

# Project 1: Function to Summarize and Plot Gene Expression

2024-08-07

## Import and summarize data

```
genes <- read.csv('/Users/isabellabaldacci/Desktop/QBS103/Project/QBS103_GSE157103_genes.csv',
                  sep = ',', header = TRUE)
metadata <- read.csv('/Users/isabellabaldacci/Desktop/QBS103/Project/QBS103_GSE157103_series_matrix.csv',
                    header = TRUE)
dim(genes)
```

```
## [1] 100 127
```

## Clean up metadata

```
## cleaning up metadata
# metadata were removed if they were not helpful or if they had too many unknowns
metadata_clean <- dplyr::select(metadata,
                                -c("X.Sample_submission_date", "channel_count", "status", "last_update_date"))
#head(metadata_clean)
```

## Transpose Genes dataframe

```
#transposing the genes x samples matrix to samples x genes
rownames(genes) <- genes$X #assigning row names
genes_t <- as.data.frame(t(dplyr::select(genes, -c('X')))) #transpose as a dataframe
genes_t$participant_id <- rownames(genes_t) #assigning participant id column for later joining
#which(genes_t$participant_id == "COVID_06.y_male_NonICU") #finding column with mismatch to metadata
genes_t$participant_id[6] <- "COVID_06:y_male_NonICU" #reassigning value to match metadata
```

## Combining metadata and genes dataframe

```
#combining the samples x genes matrix to the metadata
all_data <- dplyr::inner_join(genes_t, metadata_clean, by = 'participant_id')
```

## Write function for plotting

```

plots <- function(df, genes_list, cont_cov, cat_cov) {
  df_filtered <- df %>% dplyr::select(participant_id, all_of(genes_list), all_of(cont_cov), all_of(cat_

  df_filtered_long <- df_filtered %>% tidyr::pivot_longer(cols = genes_list, names_to = 'Gene', values_
    dplyr::arrange(Gene, .data[[cont_cov]])

  for (gene in genes_list){

    one_gene <- df_filtered_long %>% dplyr::filter(Gene == gene)

    mean_expression <- round(mean(one_gene$Expression), 2)
    sd_expression <- round(sd(one_gene$Expression), 2)
    x_min_box <- quantile(one_gene$Expression, 3/4)

    n <- length(one_gene$Expression)
    hist_title <- paste('Distribution of ', gene, ' expression across samples')
    hist_x <- paste(gene, ' Expression')

    hist <- one_gene %>% ggplot(aes(Expression)) +
      geom_histogram(bins = 30, fill = 'slateblue1') +
      theme_minimal() +
      theme(plot.title = element_text(hjust = .4, size = 20)) +
      labs(title = hist_title,
        x = hist_x,
        y = 'Frequency',
        size = 20) +
      scale_color_manual(values = c('slateblue1')) +
      annotate('rect',
        xmin = x_min_box + 5,
        xmax = x_min_box + 25,
        ymin = 10,
        ymax = 14, fill = 'grey', alpha = .9)+
      annotate(geom = 'text', x = x_min_box + 15, y = 12, label = paste('N:', n, ' samples',
        '\n Mean Expression:', mean_expression,
        '\n sd:', sd_expression), size = 3)

    print(hist)

    scatter_filtered <- one_gene %>% filter(one_gene[[cont_cov]] != ' unknown') #filter out unknowns
    scatter_filtered[cont_cov] <- as.numeric(scatter_filtered[[cont_cov]])
    if (cont_cov == 'crp.mg.l.') {
      scatter_filtered$crp.level <- cut(scatter_filtered$crp.mg.l., breaks = c(0, 3, 10, 100, 500000),
        labels = c('normal', 'moderate', 'high', 'severe'))
    }
    colorPalette_scatter <- c('azure4', 'skyblue2', 'slateblue2', 'midnightblue')
    title_scatter <- paste(toupper(cont_cov), ' vs ', gene, ' Expression')
    x_scatter <- paste(toupper(cont_cov), '(mg/L)')
    y_scatter <- paste(gene, ' Expression')

    scatter <- scatter_filtered %>% ggplot(aes(x = .data[[cont_cov]], y = Expression, color = crp.level)) +
      geom_point(na.rm = TRUE) +
      theme(
        axis.text.x = element_text(angle=90),

