

PROJECT 4

Marie-Claire Indrilevitch, Anna Toidze, Sophia Houhamdi, Ieva Pudziuvelyte

Problem 8:
$$e_i(a) = \frac{t_i(a) + 1}{\sum_{a' \in A} (E_i(a') + 1)}$$

$$e_1(A) = \frac{4+1}{(4+1) + (0+1) + (0+1) + (0+1)} = 5/8$$

$$e_2(A) = \frac{0+1}{(0+1) + (0+1) + (0+1) + (0+1)} = 1/8 = e_2(C) = e_2(T)$$

$$e_1(C) = \frac{0+1}{(4+1) + 3(0+1)} = 1/8 \quad e_1(h) = 1/8 \quad e_1(T) = 1/8$$

$$e_2(h) = \frac{4+1}{8} = 5/8$$

$$e_3(A) = \frac{0+1}{(0+1) + (5+1) + (0+1) + (0+1)} = 1/9 = e_3(a) = e_3(T)$$

$$e_3(C) = \frac{5+1}{9} = 6/9$$

$$e_{\text{match}} = \begin{matrix} & \begin{matrix} 1 & 2 & 3 \end{matrix} \\ \begin{matrix} A \\ C \\ G \\ T \end{matrix} & \left\{ \begin{array}{ccc} 5/8 & 1/8 & 1/9 \\ 1/8 & 1/8 & 6/9 \\ 1/8 & 5/8 & 1/9 \\ 1/8 & 1/8 & 1/9 \end{array} \right\} \end{matrix}$$

Problem 9:
$$e_i(a) = \frac{t_i(a) + 1}{\sum_{a' \in A} (E_i(a') + 1)}$$

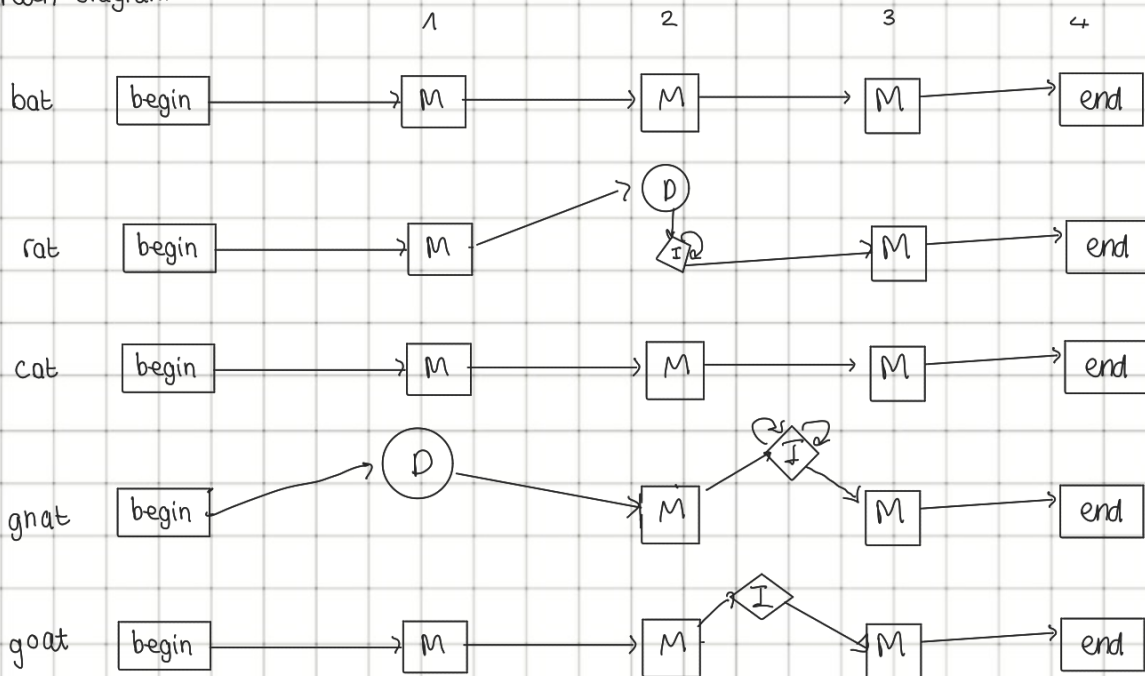
$$e_1(A) = \frac{0+1}{4 \cdot (0+1)} = 1/4 = e_1(C) = e_1(h) = e_1(T) = e_3(A) = e_3(C) = e_3(h) = e_3(T)$$

