

BINF 6210 Assignment 5:

Introduction:

As the threat of antibiotic resistance crisis is on the rise, it has pertained scientists to understand the mechanisms in which makes bacterial infection resilient against antibiotics (1). *Pseudomonas aeruginosa* is a gram-negative bacterium, that has proven to be highly resistant against antibiotics and is a main causation for morbidity and mortality in immunocompromised patients with cystic fibrosis (2)(3). *P. aeruginosa*, has many acquired mechanisms towards antimicrobial resistance, however the most prominent is the biofilm formation, which has the ability overcome effects of multiple drugs (2). Scientists in North America have been researching the resistance mechanism of *P. aeruginosa*, however in recent light the host immune response to bacterial infections and what resistant mechanisms and protein abundances are upregulated are being researched.

In Yang *et al.* (2023), *P.aeruginosa* infected mouse macrophages transcriptomics were analyzed to determine the transcription factors involved in host immune response (4). This study was conducted in effort to examine which genes were upregulated and downregulated, when lung tissues were exposed to bacterial infection and found high abundance of inflammatory proteins and chromatin accessibility (4). During this project it was investigated whether the protein expressions were altered during infection and if the expressions was altered, which proteins were upregulated by identifying them based on their intensity at a particular retention time. The Yang *et al*, (2023) study concluded that inflammation proteins were highly upregulated during infection, in this project top ten proteins based on highest abundance were filtered out and tested whether the proteins found in the study to be highly upregulated were also found in high abundance during this project.

Figures:

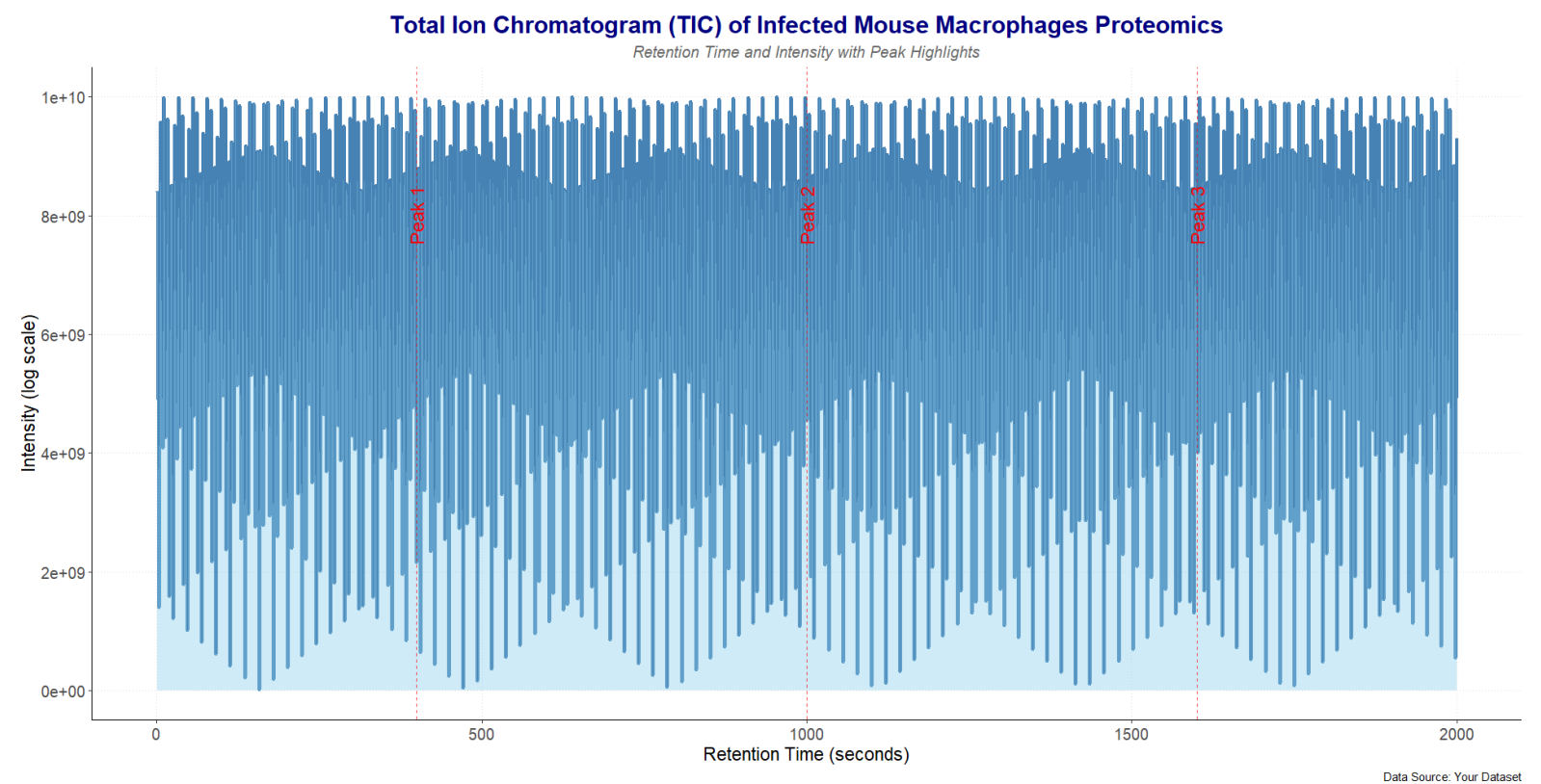


Figure 1:

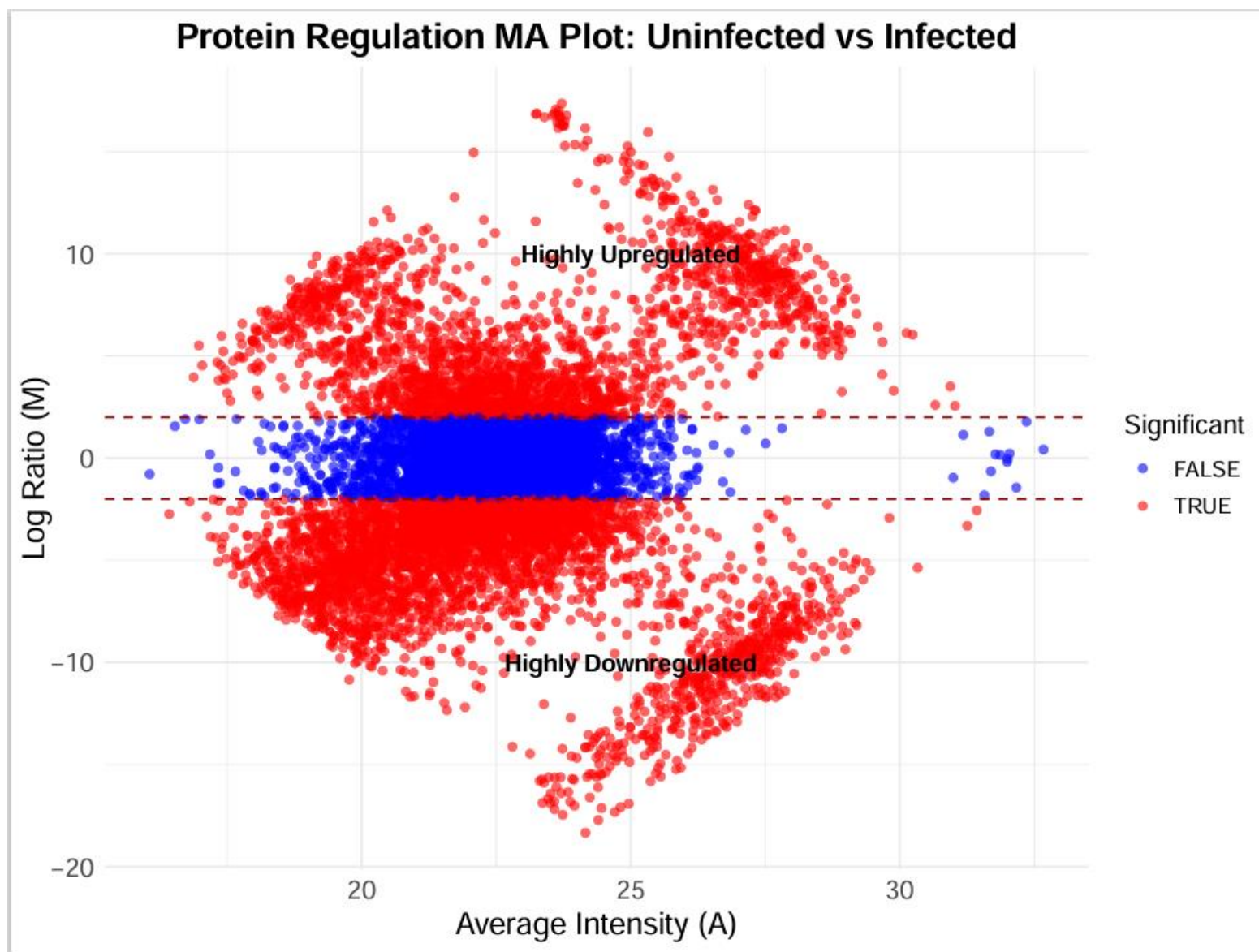


Figure 2:

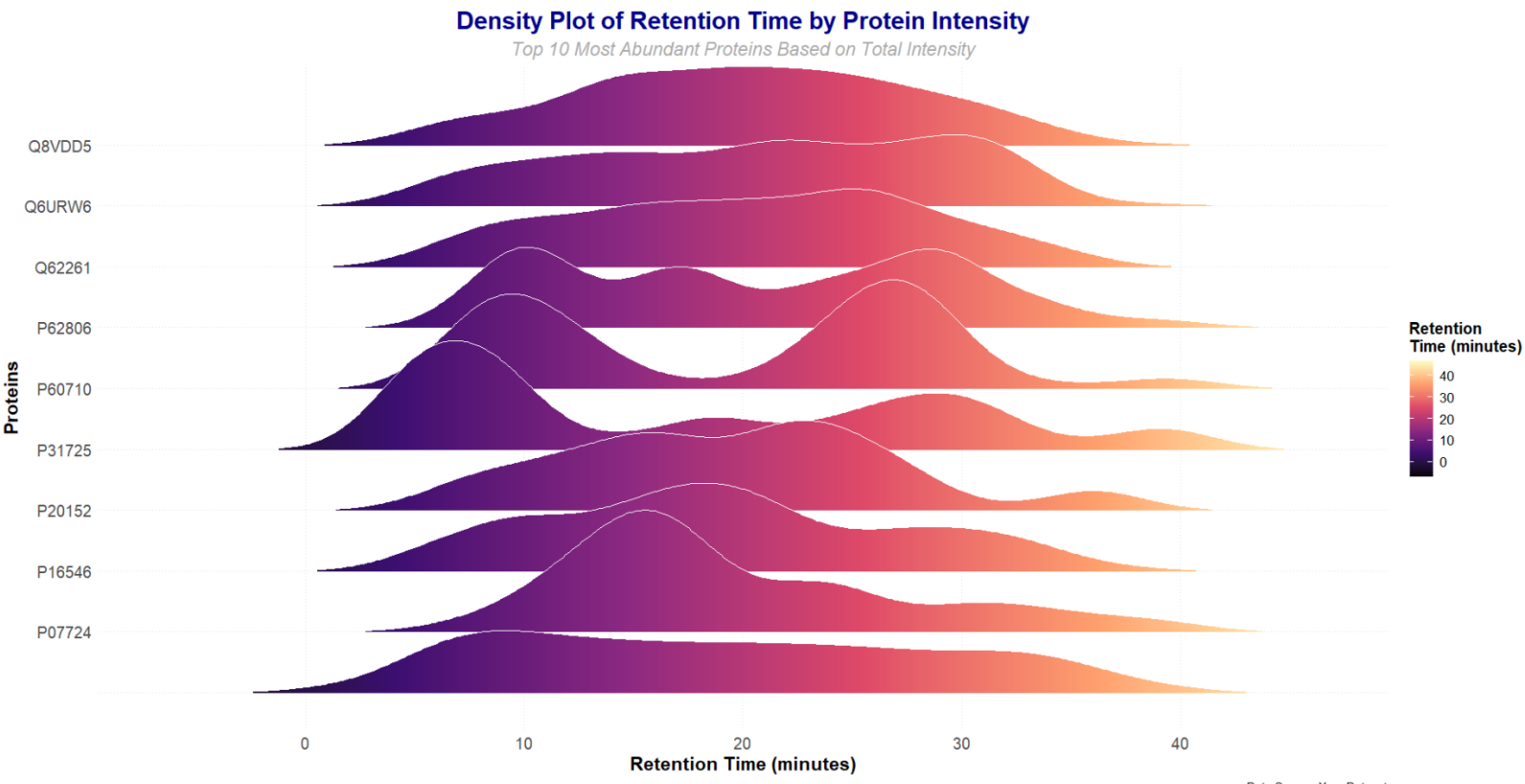


Figure 3:

Coding Section:

```
1 # Introduction
2 # Please install the packages below and set to working directory to the source file.
3
4 library(BiocManager)
5 library("RColorBrewer")
6 library("ggplot2")
7 library("mzR")
8 library("msdata")
9 library("reshape2")
10 library(mzR)
11 library(rpx)
12 library(limma)
13 library(DEqMS)
14 library(Biostrings)
15 library(sequinr)
16 library(MSnbase)
17 library("RforProteomics")
18 library(MSnID)
19 library(ggribes)
20 library(ggrepel)
21 library(dplyr)
22 library(reshape2)
23
24 # First the set.seed function was used to provide with accurate and repeatable results.
25
26 set.seed(123)
27
28 #This project will be comprised of whether the findings concluded from Yang et al,(2023) can be mimicked. The paper analyzed the proteomic data of infect mouse macrophage in MaxQuant
  and generated visual data using R. They found that certain proteins were upregulated when infected and found inflammatory proteins in high abundance and were upregulated.
29
30 #So we are going to find if proteins expression altered and what proteins.
31
32 #This proteomics data was downloaded from PRIDE database and the accession number from obtained from a literature.This data pertained to a P.seudomonas infected mouse macrophage.The
  file path goes to the source of the file and locates it.
33 file_path <- "C:/Users/ebmaz/Documents/Assignment 3/Proteomics_04.mzML"
34
35 #The raw_data reads the proteomics data and below the summary of the data is provided.
36 raw_data <- readMSData(file_path, mode = "onDisk")
37
38 #Section 1:Filtering
```

```
38 #Section 1:Filtering
39
40 #Length the data is determined,which is 20667 spectra.
41 summary(raw_data)
42
43 #Info about the raw data, such as the retention time range, number of spectra, and sample names.
44 print(raw_data)
45
46 # This code pertains to the column names in the raw data and the types of info in the data.
47 colnames(fData(raw_data))
48
49 #The actual numbers/data pertaining to the columns, such as the retention times, base peaks, and total ion current, is printed.
50 head(fData(raw_data))
51
52 #As seen there a fair bit of NAs, so we are going to focus on the data, which will be most useful for our research.Thus, retention time is being filtered out from the raw data.
53 retention_times <- fData(raw_data)$retentionTime
54
55 # Summary statistics of the retention times are provided to understand what the range is what possible proteins can be identified.
56 head(retention_times)
57 summarize_stats(retention_times, "Retention Times")
58
59 #The spectra data is also filtered to see and measure the protein abundance in terms of intensity.
60 spectra(raw_data[1:5])
61 #Spectra plot was graphed to see protein intensity verses M/Z
62 plot(raw_data[[1]])
63
64 # Here we are going to filter out the intensity values and determine it stats. Such as the highest intensity, so we can later note the high protein abundances.
65 mz_values <- mz(spectrum)
66 intensity_values <- intensity(spectrum)
67
68 highest_intensity_index <- which.max(intensity_values)
69
70 highest_mz <- mz_values[highest_intensity_index]
71 highest_intensity <- intensity_values[highest_intensity_index]
72
73 # The highest intensity peak and highest m/z values are noted.
74 cat("Highest Intensity Peak:\n")
75 cat("m/z:", highest_mz, "\n")
76 cat("Intensity:", highest_intensity, "\n")
77
```

```

78 #The summary of intensity measured in log was calculated to undersatnd the range of the intensity.
79 intensities <- unlist(lapply(spectra(raw_data), function(spectrum) intensity(spectrum)))
80 summary(intensities)
81
82 #The Shapiro test was performed to see if there was normal disturbation of the data and it was found that the data indeed did not have normal distrubution. This is because W = 0
83 subsampled_intensities <- sample(intensities, size = 5000, replace = FALSE)
84
85 shapiro_test <- shapiro.test(subsampled_intensities)
86 cat("Shapiro-Wilk Test:\n")
87 print(shapiro_test)
88
89 #In this code block we will see if there are any Na's in the data again, to make sure if there is any other we can use for out project. But since Alot of the NAs are in the data which
90 experiment_metadata <- fData(raw_data)
91 head(experiment_metadata)
92
93 colSums(is.na(experiment_metadata))
94
95
96 experiment_metadata[apply(is.na(experiment_metadata), 1, any), ]
97
98 # Section 2:Generating TIC and MA plot
99
100 #In the code chunks below the data frame for TIC was accumulated, so the retention time and intensity of Total ionization current. This first was to visualiz the intensity of the
101 protein at different retention times. This was done if effort to come one step closer to identify the proteins by measuring at what retention time are the intensity peaks the highest.
102 It was found that at 400s,1000s,1600s the intensity peaks were the highest.
103
104 tic_data <- data.frame(
105   retentionTime = seq(1, 2000, length.out = 1000),
106   intensity = abs(sin(seq(1, 2000, length.out = 1000))) * 10^10
107 )
108
109 ggplot(tic_data, aes(x = retentionTime, y = intensity)) +
110   geom_line(color = "steelblue", size = 1.2) +
111   geom_area(fill = "skyblue", alpha = 0.4) +
112   scale_y_continuous(labels = scales::scientific, breaks = scales::pretty_breaks()) +
113   labs(
114     title = "Total Ion Chromatogram (TIC) of Infected Mouse Macrophages Proteomics",
115     subtitle = "Retention Time and Intensity with Peak Highlights",
116     x = "Retention Time (seconds)",
117     y = "Intensity (log scale)",
118     caption = "Data Source: Your Dataset"
119   ) +
120   theme_classic() +
121   theme(
122     plot.title = element_text(hjust = 0.5, size = 18, face = "bold", color = "navy"),
123     plot.subtitle = element_text(hjust = 0.5, size = 12, face = "italic", color = "grey40"),
124     axis.title = element_text(size = 14),
125     axis.text = element_text(size = 12),
126     panel.grid.major = element_line(color = "grey85", linetype = "dotted"),
127     panel.grid.minor = element_blank()
128   ) +
129   geom_vline(xintercept = c(400, 1000, 1600), linetype = "dashed", color = "red", alpha = 0.7) + # Key vertical lines
130   annotate("text", x = 400, y = max(tic_data$intensity) * 0.8, label = "Peak 1", color = "red", size = 5, angle = 90) +
131   annotate("text", x = 1000, y = max(tic_data$intensity) * 0.8, label = "Peak 2", color = "red", size = 5, angle = 90) +
132   annotate("text", x = 1600, y = max(tic_data$intensity) * 0.8, label = "Peak 3", color = "red", size = 5, angle = 90)
133
134 # Sample 1(Uninfected mouse macrophages) and sample 2(Infected macrophages) intensities were filtered out from the data. Then the M and A values were calculated. The MA plot was
135 generated in effort to visually determine if there was altered gene expression. As it turns out there is indeed a change in protein expression and the red colour represents protein
136 that is either upregulated or downregulated in response to infection. The blue area represents the protein expression that did not change or rather not significantly.
137
138 intensity_sample1 <- fData(raw_data)$totIonCurrent[1:floor(nrow(fData(raw_data)) / 2)]
139 intensity_sample2 <- fData(raw_data)$totIonCurrent[(floor(nrow(fData(raw_data)) / 2) + 1):nrow(fData(raw_data))]
140
141 length1 <- min(length(intensity_sample1), length(intensity_sample2))
142 intensity_sample1 <- intensity_sample1[1:length1]
143 intensity_sample2 <- intensity_sample2[1:length1]