```

```

    plot.title = element_text(hjust = .4, size = 20),
    legend.position = c(.85,.75)) +
  labs(title = title_scatter, x = x_scatter, y = y_scatter) +
  scale_color_manual(labels = c('Normal (<3mg/L)', 'Moderate (3-10mg/L)', 'High (10 -100mg/L)', 'Se
    colorPalette_scatter)

print(scatter)

title_box <- paste(gene, ' Expression across ', cat_cov[1], ' and ', cat_cov[2])
x_lab_box <- cat_cov[1]
y_lab_box <- paste(gene, ' Expression')
n_covid_pos <- length(dplyr::filter(
  one_gene, disease_status == 'disease state: COVID-19')$Expression)
n_covid_neg <- length(dplyr::filter(
  one_gene, disease_status != 'disease state: COVID-19')$Expression)
y_annot <- max(one_gene$Expression)

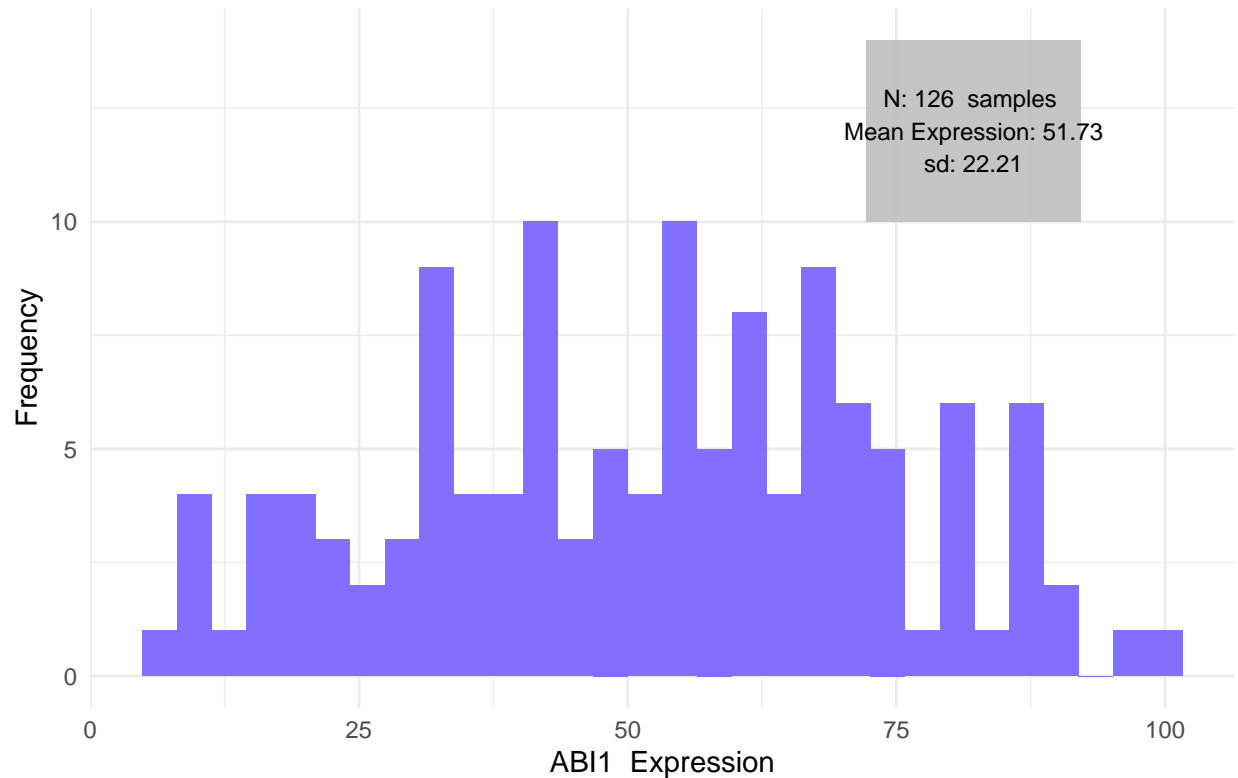
box <- one_gene %>% ggplot(aes_string(x = cat_cov[1], y = 'Expression', fill = cat_cov[2])) +
  geom_boxplot() +
  theme_minimal() +
  theme(legend.position = 'bottom') +
  labs(title = title_box,
    x = x_lab_box,
    y = y_lab_box) +
  scale_fill_manual(labels = c('no ventilator', 'ventilator'),
