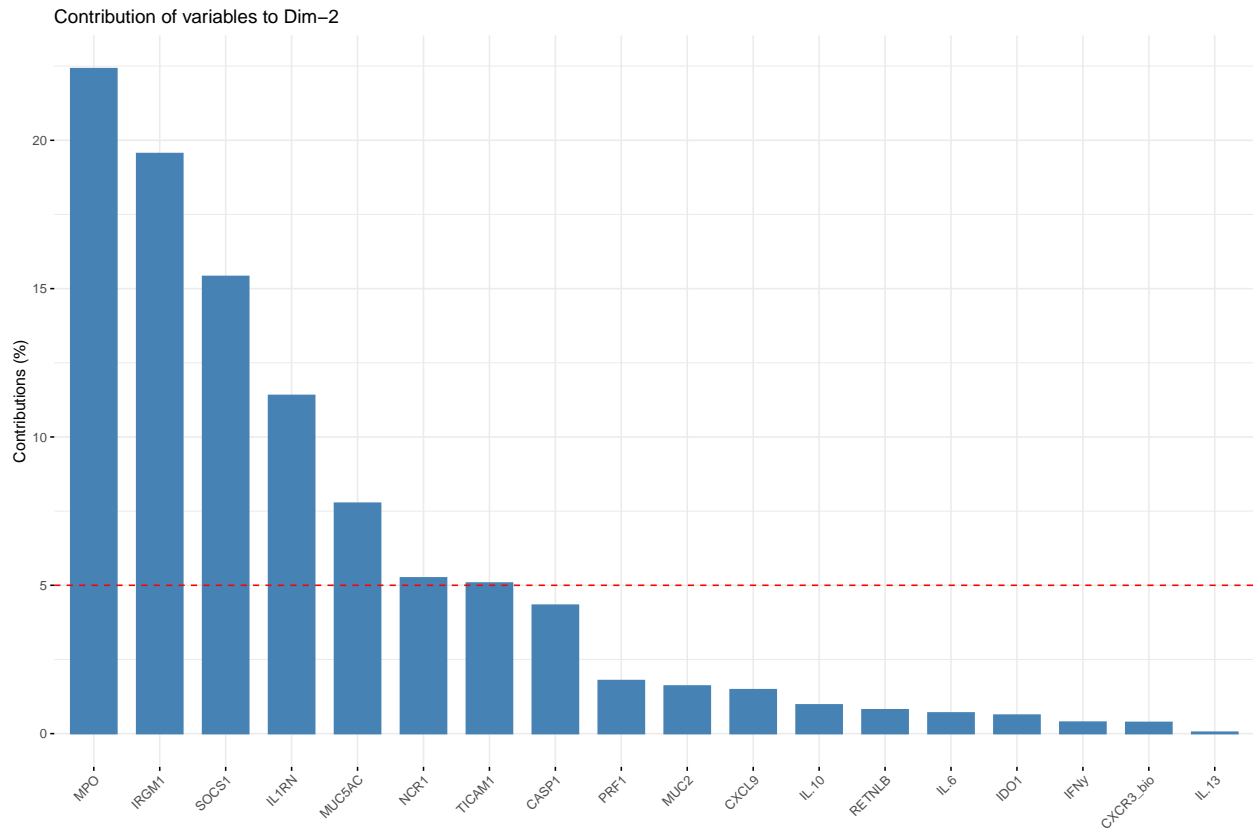


The function `fviz_contrib()` [factoextra package] can be used to draw a bar plot of variable contributions. If your data contains many variables, you can decide to show only the top contributing variables. The R code below shows the top 10 variables contributing to the principal components:





```
# Contributions of variables to PC2  
fviz_contrib(res.pca, choice = "var", axes = 2, top = 18)
```

