```r
setwd("~/04_Champion_Lab/02_N-terminal_Acetylation/OnePot")
library(readxl)
library(stringr)
library(tidyr)
library(dplyr)
library(seqinr)
library(Peptides)
library(cleaver)

#One pot import and cleaning
opDF <- read.csv("Data/db.protein-peptides.csv")

opDF <- opDF[!grepl("CONTAM", opDF$Accession),]
opDF$Peptide <- str_remove(opDF$Peptide, "^.\\.")
opDF$Peptide <- str_remove(opDF$Peptide, "\\..$")


opDF$NtermAcet_L <- grepl("^[A-Z]\\(\\+42\\.01\\)", opDF$Peptide)
opDF$NtermAcet_H <- grepl("^[A-Z]\\(\\+45\\.03\\)", opDF$Peptide)
opDF$UnLabelled <- !opDF$NtermAcet_H & !opDF$NtermAcet_L

names(opDF) <- str_replace_all(names(opDF), "Intensity.New", "Intensity_Sample")
names(opDF) <- str_replace_all(names(opDF), "Area.New", "Area_Sample")
names(opDF) <- str_replace_all(names(opDF), "X.Spec.New", "nSpec_Sample")


opDF <- opDF  %>% pivot_longer(which((grepl("^Area", names(opDF)) |
                                      grepl("^Intensity", names(opDF)) |
                                      grepl("^nSpec", names(opDF)))),
                                names_to = c(".value", "measure"),
                                names_sep = "_")
drops <- c("Unique","Top","scan", "Source File", "AScore")
opDF <- opDF[,!names(opDF) %in% drops]


opDF$Accession2 <- str_remove(str_extract(opDF$Accession, "^[^\\|]*\\|"), "\\|")


opDFsummary <- opDF %>% group_by(Peptide, Mass, m.z, z,
                                  RT, Start, End, Length) %>%
  summarise(Accessions = paste(unique(Accession2), collapse = ";"),
            rowsN = n(),
            valuesArea = sum(!is.na(Area)),
            meanArea = mean(Area, na.rm = T),
            sdArea = sd(Area, na.rm = T),
            maxArea = max(Area, na.rm = T),
            valuesIntensity = sum(!is.na(Intensity)),
            meanIntensity = mean(Intensity, na.rm = T),
            sdIntensity = sd(Intensity, na.rm = T),
            maxIntensity = max(Intensity, na.rm = T),
            valuesSpec = sum(nSpec > 0))


opDFsummary$str_seq <- str_remove_all(opDFsummary$Peptide,
"\\([\\+\\-][0-9]+\\.[0-9]+\\)")

opDFsummary$IEP <- NA
for (i in 1:nrow(opDFsummary)) {
  opDFsummary$IEP[i] <- computePI(s2c(opDFsummary$str_seq[i]))
  if (i %% 500 == 0) {
    print(i/ nrow(opDFsummary) *100)
  }
}
#save as dataframe
write.csv(opDFsummary, "Data/OnePotSummarisedData.csv",
          row.names = F, quote = F)


#enriched data (Thompson et al.) import and cleaning (trypsin only)
```

```r
69    enrDF <- read_excel("Data/Cristal Marinum.xlsx", sheet = "Raw NTA Peptides")
70
71    enrDF <- enrDF[!grepl("GluC", enrDF$File),]
72
73    enrDF$Peptide <- paste0(enrDF$Sequence, "|", enrDF$Modifications)
74
75    enrDFsummary <- enrDF %>% group_by(Peptide, Sequence, `Theor m/z`, `Theor z`,
76                                        `Theor MW`, Length) %>%
77      summarise(Genes = paste(unique(Gene), collapse = ";"),
78                rowsN = n(),
79                meanIntensity = mean(`Intensity (Peptide)`, na.rm = T),
80                maxIntensity = max(`Intensity (Peptide)`, na.rm = T),
81                valuesIntensity = sum(!is.na(`Intensity (Peptide)`)))
82
83    #function to return start position of peptide from proteome db
84    return_start_position <- function(peptide, accession, database){
85      protein_names <- names(database)
86      protein_index <- which(grepl(accession, protein_names))
87      if(length(protein_index) > 1) {
88        return(NA)
89      } else if (length(protein_index) < 1) {
90        return(NA)
91      } else {
92        protein <- toupper(c2s(database[[protein_index[1]]]))
93      }
94      matches_df <- str_locate(protein, peptide)
95      if (length(matches_df) < 1) {
96        return(NA)
97      } else {
98        start_position <- matches_df[[1,1]]
99      }
100     return(start_position)
101   }
102
103   #read in database
104   db <- read.fasta("Data/SDW_codon_MARINUM_MetStart.fasta", seqtype = "AA")
105
106
107   enrDFsummary$IEP <- NA
108   enrDFsummary$Start <- NA
109   for (i in 1:nrow(enrDFsummary)) {
110     enrDFsummary$IEP[i] <- computePI(s2c(enrDFsummary$Sequence[i]))
111     enrDFsummary$Start[i] <- return_start_position(enrDFsummary$Sequence[i],
112                                                    enrDFsummary$Genes[i],
113                                                    database = db)
114     if (i %% 500 == 0) {
115       print(i/ nrow(enrDFsummary) *100)
116     }
117   }
118
119   #saving
120   write.csv(enrDFsummary, "Data/EnrichedSummarisedData.csv",
121             row.names = F, quote = F)
122
123
124   #same but with GluC
125   enrDF <- read_excel("Data/Cristal Marinum.xlsx", sheet = "Raw NTA Peptides")