    values = c('ghostwhite', 'mediumvioletred')) +
  scale_x_discrete(labels = c('COVID', 'NON-COVID')) +
  annotate(geom = 'text', x = 1, y = y_annot + 5,
    label = paste('Covid positive: N = ', n_covid_pos),
    color = 'black') +
  annotate(geom = 'text', x = 2, y = y_annot + 5,
    label = paste('Covid negative: N = ', n_covid_neg),
    color = 'black')
print(box)
}
}

plots(all_data, c('ABI1', 'AAMP', 'ABHD17A'), c('crp.mg.l.'), c('disease_status', 'mechanical_ventilati

## Warning: Using an external vector in selections was deprecated in tidysselect 1.1.0.
## i Please use 'all_of()' or 'any_of()' instead.
## # Was:
## data %>% select(genes_list)
##
## # Now:
## data %>% select(all_of(genes_list))
##
## See <https://tidysselect.r-lib.org/reference/faq-external-vector.html>.
## This warning is displayed once every 8 hours.
## Call 'lifecycle::last_lifecycle_warnings()' to see where this warning was
## generated.

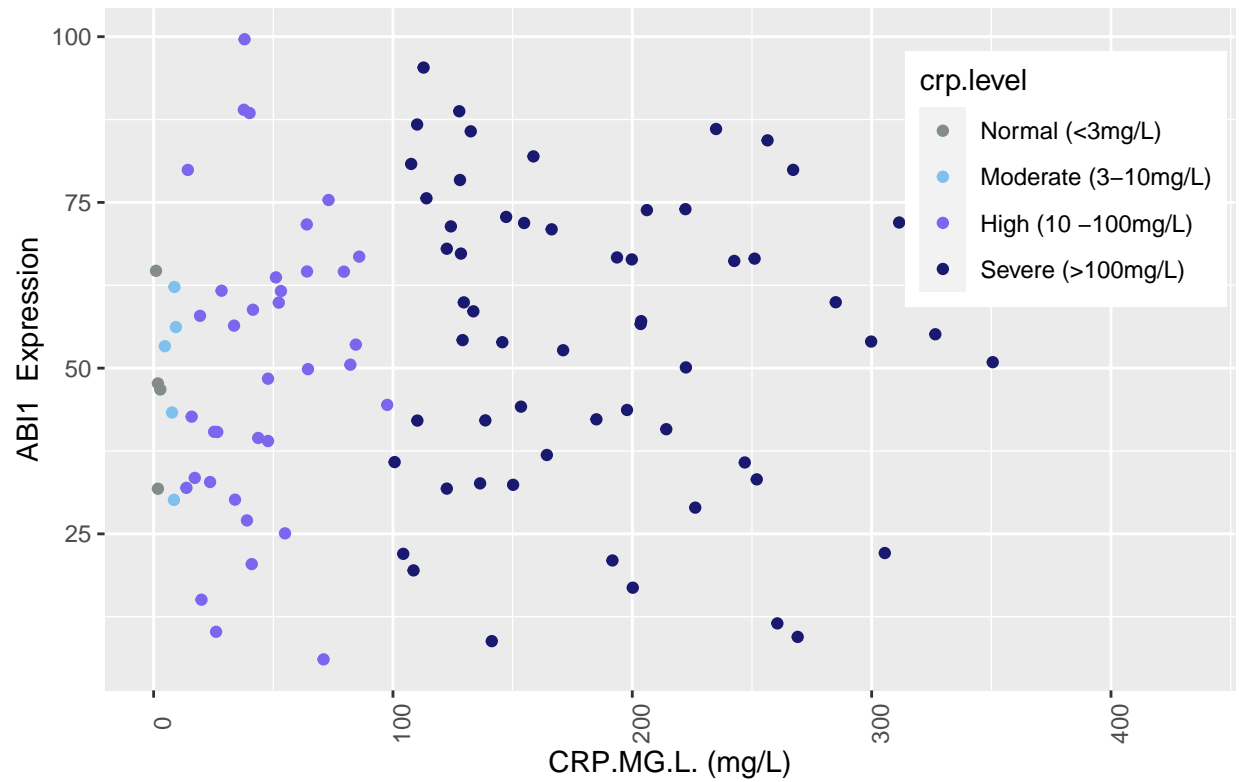
```