$$e_2(A) = \frac{5+1}{(5+1) + (1+1) + (0+1) + (0+1)} = 6/10 \quad e_2(h) = \frac{1+1}{10} = \frac{2}{10}$$

$$e_2(C) = e_2(T) = 1/10$$

$$e_{\text{insert}} = \begin{matrix} & \begin{matrix} 1 & 2 & 3 \end{matrix} \\ \begin{matrix} A \\ C \\ G \\ T \end{matrix} & \left\{ \begin{array}{ccc} 1/4 & 6/10 & 1/4 \\ 1/4 & 1/10 & 1/4 \\ 1/4 & 2/10 & 1/4 \\ 1/4 & 1/10 & 1/4 \end{array} \right\} \end{matrix}$$

Problem 3 Path diagrams:



$$t_i(k \rightarrow l) = \frac{T_i(k \rightarrow l) + 1}{\sum_{l' \in S} (T_i(k \rightarrow l') + 1)}$$

0→1:

$$t_0(M \rightarrow M) = \frac{4+1}{(4+1)+(1+1)+(0+1)} = 5/8$$

$$t_0(I \rightarrow M) = \frac{0+1}{3 \cdot (0+1)} = 1/3 = t_0(I \rightarrow I) = t_0(I \rightarrow D)$$

$$t_0(M \rightarrow I) = \frac{0+1}{8} = 1/8$$

$$t_0(D \rightarrow M) = t_0(D \rightarrow I) = t_0(D \rightarrow D) = 1/3$$

$$t_0(M \rightarrow D) = \frac{1+1}{8} = \frac{2}{8}$$

$$T_0 = \begin{matrix} & \begin{matrix} M & I & D \end{matrix} \\ \begin{matrix} M \\ I \\ D \end{matrix} & \begin{pmatrix} 5/8 & 1/8 & 2/8 \\ 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \end{pmatrix} \end{matrix}$$

1→2 $t_1(M \rightarrow M) = \frac{3+1}{(3+1)+(1+1)+(0+1)} = 4/7$

$$t_1(I \rightarrow M) = \frac{0+1}{3 \cdot (0+1)} = 1/3 = t_1(I \rightarrow I) = t_1(I \rightarrow D)$$

$$t_1(M \rightarrow I) = \frac{0+1}{7} = 1/7$$

$$t_1(D \rightarrow M) = \frac{1+1}{1+1+0+1+0+1} = 1/2$$

$$t_1(M \rightarrow D) = \frac{1+1}{7} = \frac{2}{7}$$

$$t_1(D \rightarrow D) = t_1(D \rightarrow I) = 1/4$$

$$T_1 = \begin{matrix} & \begin{matrix} M & I & D \end{matrix} \\ \begin{matrix} M \\ I \\ D \end{matrix} & \begin{pmatrix} 4/7 & 1/7 & 2/7 \\ 1/3 & 1/3 & 1/3 \\ 1/2 & 1/4 & 1/4 \end{pmatrix} \end{matrix}$$

2→3 $t_2(M \rightarrow M) = \frac{2+1}{(2+1)+(2+1)+(0+1)} = 3/7$

$$t_2(I \rightarrow M) = \frac{3+1}{3+1+3+1+0+1} = 4/9$$

$$t_2(M \rightarrow I) = \frac{2+1}{7} = 3/7$$

$$t_2(I \rightarrow I) = \frac{3+1}{9} = 4/9$$

$$t_2(M \rightarrow D) = \frac{0+1}{7} = 1/7$$

$$t_2(I \rightarrow D) = \frac{0+1}{9} = 1/9$$

$$t_2(D \rightarrow M) = \frac{0+1}{0+1+1+1+0+1} = 1/4$$

$$t_2(D \rightarrow I) = \frac{1+1}{4} = 1/2$$

$$t_2(D \rightarrow D) = \frac{0+1}{4} = 1/4$$

$$T_2 = \begin{matrix} & \begin{matrix} M & I & D \end{matrix} \\ \begin{matrix} M \\ I \\ D \end{matrix} & \begin{pmatrix} 3/7 & 3/7 & 1/7 \\ 4/9 & 4/9 & 1/9 \\ 1/4 & 1/2 & 1/4 \end{pmatrix} \end{matrix}$$