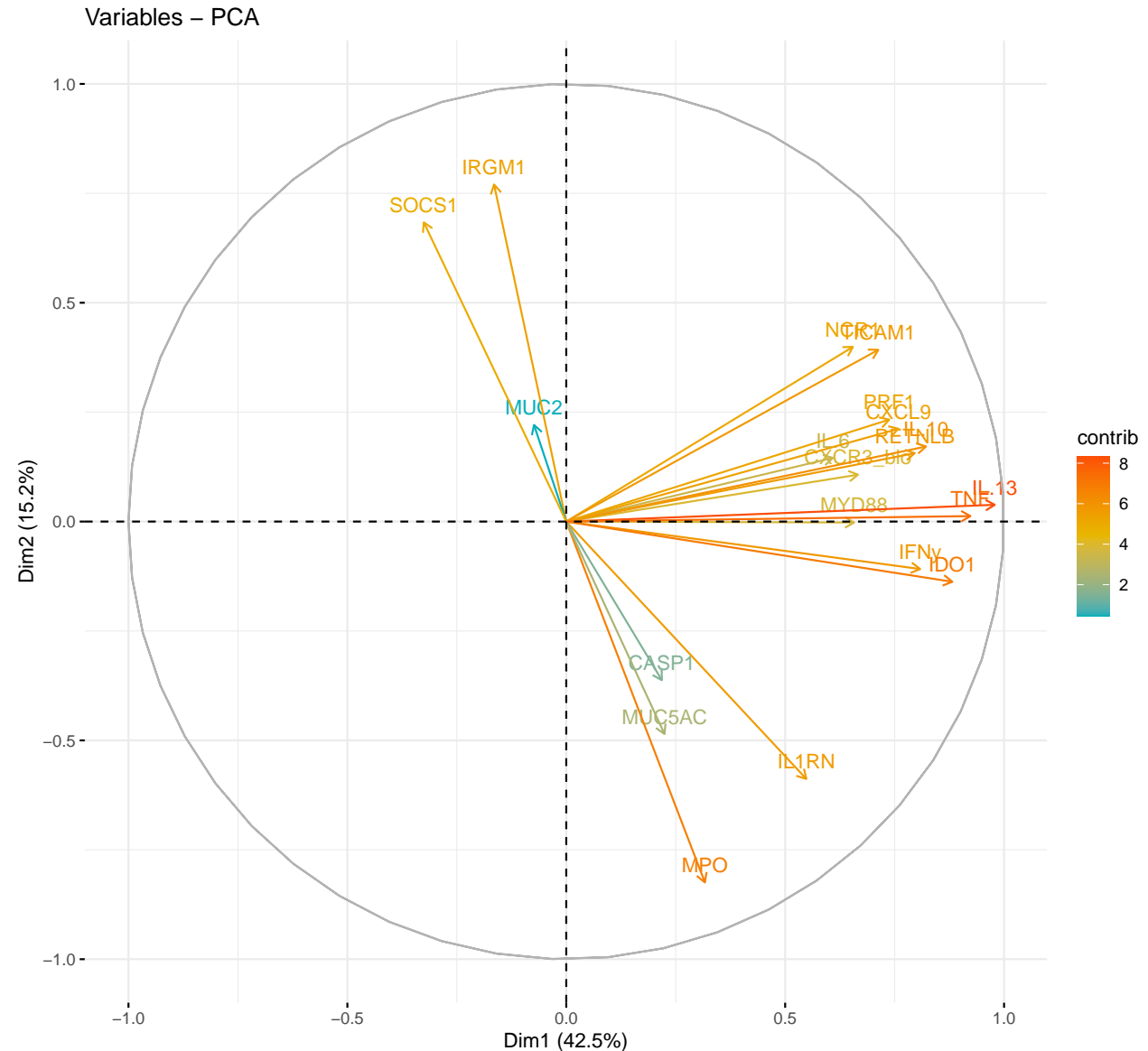


The red dashed line on the graph above indicates the expected average contribution. If the contribution of the variables were uniform, the expected value would be  $1/\text{length}(\text{variables}) = 1/10 = 10\%$ . For a given

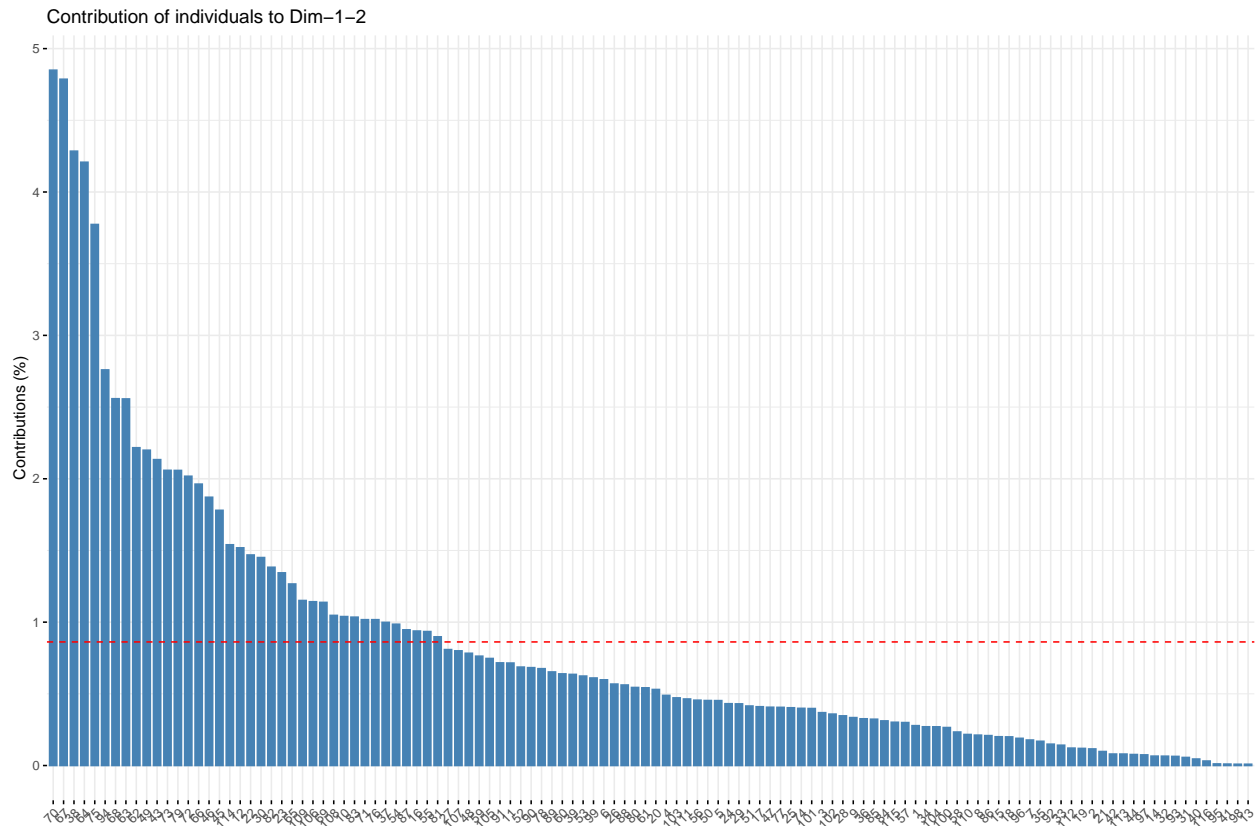
component, a variable with a contribution larger than this cutoff could be considered as important in contributing to the component.