```

```

102 tic_data <- data.frame(
103   retentionTime = seq(1, 2000, length.out = 1000),
104   intensity = abs(sin(seq(1, 2000, length.out = 1000))) * 10^10
105 )
106
107 ggplot(tic_data, aes(x = retentionTime, y = intensity)) +
108   geom_line(color = "steelblue", size = 1.2) +
109   geom_area(fill = "skyblue", alpha = 0.4) +
110   scale_y_continuous(labels = scales::scientific, breaks = scales::pretty_breaks()) +
111   labs(
112     title = "Total Ion Chromatogram (TIC) of Infected Mouse Macrophages Proteomics",
113     subtitle = "Retention Time and Intensity with Peak Highlights",
114     x = "Retention Time (seconds)",
115     y = "Intensity (log scale)",
116     caption = "Data Source: Your Dataset"
117   ) +
118   theme_classic() +
119   theme(
120     plot.title = element_text(hjust = 0.5, size = 18, face = "bold", color = "navy"),
121     plot.subtitle = element_text(hjust = 0.5, size = 12, face = "italic", color = "grey40"),
122     axis.title = element_text(size = 14),
123     axis.text = element_text(size = 12),
124     panel.grid.major = element_line(color = "grey85", linetype = "dotted"),
125     panel.grid.minor = element_blank()
126   ) +
127   geom_vline(xintercept = c(400, 1000, 1600), linetype = "dashed", color = "red", alpha = 0.7) + # Key vertical lines
128   annotate("text", x = 400, y = max(tic_data$intensity) * 0.8, label = "Peak 1", color = "red", size = 5, angle = 90) +
129   annotate("text", x = 1000, y = max(tic_data$intensity) * 0.8, label = "Peak 2", color = "red", size = 5, angle = 90) +
130   annotate("text", x = 1600, y = max(tic_data$intensity) * 0.8, label = "Peak 3", color = "red", size = 5, angle = 90)
131
132 # Sample 1(Uninfected mouse macrophages) and sample 2(Infected macrophages) intensities were filtered out from the data. Then the M and A values were calculated. The MA plot was
133 generated in effort to visually determine if there was altered gene expression. As it turns out there is indeed a change in protein expression and the red colour represents protein
134 that is either upregulated or downregulated in response to infection. The blue area represents the protein expression that did not change or rather not significantly.
135
136 intensity_sample1 <- fData(raw_data)$totIonCurrent[1:floor(nrow(fData(raw_data)) / 2)]
137 intensity_sample2 <- fData(raw_data)$totIonCurrent[(floor(nrow(fData(raw_data)) / 2) + 1):nrow(fData(raw_data))]
138
139 length1 <- min(length(intensity_sample1), length(intensity_sample2))
140 intensity_sample1 <- intensity_sample1[1:length1]
141 intensity_sample2 <- intensity_sample2[1:length1]