3→4

$$t_3(M \rightarrow M) = \frac{5+1}{(5+1)+(0+1)+(0+1)} = 6/8$$

$$t_3(I \rightarrow M) = \frac{0+1}{3 \cdot (0+1)} = 1/3 = t_3(I \rightarrow I) = t_3(I \rightarrow D)$$

$$t_3(M \rightarrow I) = \frac{0+1}{8} = 1/8$$

$$t_3(D \rightarrow M) = t_3(D \rightarrow I) = t_3(D \rightarrow D) = 1/3$$

$$t_3(M \rightarrow D) = \frac{0+1}{8} = \frac{1}{8}$$

$$T_3 = \begin{matrix} & & M & I & D \\ \begin{matrix} M \\ I \\ D \end{matrix} & \begin{pmatrix} 6/8 & 1/8 & 1/8 \\ 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \end{pmatrix} \end{matrix}$$

| | 0 | 1 | 2 | 3 |
|-----|-----|-----|-----|-----|
| M→M | 5/8 | 4/7 | 3/7 | 6/8 |
| M→I | 2/8 | 1/7 | 3/7 | 1/8 |
| M→D | 1/8 | 2/7 | 1/7 | 1/8 |
| I→M | 1/3 | 1/3 | 4/9 | 1/3 |
| I→I | 1/3 | 1/3 | 4/9 | 1/3 |
| I→D | 1/3 | 1/3 | 1/9 | 1/3 |
| D→M | 1/3 | 1/2 | 1/4 | 1/3 |
| D→I | 1/3 | 1/4 | 1/2 | 1/3 |
| D→D | 1/3 | 1/4 | 1/4 | 1/3 |

Project 4

Problem 11

You are given multiple alignments of protein sequences for two protein families: GTP binding proteins, and a family of ATPases. The task is to determine to which family certain unclassified proteins belong.

1. Run `source("profileHMM.R")` to import functions which you will use below

```
library(ggplot2)
source("profileHMM.R")
```

2. Read the two alignments 'GTP binding proteins.txt' and 'ATPases.txt' into memory using the function `parseAlignment()`.

```
GTP_alignment <- parseAlignment("C:/Users/marie/Desktop/CBB/SMCB/Project_4_student/data/GTP_binding_pro
ATP_alignment <- parseAlignment("C:/Users/marie/Desktop/CBB/SMCB/Project_4_student/data/ATPases.txt")
```

3. Use the function `learnHMM()` to parametrise two profile HMMs: one for each protein family (multiple alignment).

```
GTP_HMM <- learnHMM(GTP_alignment)
ATP_HMM <- learnHMM(ATP_alignment)
```