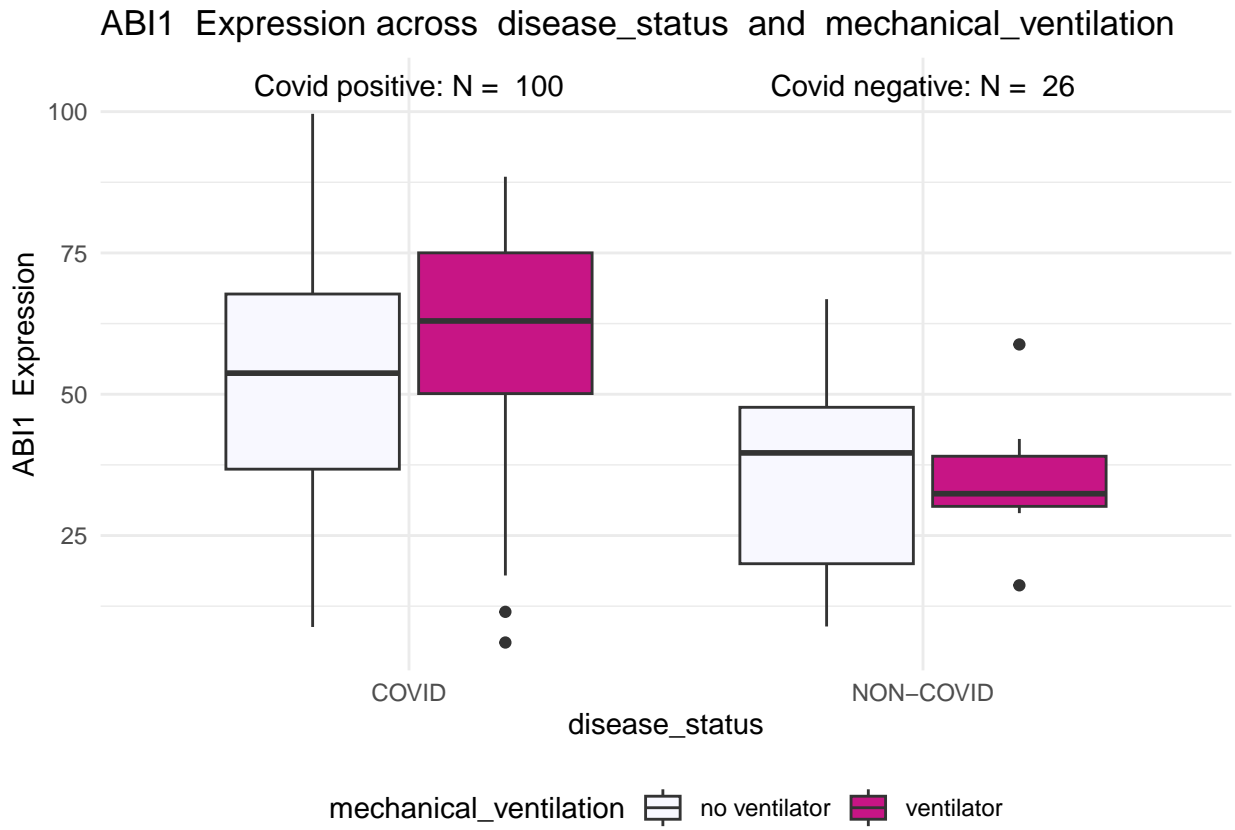
## Distribution of ABI1 expression across samples