Note that, the total contribution of a given variable, on explaining the variations retained by two principal components, say PC1 and PC2, is calculated as  $\text{contrib} = [(C1 * \text{Eig1}) + (C2 * \text{Eig2})]/(\text{Eig1} + \text{Eig2})$ , where

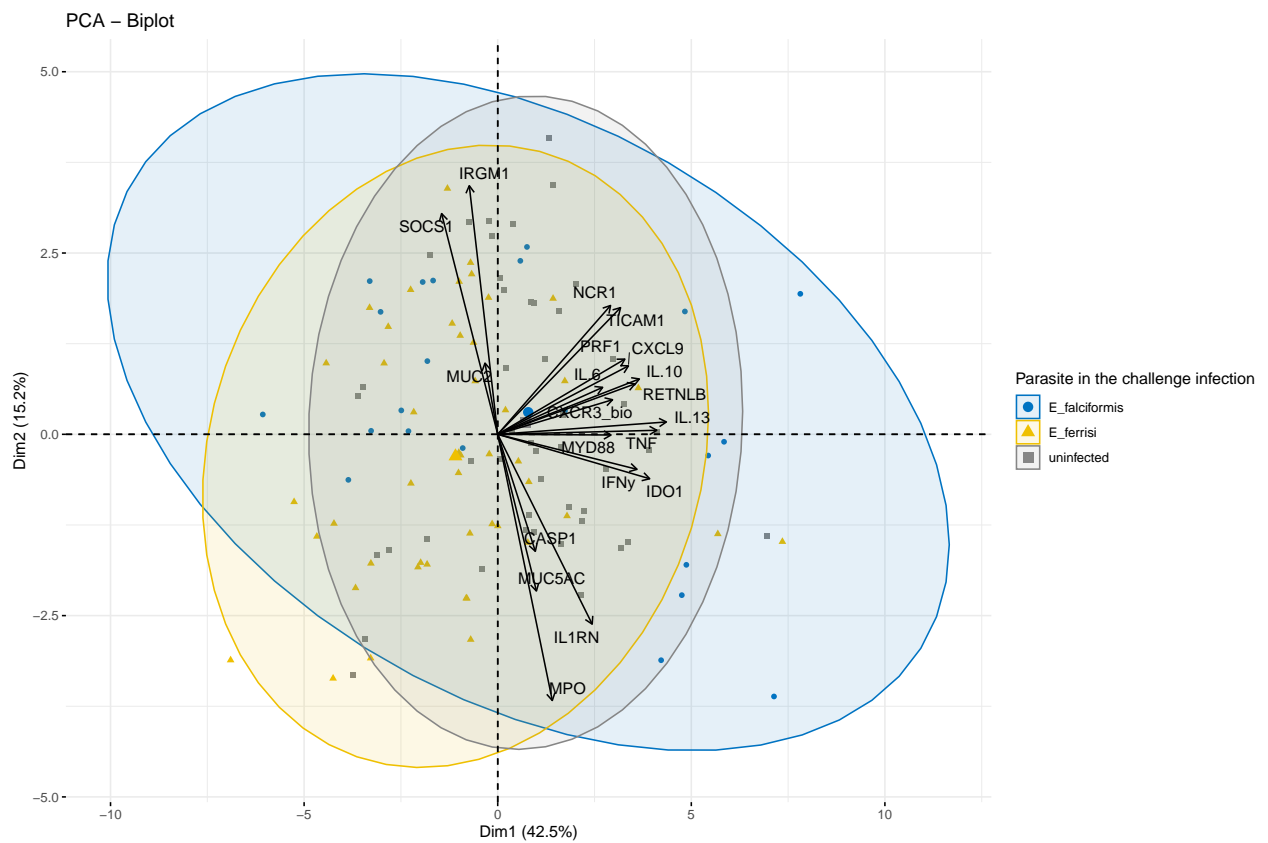
C1 and C2 are the contributions of the variable on PC1 and PC2, respectively Eig1 and Eig2 are the eigenvalues of PC1 and PC2, respectively. Recall that eigenvalues measure the amount of variation retained by each PC. In this case, the expected average contribution (cutoff) is calculated as follow: As mentioned above, if the contributions of the 10 variables were uniform, the expected average contribution on a given PC would be  $1/10 = 10\%$ . The expected average contribution of a variable for PC1 and PC2 is :  $[(10 * \text{Eig1}) + (10 * \text{Eig2})]/(\text{Eig1} + \text{Eig2})$