4. Identify the position(s) with the highest match and with the highest insert emission frequencies over all symbols. Plot the respective match and insert emission frequencies for the identified positions.

```
max_emission_GTP <- apply(GTP_HMM$mE, 2, max)
max_pos_emi_GTP <- which.max(max_emission_GTP)

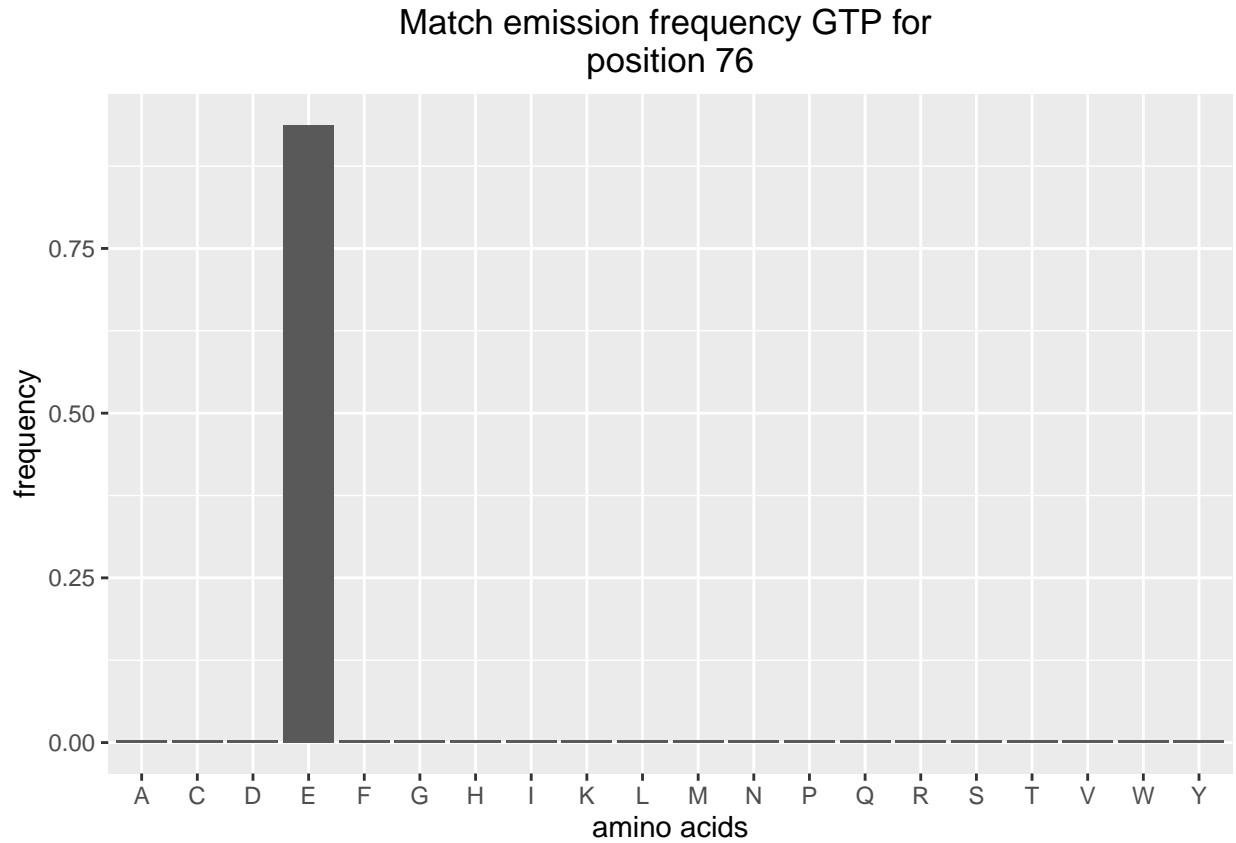
max_emission_ATP <- apply(ATP_HMM$mE, 2, max)
max_pos_emi_ATP <- which.max(max_emission_ATP)

max_insertion_GTP <- apply(GTP_HMM$iE, 2, max)
max_pos_ins_GTP <- which.max(max_insertion_GTP)