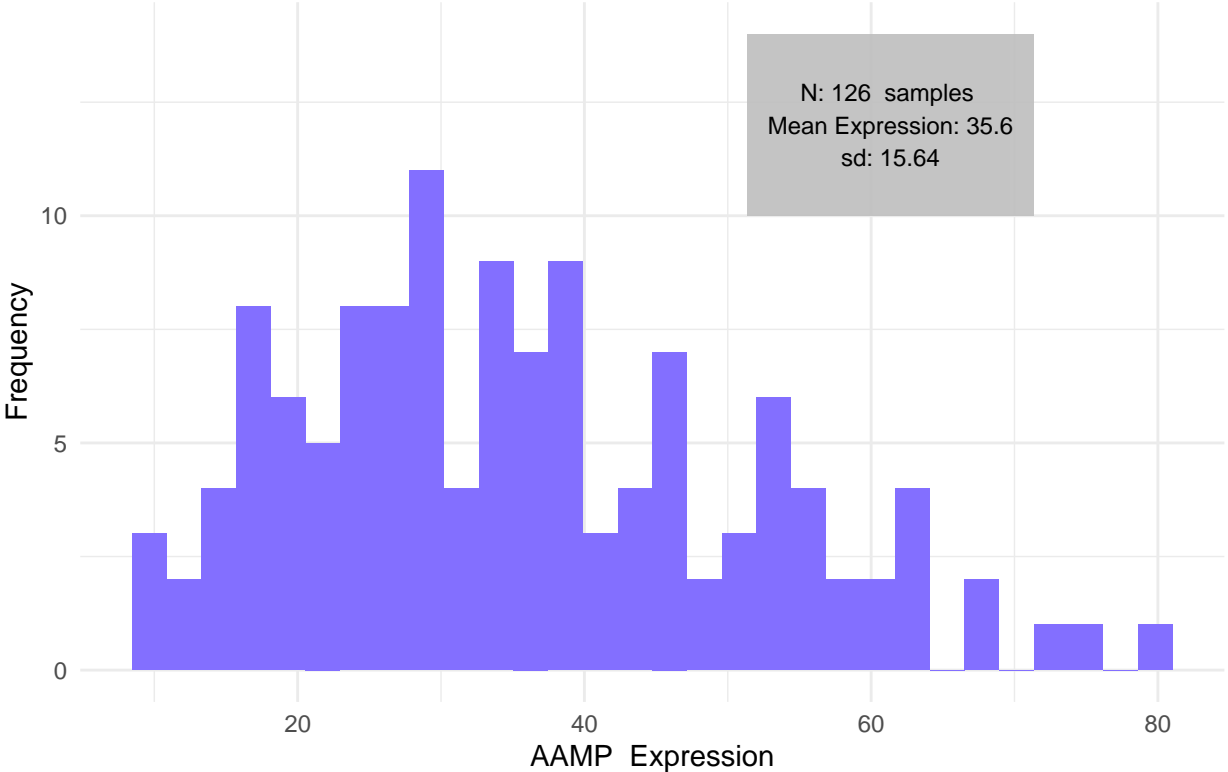
```
## Warning: 'aes_string()' was deprecated in ggplot2 3.0.0.  
## i Please use tidy evaluation idioms with 'aes()'.  
## i See also 'vignette("ggplot2-in-packages")' for more information.  
## This warning is displayed once every 8 hours.  
## Call 'lifecycle::last_lifecycle_warnings()' to see where this warning was  
## generated.
```