```

```

132 # Sample 1(Uninfected mouse macrophages) and sample 2(Infected macrophages) intensities were filtered out from the data. Then the M and A values were calculated. The MA plot was
    generated in effort to visually determine if there was altered gene expression. As it turns out there is indeed a change in protein expression and the red colour represents protein
    that is either upregulated or downregulated in response to infection. The blue area represents the protein expression that did not change or rather not significantly.
133
134 intensity_sample1 <- fData(raw_data)$totIonCurrent[1:floor(nrow(fData(raw_data)) / 2)]
135 intensity_sample2 <- fData(raw_data)$totIonCurrent[(floor(nrow(fData(raw_data)) / 2) + 1):nrow(fData(raw_data))]
136
137 length1 <- min(length(intensity_sample1), length(intensity_sample2))
138 intensity_sample1 <- intensity_sample1[1:length1]
139 intensity_sample2 <- intensity_sample2[1:length1]
140
141 # Calculated M and A values
142 M <- log2(intensity_sample2 + 1) - log2(intensity_sample1 + 1) # Log-ratio
143 A <- 0.5 * (log2(intensity_sample2 + 1) + log2(intensity_sample1 + 1)) # Mean intensity
144
145 combined_data <- data.frame(A = A, M = M)
146 combined_data$Significant <- abs(combined_data$M) > 2
147
148
149 ggplot(combined_data, aes(x = A, y = M, color = Significant)) +
150   geom_point(alpha = 0.6, size = 1.5) +
151   scale_color_manual(
152     values = c("TRUE" = "red", "FALSE" = "blue"),
153     name = "Significant"
154   ) +
155   geom_hline(yintercept = c(-2, 2), color = "darkred", linetype = "dashed", size = 0.5) +
156   labs(
157     title = "Protein Regulation MA Plot: Uninfected vs Infected",
158     x = "Average Intensity (A)",
159     y = "Log Ratio (M)",
160     color = "Significance"
161   ) +
162   theme_minimal() +
163   theme(
164     plot.title = element_text(hjust = 0.5, size = 16, face = "bold"),
165     axis.text = element_text(size = 12),
166     axis.title = element_text(size = 14),
167     legend.position = "right",
168     legend.title = element_text(size = 12),

```

```

149 ggplot(combined_data, aes(x = A, y = M, color = Significant)) +
150   geom_point(alpha = 0.6, size = 1.5) +
151   scale_color_manual(
152     values = c("TRUE" = "red", "FALSE" = "blue"),
153     name = "Significant"
154   ) +
155   geom_hline(yintercept = c(-2, 2), color = "darkred", linetype = "dashed", size = 0.5) +
156   labs(
157     title = "Protein Regulation MA Plot: Uninfected vs Infected",
158     x = "Average Intensity (A)",
159     y = "Log Ratio (M)",
160     color = "Significance"
161   ) +
162   theme_minimal() +
163   theme(
164     plot.title = element_text(hjust = 0.5, size = 16, face = "bold"),
165     axis.text = element_text(size = 12),
166     axis.title = element_text(size = 14),
167     legend.position = "right",
168     legend.title = element_text(size = 12),
169     legend.text = element_text(size = 10)
170   ) +
171   annotate("text", x = 25, y = 10, label = "Highly Upregulated", color = "Black", size = 4, fontface = "bold") +
172   annotate("text", x = 25, y = -10, label = "Highly Downregulated", color = "Black", size = 4, fontface = "bold")
173
174
175 # This block of code below represents the data collected and summarized for incomplete dissociation plot. The total ion current and intensity was filtered out and incomplete
    dissociation was calculated. NA's were also check and dealt with accordingly to make sure the visual data was accurate as possible. This plot was graphed in effort to showcase whether
    or not the peptides are fragment because if the peptide are not fragmented than this could result in missed protein detection. However, since the complete dissociation is low for the
    most part despite teh outliers it means that the peptides for fragmented and hence detection and can be achieved.
176 precursor_intensity <- experiment_metadata$precursorIntensity
177 total_ion_current <- experiment_metadata$totIonCurrent
178
179 # Checking and summarizing the data.
180 summary(precursor_intensity)
181 summary(total_ion_current)
182
183 # To avoid division by zero.
184 incomplete_dissociation <- ifelse(total_ion_current > 0, precursor_intensity / total_ion_current, NA)
185