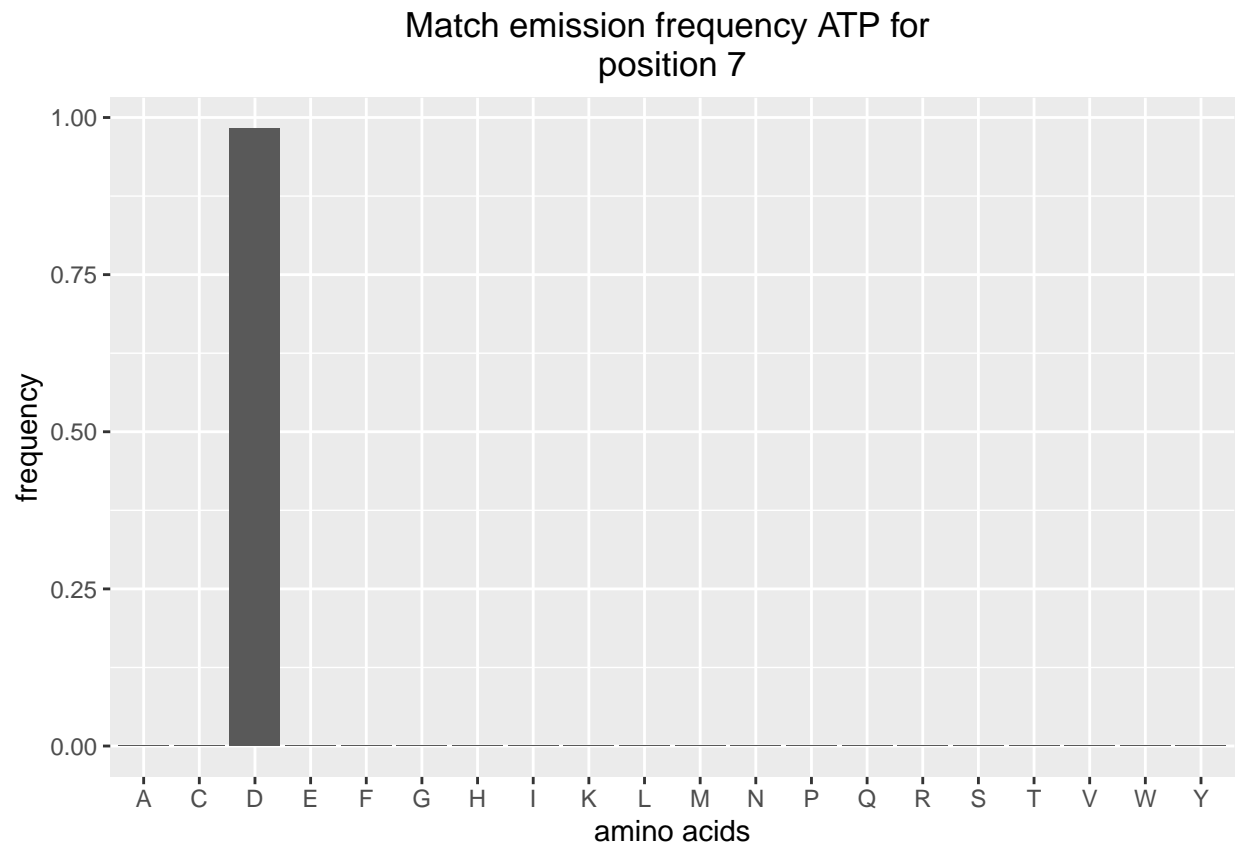
max_insertion_ATP <- apply(ATP_HMM$iE, 2, max)
max_pos_ins_ATP <- which.max(max_insertion_ATP)

df_emi_GTP <- data.frame(freq = GTP_HMM$mE[,max_pos_emi_GTP], rownames = rownames(GTP_HMM$mE))
df_emi_ATP <- data.frame(freq = ATP_HMM$mE[,max_pos_emi_ATP], rownames = rownames(ATP_HMM$mE))
df_ins_GTP <- data.frame(freq = GTP_HMM$iE[,max_pos_ins_GTP], rownames = rownames(GTP_HMM$mE))
df_ins_ATP <- data.frame(freq = ATP_HMM$iE[,max_pos_ins_ATP], rownames = rownames(ATP_HMM$mE))
```