126
127   enrDF <- enrDF[grepl("GluC", enrDF$File),]
128
129   enrDF$Peptide <- paste0(enrDF$Sequence, "|", enrDF$Modifications)
130
131   enrDFsummary <- enrDF %>% group_by(Peptide, Sequence, `Theor m/z`, `Theor z`,
132                                       `Theor MW`, Length) %>%
133     summarise(Genes = paste(unique(Gene), collapse = ";"),
134               rowsN = n(),
135               meanIntensity = mean(`Intensity (Peptide)`, na.rm = T),
136               maxIntensity = max(`Intensity (Peptide)`, na.rm = T),
137               valuesIntensity = sum(!is.na(`Intensity (Peptide)`)))
```

```r
return_start_position <- function(peptide, accession, database){
  protein_names <- names(database)
  protein_index <- which(grepl(accession, protein_names))
  if(length(protein_index) > 1) {
    return(NA)
  } else if (length(protein_index) < 1) {
    return(NA)
  } else {
    protein <- toupper(c2s(database[[protein_index[1]]]))
  }
  matches_df <- str_locate(protein, peptide)
  if (length(matches_df) < 1) {
    return(NA)
  } else {
    start_position <- matches_df[[1,1]]
  }
  return(start_position)
}
db <- read.fasta("Data/SDW_codon_MARINUM_MetStart.fasta", seqtype = "AA")


enrDFsummary$IEP <- NA
enrDFsummary$Start <- NA
for (i in 1:nrow(enrDFsummary)) {
  enrDFsummary$IEP[i] <- computePI(s2c(enrDFsummary$Sequence[i]))
  enrDFsummary$Start[i] <- return_start_position(enrDFsummary$Sequence[i],
                                                 enrDFsummary$Genes[i],
                                                 database = db)

  if (i %% 500 == 0) {
    print(i/ nrow(enrDFsummary) *100)
  }
}


write.csv(enrDFsummary, "Data/EnrichedSummarisedDataGLUC.csv",
          row.names = F, quote = F)




#list of Nterminal peptides from genome
NtermPeptides <- list()

for (i in 1:length(db)) {
  protein <- db[[i]]
  seq <- toupper(c2s(protein))
  MMAR <- str_remove(str_extract(attr(protein, 'name'), "^[^\\|]+|"), "\\|")
  peptide <- cleave(seq, enzym = "trypsin")[[1]][1]
  NtermPeptides[[MMAR]] <- peptide
}

NtermGenome <- data.frame(Accession = names(NtermPeptides),
                          Peptide = unlist(NtermPeptides))

NtermGenome$Length <- nchar(NtermGenome$Peptide)
NtermGenome$secondAA <- substr(NtermGenome$Peptide, 2,2)
NtermGenome$metCleavedPeptide <- str_remove(NtermGenome$Peptide, "^M")


NtermGenome$IEP <- NA
for (i in 1:nrow(NtermGenome)) {
  if (NtermGenome$secondAA[i] %in% c("G", "A", "S", "T", "V", "C", "P")) {
    NtermGenome$IEP[i] <- computePI(s2c(NtermGenome$metCleavedPeptide[i]))
  } else {
    NtermGenome$IEP[i] <- computePI(s2c(NtermGenome$Peptide[i]))
  }
```

```r
207        if (i %% 500 == 0) {
208          print(i/ nrow(NtermGenome) *100)
209        }
210      }
211
212
213      write.csv(NtermGenome, "Data/GenomeSummarisedData_GASTVCP.csv", row.names = F, quote = F)
214
215
216
217
218
219
220      #all genome peptides (not just n terminal)
221
222
223      proteinDFList <- list()
224
225      #filters
226      minLength <- 5
227      maxLength <- 50
228
229
230
231      for (i in 1:length(db)) {
232        if (i %% 100 == 0) {
233          print(i / length(db))
234        }
235        protein <- db[[i]]
236        seq <- toupper(c2s(protein))
237        seq <- str_remove(seq, "\\*")
238        MMAR <- str_remove(str_extract(attr(protein, 'name'), "^[^\\|]+|"), "\\|")
239        peptides <- cleave(seq, enzym = "trypsin")[[1]]
240        peptideDF <- data.frame(Accession = MMAR,
241                                Peptide = peptides)
242        peptideDF$NtermPeptide <- c(T, rep(F, nrow(peptideDF) - 1))
243        NtermPeptide <- peptideDF$Peptide[1]
244        secondAA <- substr(NtermPeptide, 2,2)
245        if (secondAA %in% c("G", "A", "S", "T", "V", "C", "P")) {
246          peptideDF$Peptide[1] <- str_remove(peptideDF$Peptide[1], "^M")
247        }
248        peptideDF$Length <- nchar(peptideDF$Peptide)
249
250        peptideDF <- peptideDF[peptideDF$Length >= minLength &
251                                 peptideDF$Length <= maxLength,]
252
253
254        if (nrow(peptideDF) > 0) {
255          peptideDF$IEP <- NA
256          peptideDF$mass <- NA
257          peptideDF$IEP <- sapply(peptideDF$Peptide, function(peptide) computePI(s2c(peptide)))
258          peptideDF$mass <- sapply(peptideDF$Peptide, mw)
259
260          proteinDFList[[MMAR]] <- peptideDF
261        } else {
262          print(i)
263          print("noPeps")
264        }
265
266      }
267
268      genomePeps <- bind_rows(proteinDFList)
269
270      write.csv(genomePeps, "Data/GenomeSummarisedData_allPeptides_GASTVCP.csv", row.names = F
271      , quote = F)
```