To visualize the contribution of individuals to the first two principal components:



## PCA + Biplot combination



In the following example, we want to color both individuals and variables by groups. The trick is to use `pointshape = 21` for individual points. This particular point shape can be filled by a color using the argument `fill.ind`. The border line color of individual points is set to “black” using `col.ind`. To color variable by groups, the argument `col.var` will be used.

Linear models:

```
##
## Call:
## lm(formula = max_WL ~ pc1 + pc2 + Parasite_challenge, data = imputed_expr)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -14.4014  -3.0944   0.1175   3.5262  14.3050
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      85.6719     1.1323  75.664 < 2e-16 ***
## pc1               0.1272     0.1762   0.722  0.4718
## pc2              -0.7278     0.2834  -2.568  0.0116 *
## Parasite_challengeE_ferrisi    6.0895     1.4092   4.321 3.40e-05 ***
## Parasite_challengeuninfected  10.4534     1.3576   7.700 6.09e-12 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 5.253 on 111 degrees of freedom
## Multiple R-squared:  0.3818, Adjusted R-squared:  0.3596
## F-statistic: 17.14 on 4 and 111 DF, p-value: 5.657e-11
## [1] 720.9133
##
## Call:
## lm(formula = max_WL ~ pc1 + pc2 + Parasite_challenge + hybrid_status,
##     data = imputed_expr)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -12.7141  -3.5997   0.4672   3.5380  13.9501
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      86.0601     1.3838  62.193 < 2e-16 ***
## pc1               0.1364     0.2199   0.620  0.5364
## pc2              -0.5959     0.3144  -1.896  0.0607 .
## Parasite_challengeE_ferrisi    5.9059     1.4538   4.062 9.34e-05 ***
## Parasite_challengeuninfected  10.0684     1.4516   6.936 3.30e-10 ***
## hybrid_statusF0 M. m. musculus -1.1985     1.4579  -0.822  0.4129
## hybrid_statusF1 hybrid         1.4620     1.6821   0.869  0.3867
## hybrid_statusF1 M. m. domesticus -1.7765     2.2126  -0.803  0.4238
## hybrid_statusF1 M. m. musculus   1.7684     2.6843   0.659  0.5115
## hybrid_statusother            -0.3591     1.4859  -0.242  0.8095
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 5.288 on 106 degrees of freedom
```

```
## Multiple R-squared:  0.4017, Adjusted R-squared:  0.3509
## F-statistic: 7.907 on 9 and 106 DF,  p-value: 7.337e-09
## [1] 727.1249
```

Try instead: LLR test (likelihood ration) (LM4 package )?

<https://www.rdocumentation.org/packages/lmtest/versions/0.9-38/topics/lrtest>

In this way you compare each model, with the different variables used to predict.