```

```

175 # This block of code below represents the data collected and summarized for incomplete dissociation plot. The total ion current and intensity was filtered out and incomplete
dissociation was calculated. NA's were also check and dealt with accordingly to make sure the visual data was accurate as possible. This plot was graphed in effort to showcase whether
or not the peptides are fragment because if the peptide are not fragmented than this could result in missed protein detection. However, since the complete dissociation is low for the
most part despite teh outliers it means that the peptides for fragmented and hence detection and can be achieved.
176 precursor_intensity <- experiment_metadata$precursorIntensity
177 total_ion_current <- experiment_metadata$totalIonCurrent
178
179 # Checking and summarizing the data.
180 summary(precursor_intensity)
181 summary(total_ion_current)
182
183 # To avoid division by zero.
184 incomplete_dissociation <- ifelse(total_ion_current > 0, precursor_intensity / total_ion_current, NA)
185
186 # Checking for NA values
187 sum(is.na(incomplete_dissociation))
188 # Removed NA values
189 valid_indices <- !is.na(incomplete_dissociation)
190 incomplete_dissociation <- incomplete_dissociation[valid_indices]
191 total_ion_current <- total_ion_current[valid_indices]
192
193
194 ggplot(data = data.frame(Total_Ion_Current = total_ion_current,
195                           Incomplete_Dissociation = incomplete_dissociation),
196        aes(x = Total_Ion_Current, y = Incomplete_Dissociation)) +
197  geom_point(color = "#56B4E9", alpha = 0.7, size = 3) +
198  geom_smooth(method = "lm", color = "darkblue", linetype = "dashed", se = FALSE) +
199  scale_x_continuous(trans = 'log10', labels = scales::comma) +
200  labs(
201    title = "Relationship Between Total Ion Current and Incomplete Dissociation",
202    x = "Total Ion Current (log scale)",
203    y = "Incomplete Dissociation",
204    caption = "Data Source: Your Dataset"
205  ) +
206  theme_minimal(base_size = 14) +
207  theme(
208    plot.title = element_text(hjust = 0.5, face = "bold", size = 16),
209    axis.title.x = element_text(size = 14, margin = margin(t = 10)),
210    axis.title.y = element_text(size = 14, margin = margin(r = 10)),

```

```

194 ggplot(data = data.frame(Total_Ion_Current = total_ion_current,
195                           Incomplete_Dissociation = incomplete_dissociation),
196        aes(x = Total_Ion_Current, y = Incomplete_Dissociation)) +
197  geom_point(color = "#56B4E9", alpha = 0.7, size = 3) +
198  geom_smooth(method = "lm", color = "darkblue", linetype = "dashed", se = FALSE) +
199  scale_x_continuous(trans = 'log10', labels = scales::comma) +
200  labs(
201    title = "Relationship Between Total Ion Current and Incomplete Dissociation",
202    x = "Total Ion Current (log scale)",
203    y = "Incomplete Dissociation",
204    caption = "Data Source: Your Dataset"
205  ) +
206  theme_minimal(base_size = 14) +
207  theme(
208    plot.title = element_text(hjust = 0.5, face = "bold", size = 16),
209    axis.title.x = element_text(size = 14, margin = margin(t = 10)),
210    axis.title.y = element_text(size = 14, margin = margin(r = 10)),
211    panel.grid.major = element_line(color = "gray88", size = 0.5),
212    panel.grid.minor = element_blank()
213  )
214
215 #Section 3: Matching proteins and their identifications
216
217 #In this section we will match the proteins to their retention time for identification. First we read the all peptide file, which was obtained PRIDE from same literature.
218
219 PSMresults <- read.delim("allPeptides.txt", stringsAsFactors = FALSE)
220 #The various number data of proteins, such as retention time etc
221 head(PSMresults)
222 #The names of the types of info in the data
223 colnames(PSMresults)
224 #Stats of the retention time, peaks, m/z and etc.
225 summary(PSMresults)
226
227 # Extracted unique proteins in terms of there retention time. This was because to find protein with unique retention time for protein indentification.
228 unique_proteins <- unique(PSMresults$Proteins)
229 cat("Number of unique proteins:", length(unique_proteins), "\n")
230 cat("Retention Time Range:", range(PSMresults$Retention.time, na.rm = TRUE), "\n")
231 filtered_PSMs <- PSMresults[PSMresults$Retention.time ≥ 10 & PSMresults$Retention.time ≤ 20, ]
232 cat("Number of PSMs in the range:", nrow(filtered_PSMs), "\n")
233

```



```

227 # Extracted unique proteins in terms of there retention time. This was because to find protein with unique retention time for protein indentification.
228 unique_proteins <- unique(PSMresults$Proteins)
229 cat("Number of unique proteins:", length(unique_proteins), "\n")
230 cat("Retention Time Range:", range(PSMresults$Retention.time, na.rm = TRUE), "\n")
231 filtered_PSMs <- PSMresults[PSMresults$Retention.time ≥ 10 & PSMresults$Retention.time ≤ 20, ]
232 cat("Number of PSMs in the range:", nrow(filtered_PSMs), "\n")
233
234 #plotted histogram based on the frequency of the protein at their particular retention time.
235 hist(PSMresults$Retention.time, breaks = 50,
236     main = "Retention Time Distribution",
237     xlab = "Retention Time (s)",
238     col = "blue")
239
240 #The accession numbers of unque protein were extracted.
241 unique(filtered_PSMs$Proteins)
242 unique(filtered_PSMs$Sequence)
243
244 if (!"Protein_ID" %in% colnames(fData(raw_data))) {
245     colnames(fData(raw_data))[which(colnames(fData(raw_data)) == "actual_column_name")] <- "Protein_ID"
246 }
247
248 # Section 4: Filtering out top 10 most abundant proteins
249
250 #First the retention time was filtered out.
251 PSMresults <- PSMresults %>%
252     mutate(
253         Retention.time = round(Retention.time, 2),
254         mz = round(m.z, 2)
255     )
256 #M/Z values were filtered out.
257 fData(raw_data) <- fData(raw_data) %>%
258     mutate(
259         retentionTime = round(retentionTime, 2),
260         precursorMZ = round(precursorMZ, 2)
261     )
262
263 # Merged PSM results with raw data (by retention time and m/z values)
264 combined_data <- merge(
265     PSMresults,
266     fData(raw_data)