## CRP.MG.L. vs ABI1 Expression

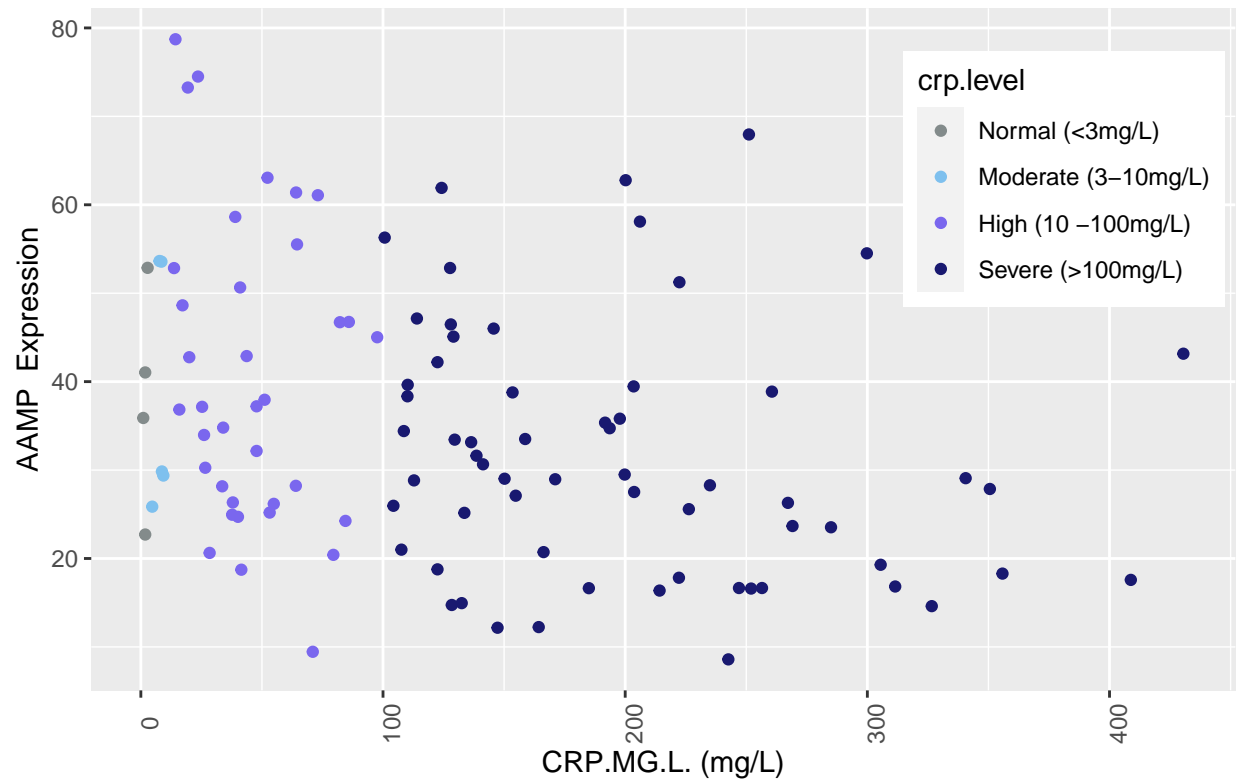




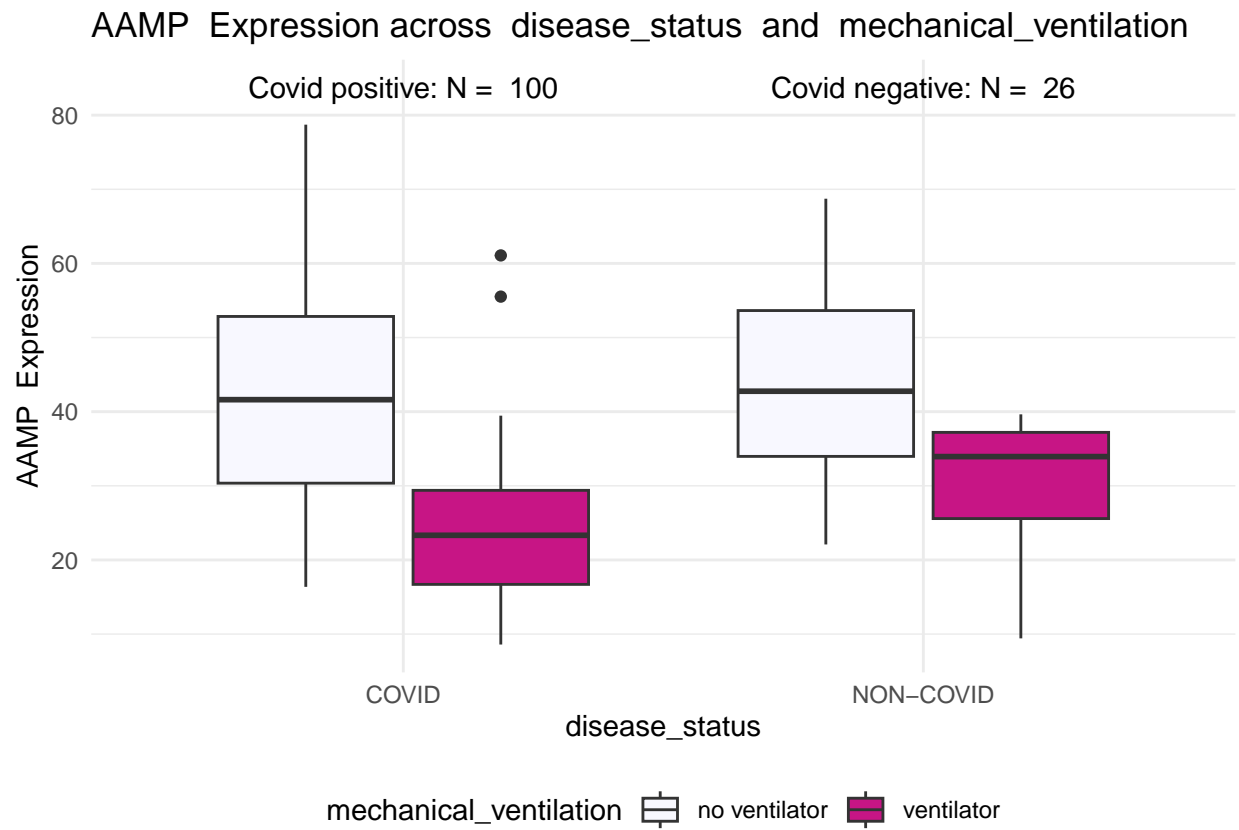
# Distribution of AAMP expression across samples