Another way is to compare the AIC. (function : step)

```
weight_lm3 <- lm(max_WL ~ pc1 + pc2 + hybrid_status, data = imputed_expr)
weight_no_pc1 <- lm(max_WL ~ pc2 + hybrid_status, data = imputed_expr)
weight_no_pc2 <- lm(max_WL ~ pc1 + hybrid_status, data = imputed_expr)
weight_no_hybrid <- lm(max_WL ~ pc1 + pc2, data = imputed_expr)
lrtest(weight_lm3, weight_no_pc1)
```

```
## Likelihood ratio test
##
## Model 1: max_WL ~ pc1 + pc2 + hybrid_status
## Model 2: max_WL ~ pc2 + hybrid_status
##   #Df LogLik Df  Chisq Pr(>Chisq)
## 1    9 -374.55
## 2    8 -375.87 -1  2.6379    0.1043
lrtest(weight_lm3, weight_no_pc2)
```

```
## Likelihood ratio test
##
## Model 1: max_WL ~ pc1 + pc2 + hybrid_status
## Model 2: max_WL ~ pc1 + hybrid_status
##   #Df LogLik Df  Chisq Pr(>Chisq)
## 1    9 -374.55
## 2    8 -374.96 -1  0.8221    0.3646
lrtest(weight_lm3, weight_no_hybrid)
```

```
## Likelihood ratio test
##
## Model 1: max_WL ~ pc1 + pc2 + hybrid_status
## Model 2: max_WL ~ pc1 + pc2
##   #Df LogLik Df  Chisq Pr(>Chisq)
## 1    9 -374.55
## 2    4 -379.64 -5 10.186    0.07014 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

##
## Call:
## lm(formula = max_WL ~ pc1 + pc2 + hybrid_status, data = imputed_expr)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -15.942  -3.138   0.991   4.739   9.868
##
## Coefficients:
##
##                                Estimate Std. Error t value Pr(>|t|)
```

```

## (Intercept)          92.5229      1.0319  89.664 <2e-16 ***
## pc1                   0.3934      0.2496   1.576  0.1179
## pc2                  -0.3243      0.3701  -0.876  0.3827
## hybrid_statusF0 M. m. musculus -1.1490      1.7436  -0.659  0.5113
## hybrid_statusF1 hybrid      3.7568      1.9749   1.902  0.0598 .
## hybrid_statusF1 M. m. domesticus -0.3187      2.6314  -0.121  0.9038
## hybrid_statusF1 M. m. musculus   3.9912      3.1916   1.251  0.2138
## hybrid_statusother    -2.5944      1.7376  -1.493  0.1383
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.332 on 108 degrees of freedom
## Multiple R-squared:  0.1259, Adjusted R-squared:  0.06928
## F-statistic: 2.223 on 7 and 108 DF, p-value: 0.03774

## [1] 767.095

##
## Call:
## lm(formula = max_WL ~ pc1 + pc2 + infection_history, data = imputed_expr)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -13.4909  -3.4963   0.2167   3.0235  14.3776
##
## Coefficients:
##
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    90.04699     1.95916  45.962 < 2e-16
## pc1             0.02434     0.17334   0.140  0.88862
## pc2            -0.61060     0.29415  -2.076  0.04035
## infection_historyfalciformis_ferrisi  2.11290     2.34416   0.901  0.36947
## infection_historyfalciformis_uninfected 6.75252     2.36949   2.850  0.00527
## infection_historyferrisi_falciformis  -7.57149     2.57368  -2.942  0.00401
## infection_historyferrisi_ferrisi    2.81356     2.33435   1.205  0.23080
## infection_historyferrisi_uninfected  5.45339     2.19195   2.488  0.01442
## infection_historyuninfected    7.17664     2.60969   2.750  0.00702
## infection_historyuninfected_falciformis -4.51846     2.83429  -1.594  0.11389
## infection_historyuninfected_ferrisi  -2.60534     2.67011  -0.976  0.33144
##
## (Intercept)          ***
## pc1
## pc2                   *
## infection_historyfalciformis_ferrisi
## infection_historyfalciformis_uninfected **
## infection_historyferrisi_falciformis   **
## infection_historyferrisi_ferrisi
## infection_historyferrisi_uninfected    *
## infection_historyuninfected            **
## infection_historyuninfected_falciformis
## infection_historyuninfected_ferrisi
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 5.021 on 105 degrees of freedom
## Multiple R-squared:  0.4656, Adjusted R-squared:  0.4147

```

```
## F-statistic: 9.149 on 10 and 105 DF, p-value: 1.007e-10
## [1] 716.0163
##
## Call:
## lm(formula = max_WL ~ pc1 + pc2, data = imputed_expr)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -16.913  -3.236   1.379   5.127  10.471
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  92.3746     0.6006  153.811  <2e-16 ***
## pc1           0.1702     0.2061   0.826   0.4107
## pc2          -0.7501     0.3448  -2.175   0.0317 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.468 on 113 degrees of freedom
## Multiple R-squared:  0.04572, Adjusted R-squared:  0.02883
## F-statistic: 2.707 on 2 and 113 DF, p-value: 0.07108
##
##              df      AIC
## weight_lm         6 720.9133
## weight_lm_exp_only  4 767.2808
```