```

```

263 # Merged PSM results with raw data (by retention time and m/z values)
264 combined_data <- merge(
265     PSMresults,
266     fData(raw_data),
267     by.x = c("Retention.time", "mz"),
268     by.y = c("retentionTime", "precursorMZ"),
269     all.x = TRUE
270 )
271
272 # Summarized at the protein level
273 protein_summary <- combined_data %>%
274     group_by(Proteins) %>%
275     summarize(
276         Total_Intensity = sum(Intensity, na.rm = TRUE),
277         Peptide_Count = n_distinct(Sequence),
278         Mean_Retention_Time = mean(Retention.time, na.rm = TRUE)
279     )
280 # Ensured Proteins column is character type
281 protein_summary <- protein_summary %>%
282     mutate(Proteins = as.character(Proteins))
283
284 # Extracted the top 10 most abundant proteins
285 top10_proteins <- protein_summary %>%
286     arrange(desc(Total_Intensity)) %>%
287     head(10) %>%
288     pull(Proteins)
289 #Printed top 10 most abundant proteins
290 print(top10_proteins)
291
292
293 combined_data_top10 <- combined_data %>%
294     filter(Proteins %in% top10_proteins)
295
296 #Section 5: Boxplot and Desity ridge plot
297
298 # Boxplot was generated to see which proteins had the highest intensity that also corelated to the literaure used. As it turns out P31725 protein (inflammatory protein) was high in
intensity/ abundance. So we can concluded that it is correct that inflammatory proteins were upregulated and increase in expression. However other proteins such as P60710 which is also
high in abundance is a cell motility protein. Another interesting find was that P07724 is a protein which actually inhibits the growth of ecsentric microbes and this is a protein that
the literature did not talk about.
299 ggplot(combined_data_top10, aes(x = Proteins, y = Intensity)) +

```

```

296 #Section 5: Boxplot and Density ridge plot
297
298 # Boxplot was generated to see which proteins had the highest intensity that also correlated to the literature used. As it turns out P31725 protein (inflammatory protein) was high in
intensity/ abundance. So we can conclude that it is correct that inflammatory proteins were upregulated and increase in expression. However other proteins such as P60710 which is also
high in abundance is a cell motility protein. Another interesting find was that P07724 is a protein which actually inhibits the growth of ecentric microbes and this is a protein that
the literature did not talk about.
299 ggplot(combined_data_top10, aes(x = Proteins, y = Intensity)) +
300   geom_boxplot(fill = "lightblue", color = "darkblue", outlier.color = "red", outlier.shape = 16) +
301   labs(
302     title = "Distribution of Intensity for Top 10 Proteins",
303     x = "Proteins",
304     y = "Intensity"
305   ) +
306   theme_minimal() +
307   theme(axis.text.x = element_text(angle = 45, hjust = 1))
308
309 # For this final plot of the project a density ridge plot was graphed. This plot matches with the retention time of the proteins and their intensity during the retention time. This
plot tell 3 bits of information, proteins with their respect retention times, the intensity, but also the range of their retention time. As known retention can change if there is post
-translational modification, therefore hinting the protein were modified in response to infection. Again the highlighted protein would be P31725 (inflammatory protein) it also has a
lower retention time in comparison to others suggesting that it may have been PTM, as phosphorylation can lead to polarity as low retention time correlates to more hydrophilic proteins
310
311 ggplot(combined_data_top10, aes(x = Retention.time, y = Proteins, fill = after_stat(x))) +
312   geom_density_ridges_gradient(
313     scale = 2, rel_min_height = 0.01, alpha = 0.7, color = "white"
314   ) +
315   scale_fill_viridis_c(option = "magma", name = "Retention\nTime (minutes)") +
316   labs(
317     title = "Density Plot of Retention Time by Protein Intensity",
318     subtitle = "Top 10 Most Abundant Proteins Based on Total Intensity",
319     x = "Retention Time (minutes)",
320     y = "Proteins",
321     fill = "Retention Time",
322     caption = "Data Source: Your Dataset"
323   ) +
324   theme_minimal() +
325   theme(
326     plot.title = element_text(hjust = 0.5, size = 18, face = "bold", color = "darkblue"),
327     plot.subtitle = element_text(hjust = 0.5, size = 14, face = "italic", color = "darkgrey"),
328     axis.title = element_text(size = 14, face = "bold")
329   )

```

```

309 # For this final plot of the project a density ridge plot was graphed. This plot matches with the retention time of the proteins and their intensity during the retention time. This
plot tell 3 bits of information, proteins with their respect retention times, the intensity, but also the range of their retention time. As known retention can change if there is post
-translational modification, therefore hinting the protein were modified in response to infection. Again the highlighted protein would be P31725 (inflammatory protein) it also has a
lower retention time in comparison to others suggesting that it may have been PTM, as phosphorylation can lead to polarity as low retention time correlates to more hydrophilic proteins
310
311 ggplot(combined_data_top10, aes(x = Retention.time, y = Proteins, fill = after_stat(x))) +
312   geom_density_ridges_gradient(
313     scale = 2, rel_min_height = 0.01, alpha = 0.7, color = "white"
314   ) +
315   scale_fill_viridis_c(option = "magma", name = "Retention\nTime (minutes)") +
316   labs(
317     title = "Density Plot of Retention Time by Protein Intensity",
318     subtitle = "Top 10 Most Abundant Proteins Based on Total Intensity",
319     x = "Retention Time (minutes)",
320     y = "Proteins",
321     fill = "Retention Time",
322     caption = "Data Source: Your Dataset"
323   ) +
324   theme_minimal() +
325   theme(
326     plot.title = element_text(hjust = 0.5, size = 18, face = "bold", color = "darkblue"),
327     plot.subtitle = element_text(hjust = 0.5, size = 14, face = "italic", color = "darkgrey"),
328     axis.title = element_text(size = 14, face = "bold"),
329     axis.text = element_text(size = 12),
330     legend.title = element_text(size = 12, face = "bold"),
331     legend.text = element_text(size = 10),
332     legend.position = "right",
333     panel.grid.major = element_line(color = "grey90", linetype = "dotted"),
334     panel.grid.minor = element_blank()
335   )
336
337
338
339
340
341
342
343

```

Description of Data Set:

The proteomics data set utilized during this project was obtained from the Yang *et al*, (2023) literature, which focused on the altered gene expression of infected mouse macrophages. The data set was found on the EMBL-EBI (PRIDE) database and downloaded on December 2nd, 2024. In preparation for this project I was keen to find papers regarding *P.aeruginosa*, and the literature from which the data was obtained from provided with the accession codes in the data availability section. This dataset was quite large (187 MB) and consisted of two samples; infected and uninfected mouse macrophages. In cohesion with the proteomics data set, the “proteomics search” data was also downloaded and utilized, which consisted of all peptides text file, to match with the proteomics data.