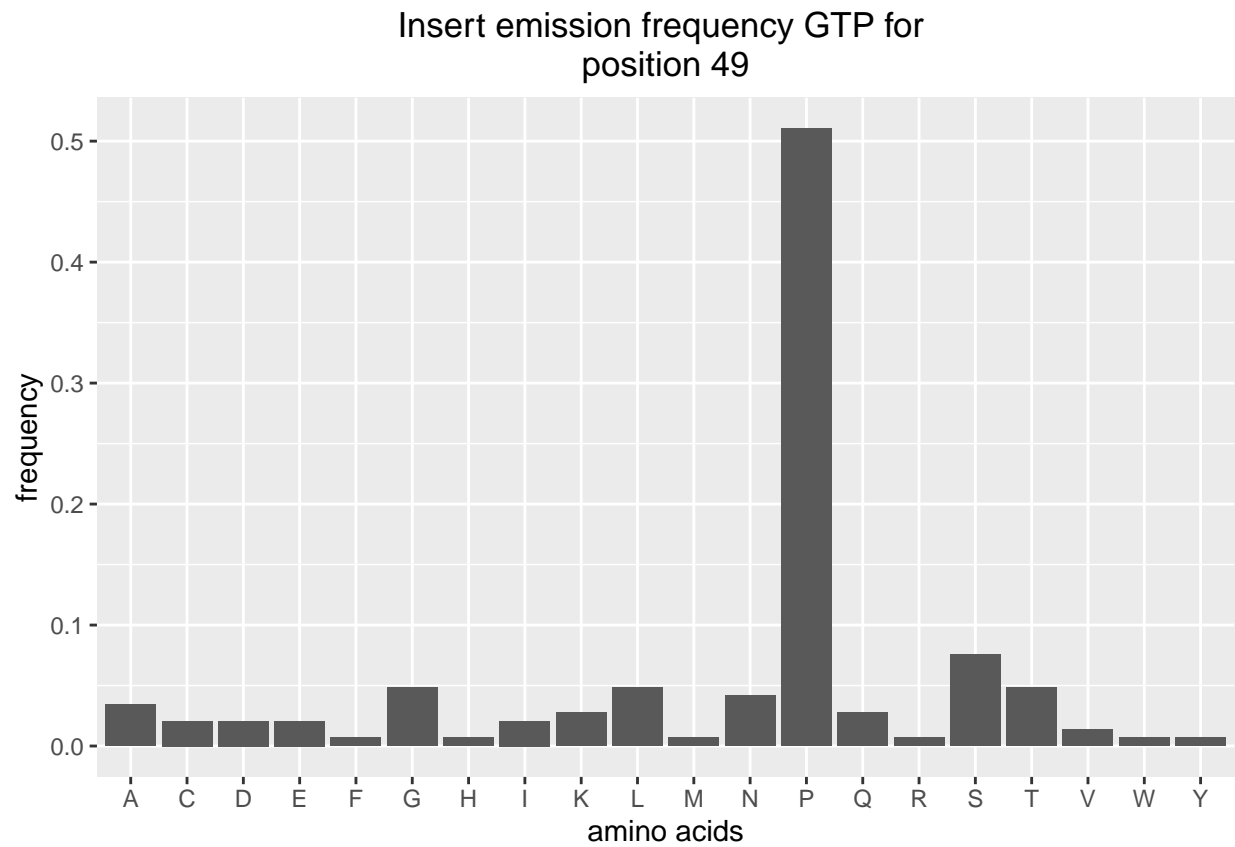
```
ggplot(df_emi_GTP, aes(x = rownames, y = freq))+
  geom_bar(stat = "identity")+
  labs(x="amino acids", y = "frequency", title = paste("Match emission frequency GTP for \nposition",ma
  theme(plot.title = element_text(hjust = 0.5))
```