repeating the heatmap on the now imputed data

```
gene <- imputed_expr %>% dplyr::select(c(EH_ID, all_of(Genes)))

# turn the data frame into a matrix and transpose it. We want to have each cell
# type as a row name
gene <- t(as.matrix(gene))

#switch the matrix back to a data frame format
gene <- as.data.frame(gene)

# turn the first row into column names
gene %>%
  row_to_names(row_number = 1) -> heatmap_data

table(rowSums(is.na(heatmap_data)) == nrow(heatmap_data))

##
## FALSE
##      20

# turn the columns to numeric other wise the heatmap function will not work
heatmap_data[] <- lapply(heatmap_data, function(x) as.numeric(as.character(x)))

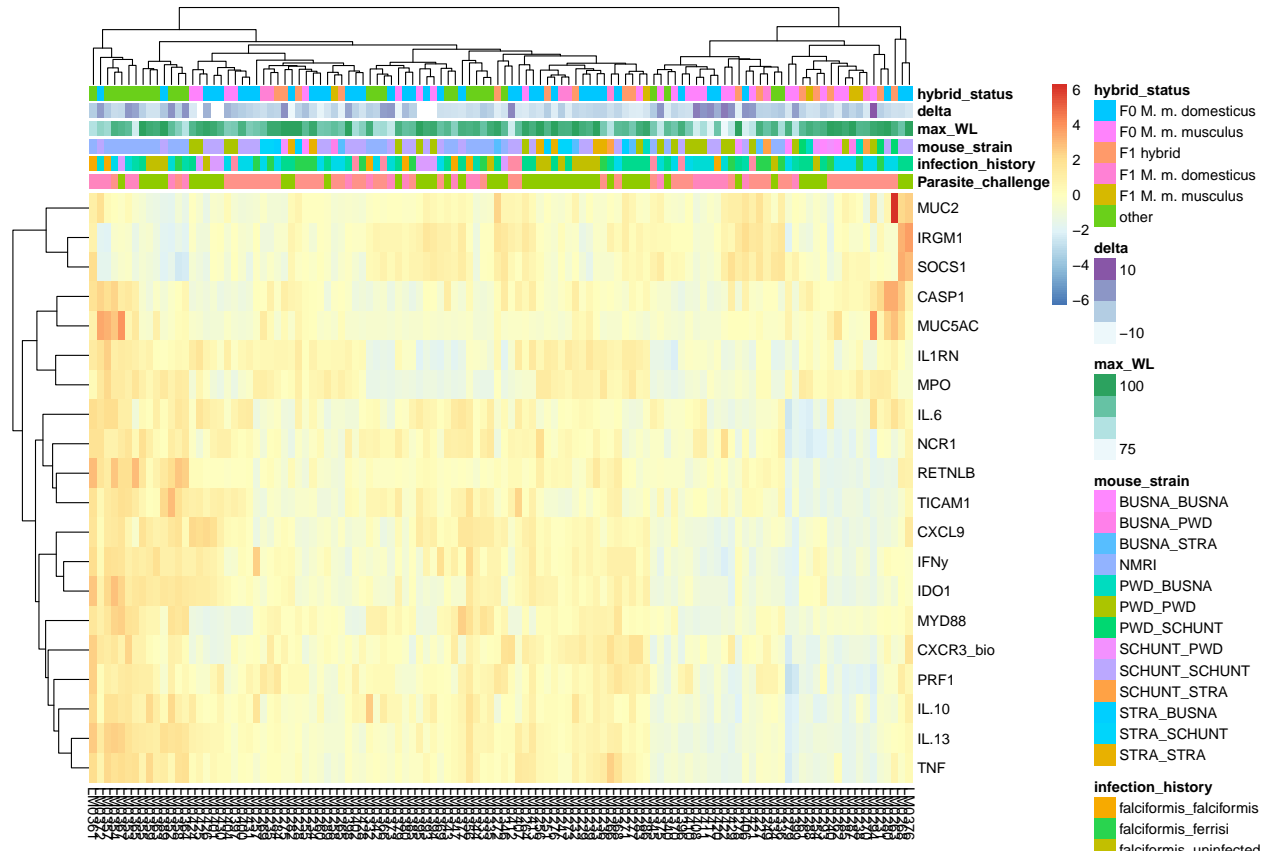
# remove columns with only NAs
heatmap_data <- Filter(function(x) !all(is.na(x)), heatmap_data)
```



```
#remove rows with only Nas
heatmap_data <- heatmap_data[, colSums(is.na(heatmap_data)) != nrow(heatmap_data)]

rownames(annotation_df) <- colnames(heatmap_data)
```