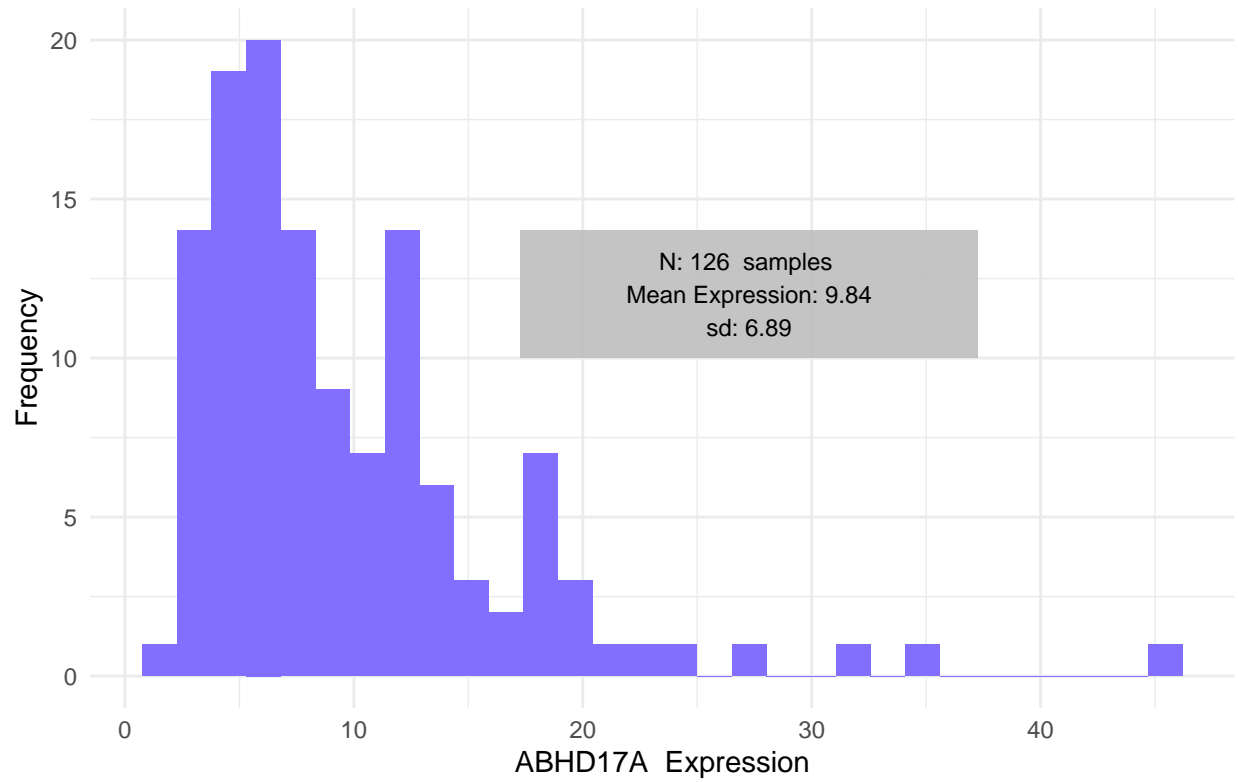
## CRP.MG.L. vs AAMP Expression







## Distribution of ABHD17A expression across sampl



## CRP.MG.L. vs ABHD17A Expression