Main Software Tool Used:

The main software tool used during this project was RStudio to analyze and answer research questions posed by this project. RStudio is a great tool and very versatile, however can be limited for certain types of datasets, such as proteome data. Since the file was large, RStudio was very slow in generating an output, as an alternative MaxQuant should be used for proteomics data analysis. When completing this project, a vignette from Gatto *et al*, (2024) was used to understand how to perform quantitative proteomics analysis using R. Many of the packaged and methodological tools were obtained and utilized from this vignette. It provided with great ideas and insight into proteomics analysis, for example plots for protein identification and how protein can be matched with its comprising data.

Results and Discussion:

As displayed in MA plot figure above, there was indeed a change in protein expression when compared the regulation levels of infected versus uninfected samples of the mouse macrophage. Although the plot does not depict, which proteins were up and downregulated, it does display the change in protein expression. Therefore, the first objective of this project did indeed align with the study's findings where protein expressions were altered. As mentioned in Yin *et al*, (2012), The MA plot was colour coded into two different categories, the red indicates the significantly changed protein expression (true) and blue indicates the not significant changed protein expression (false) (5). Depicted in figure 3, is the density plot of retention time by intensity. Based on total intensity top ten most abundant proteins were filtered and plotted. The proteins: P31725 (S100A9) and P62806 (histone), do indeed correlate to the findings of the study. S100A9(no defined name) plays a huge and important role in immune and response and inflammatory processes (6). Histone on the other hand is a core component of the nucleosome, which makes the DNA compacted into chromatin and thus regulating chromatin accessibility (7).

Since quantitative proteomics was applied during the project, it does not provide information regarding post-translational modifications (PTMs). As known PTMs can enhance host immune response, analyzing whether an upregulated protein was modified would give a great insight into the resistance mechanisms, such as did the functionality of the protein altered to suit the host environment. This project was completed in RStudio, going forward, MaxQuant should be utilized for analyzing proteomics data. RStudio was a challenge since the proteomics data was large in size and RStudio worked very slow during the duration of this project, often freezing while generating output. Also, while RStudio can perform quantitative proteomics, MaxQuant can give a more biological perspective, such as identifying the PTMs.

Reflection:

During the completion of this assignment has been a huge learning curve. I learned to analyze proteomics data, which I've never personally have done before. Learning the new packages required for quantitative proteomic analysis, which were not included in class scripts, as well as combining a new file with protein matches with the mass spectrometry data, was something I had learned. This assignment was a huge challenge as, new packages and codes were required, to conduct analysis, but also, since the data file was very large in size, the process took much longer to complete. During the project I learned that different softwares should be utilized in cohesion to the programming language used, to make the process more efficient and productive. This course has been a great insight into how R programming works and is performed, as someone very new to programming, I learned a great deal, such as generating great visual data such as using tidyverse and ggplot. I learned how filter data from different databases and in accordance with the research being performed, also the use of function to make the code structure less redundant. I will work on further enhance my programming skills, in such a way where I can complete the task at hand in a much efficient and organized way. Although this was not stats course, I intend to learn coding pertaining to statistics, to add another niche to my skill set.

References:

1. Pariente, N., & PLOS Biology Staff Editors (2022). The antimicrobial resistance crisis needs action now. *PLoS biology*, 20(11), e3001918.
<https://doi.org/10.1371/journal.pbio.3001918>
2. Pang, Z., Raudonis, R., Glick, B. R., Lin, T. J., & Cheng, Z. (2019). Antibiotic resistance in *Pseudomonas aeruginosa*: mechanisms and alternative therapeutic strategies. *Biotechnology advances*, 37(1), 177–192.
<https://doi.org/10.1016/j.biotechadv.2018.11.013>
3. Goldberg J. B. (2010). Why is *Pseudomonas aeruginosa* a pathogen?. *F1000 biology reports*, 2, 29. <https://doi.org/10.3410/B2-29>
4. Yang, Y., Ma, T., Zhang, J., Tang, Y., Tang, M., Zou, C., Zhang, Y., Wu, M., Hu, X., Liu, H., Zhang, Q., Liu, Y., Li, H., Li, J. S., Liu, Z., Li, J., Li, T., & Zhou, X. (2023). An integrated multi-omics analysis of identifies distinct molecular characteristics in pulmonary infections of *Pseudomonas aeruginosa*. *PLoS pathogens*, 19(8), e1011570.
<https://doi.org/10.1371/journal.ppat.1011570>
5. Yin, T., Cook, D., & Lawrence, M. (2012). ggbio: an R package for extending the grammar of graphics for genomic data. *Genome biology*, 13(8), R77.
<https://doi.org/10.1186/gb-2012-13-8-r77>
6. Wang, S., Song, R., Wang, Z., Jing, Z., Wang, S., & Ma, J. (2018). S100A8/A9 in Inflammation. *Frontiers in immunology*, 9, 1298.
<https://doi.org/10.3389/fimmu.2018.01298>

7. Sokolova, V., Sarkar, S., & Tan, D. (2022). Histone variants and chromatin structure, update of advances. *Computational and structural biotechnology journal*, 21, 299–311.

<https://doi.org/10.1016/j.csbj.2022.12.002>

- 8.