Heatmap on gene expression data:



repeating the heatmap on the significant genes imputed data

```
gene <- imputed_expr %>% dplyr::select(c(EH_ID, c("IFNy", "IL.13", "PRF1",
"PRF1", "TICAM1")))

# turn the data frame into a matrix and transpose it. We want to have each cell
# type as a row name
gene <- t(as.matrix(gene))

#switch the matrix back to a data frame format
gene <- as.data.frame(gene)

# turn the first row into column names
gene %>%
  row_to_names(row_number = 1) -> heatmap_data

table(rowSums(is.na(heatmap_data)) == nrow(heatmap_data))
```

```
##
```

```

## FALSE
##      4

# turn the columns to numeric other wise the heatmap function will not work
heatmap_data[] <- lapply(heatmap_data, function(x) as.numeric(as.character(x)))

# remove columns with only NAs
heatmap_data <- Filter(function(x)!all(is.na(x)), heatmap_data)

#remove rows with only NAs
heatmap_data <- heatmap_data[, colSums(is.na(heatmap_data)) != nrow(heatmap_data)]

annotation_df <- as_tibble(Challenge) %>%
  dplyr::filter(infection == "challenge", dpi == dpi_max) %>%
  dplyr::group_by(EH_ID) %>%
  dplyr::select(c("EH_ID", "Parasite_challenge",
                  "hybrid_status")) %>%
  dplyr::filter(EH_ID %in% colnames(heatmap_data))

annotation_df <- unique(annotation_df)

annotation_df <- as.data.frame(unique(annotation_df)) %>%
  dplyr::select(-EH_ID)

### Prepare the annotation columns for the heatmap
rownames(annotation_df) <- annotation_df$EH_ID

rownames(annotation_df) <- colnames(heatmap_data)

```

Heatmap on gene expression data:



```

))

g2 <- g2 %>%
  mutate(Eim_MC = case_when(
    Eim_MC == "TRUE" ~ "infected",
    Eim_MC == "FALSE" ~ "uninfected"
  ))

g2$Eim_MC <- as.factor(g2$Eim_MC)

require(devtools)

## Loading required package: devtools
## Loading required package: usethis
devtools::install_github("alicebalard/parasiteLoad@v2.0", force = TRUE)

## Downloading GitHub repo alicebalard/parasiteLoad@v2.0
##
## * checking for file '/tmp/RtmpVkJetF/remotes32ddaf5003f288/alicebalard-parasiteLoad-1b43216/DESCRIPTION' ... OK
## * preparing 'parasiteLoad':
## * checking DESCRIPTION meta-information ... OK
## * checking for LF line-endings in source and make files and shell scripts
## * checking for empty or unneeded directories
## * building 'parasiteLoad_0.1.0.tar.gz'

## Installing package into '/localstorage/fay/R/x86_64-pc-linux-gnu-library/4.2'
## (as 'lib' is unspecified)

library(parasiteLoad)

str(g2$max_WL)

## num [1:90] 4.11 0 3.57 1.26 0 ...

parasiteLoad::getParamBounds("weibull", data = g2, response = "max_WL")

##      L1start      L1LB      L1UB      L2start      L2LB      L2UB
## 7.158694975 0.000000001 26.554054054 7.158694975 0.000000001 26.554054054
## alphaStart    alphaLB    alphaUB myshapeStart    myshapeLB    myshapeUB
## 0.000000000 -5.000000000 5.000000000 1.000000000 0.000000001 5.000000000

speparam <- c(L1start = 7.158694975,
              L1LB = 0.000000001,
              L1UB = 26.554054054,
              L2start = 7.158694975,
              L2LB = 0.000000001,
              L2UB = 26.554054054,
              alphaStart = 0, alphaLB = -5, alphaUB = 5,
              myshapeStart = 1, myshapeLB = 0.000000001, myshapeUB = 5)

##All
#fitWL_MC <- parasiteLoad::analyse(data = g2,

```

```

# response = "max_WL",
# model = "weibull",
# group = "Eim_MC")

##fitWL_MC

#plot_WL_MC<- bananaPlot(mod = fitWL_MC$H3,
#      #      data = g2,
#      #      response = "max_WL",
#      #      group = "Eim_MC") +
#      # scale_fill_manual(values = c("blue", "red")) +
#      # scale_color_manual(values = c("blue", "red")) +
#      # theme_bw()
#plot_WL_MC

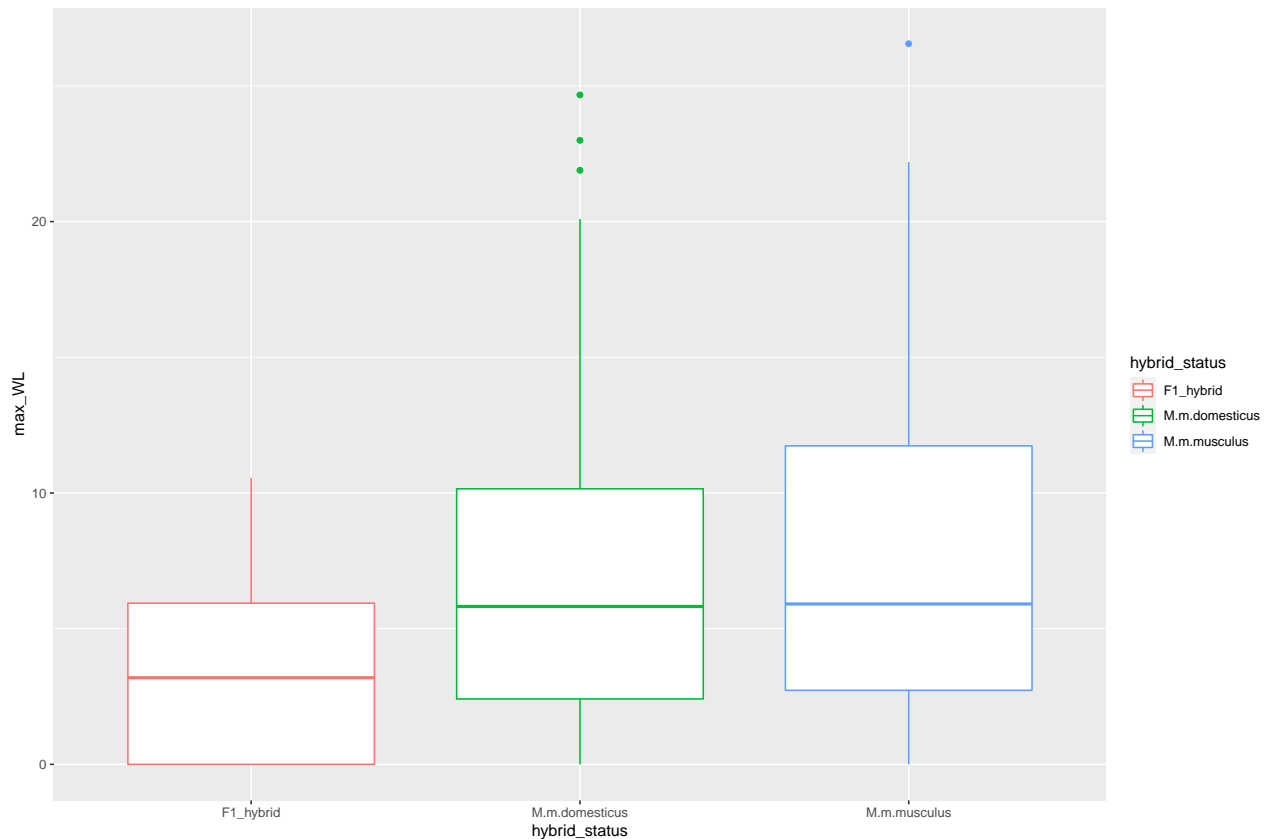
HI_weight <- lm(max_WL ~ hybrid_status, data = g2)
summary(HI_weight)

##
## Call:
## lm(formula = max_WL ~ hybrid_status, data = g2)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -8.569 -4.712 -1.856  2.824 17.985
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      3.680      1.638   2.246  0.0272 *
## hybrid_statusM.m.domesticus  3.806      1.908   1.995  0.0491 *
## hybrid_statusM.m.musculus   4.889      2.041   2.395  0.0187 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.554 on 87 degrees of freedom
## Multiple R-squared:  0.06415,    Adjusted R-squared:  0.04263
## F-statistic: 2.982 on 2 and 87 DF,  p-value: 0.05592
HI_weight$coefficients

##              (Intercept) hybrid_statusM.m.domesticus
##              3.680309              3.806160
## hybrid_statusM.m.musculus
##              4.888880

ggplot(data = g2, aes(x = hybrid_status, y = max_WL, color = hybrid_status)) +
  geom_boxplot()

```



```
AIC(HI_weight)
```

```
## [1] 598.7605
```

```
one.way <- aov(max_WL ~ hybrid_status, data = g2)
```

```
summary(one.way)
```

```
##              Df Sum Sq Mean Sq F value Pr(>F)
## hybrid_status  2    256  128.06    2.982 0.0559 .
## Residuals    87   3737   42.95
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

The Df column displays the degrees of freedom for the independent variable (the number of levels in the variable minus 1), and the degrees of freedom for the residuals (the total number of observations minus one and minus the number of levels in the independent variables). The Sum Sq column displays the sum of squares (a.k.a. the total variation between the group means and the overall mean). The Mean Sq column is the mean of the sum of squares, calculated by dividing the sum of squares by the degrees of freedom for each parameter. The F-value column is the test statistic from the F test. This is the mean square of each independent variable divided by the mean square of the residuals. The larger the F value, the more likely it is that the variation caused by the independent variable is real and not due to chance. The Pr(>F) column is the p-value of the F-statistic. This shows how likely it is that the F-value calculated from the test would have occurred if the null hypothesis of no difference among group means were true. The p-value of the fertilizer variable is low ( $p < 0.001$ ), so it appears that the type of fertilizer used has a real impact on the final crop yield.

```
write.csv(imputed_expr, "output_data/gene_expression/data_products/imputed_gene_expression.csv", row.names = FALSE)
```

```
write.csv(g2, "output_data/gene_expression/data_products/clean_gene_expression.csv", row.names = FALSE)
```