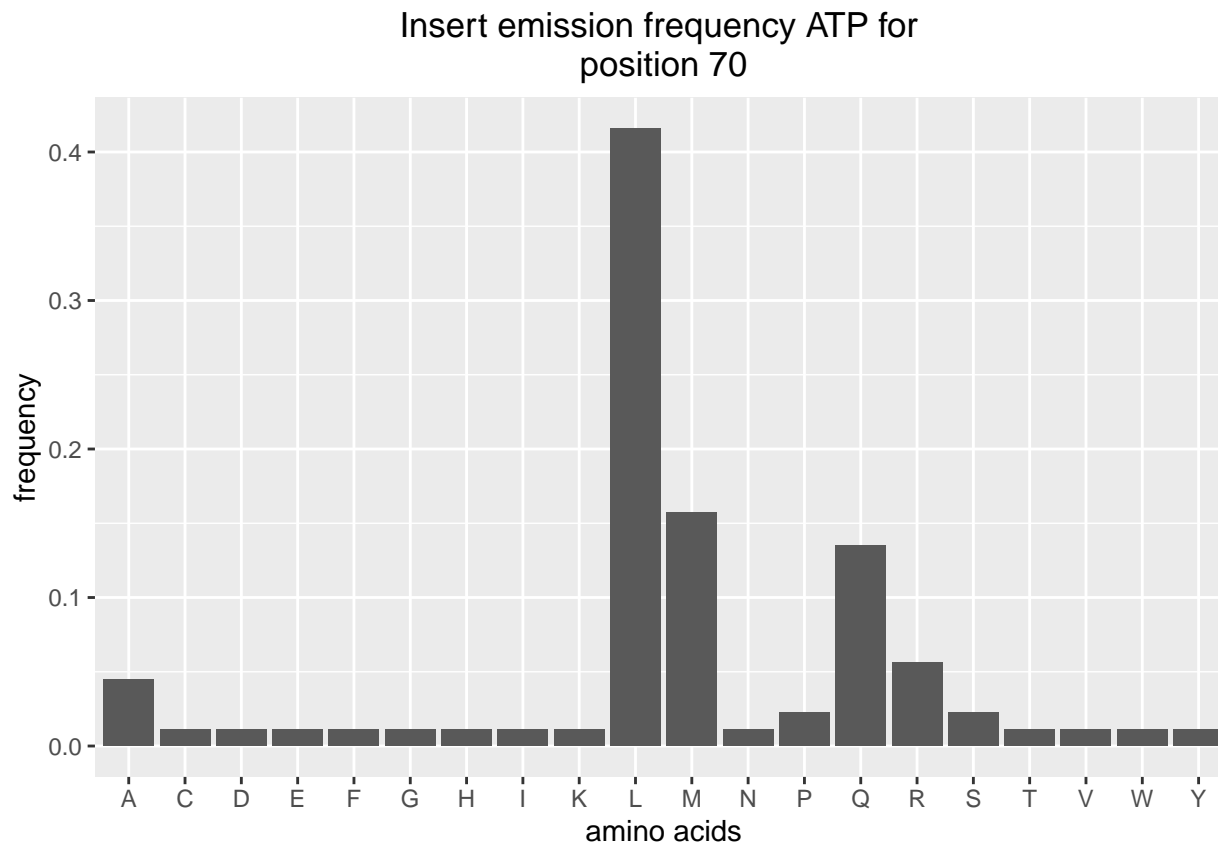
```
ggplot(df_emi_ATP, aes(x = rownames, y = freq))+
  geom_bar(stat = "identity")+
  labs(x="amino acids", y = "frequency", title = paste("Match emission frequency ATP for \nposition",ma
  theme(plot.title = element_text(hjust = 0.5))
```



```
ggplot(df_ins_GTP, aes(x = rownames, y = freq))+
  geom_bar(stat = "identity")+
  labs(x="amino acids", y = "frequency", title = paste("Insert emission frequency GTP for \nposition",m
  theme(plot.title = element_text(hjust = 0.5))
```



```
ggplot(df_ins_ATP, aes(x = rownames, y = freq))+
  geom_bar(stat = "identity")+
  labs(x="amino acids", y = "frequency", title = paste("Insert emission frequency ATP for \nposition",m
  theme(plot.title = element_text(hjust = 0.5))
```



5. The file `Unclassified proteins.txt` contains 31 protein sequences from unknown families. Load the protein sequences into a list using the `parseProteins()` function.

```
unclass_proteins <- parseProteins("C:/Users/marie/Desktop/CBB/SMCB/Project_4_student/data/Unclassified_
```

6. The function `forward()` takes as input a profile HMM `M` and a sequence `x`. It returns the log odds ratio of the probability of observing the sequence `x` given the model `M` versus the probability of observing the sequence `x` given the random model `R`. For each unclassified protein `x(i)` in the list, apply the forward algorithm for both models `M1` and `M2` to obtain the log odds ratio. Plot the values $q(x_i)$ and include this in your report. Which proteins in the list belong to which family? Can you clearly decide for each protein?

```
#apply forward to GTP and ATP
```

```
GTP_forward <- sapply(unclass_proteins, forward, HMM = GTP_HMM)
```

```
ATP_forward <- sapply(unclass_proteins, forward, HMM = ATP_HMM)
```

```
#get the log odds ratio by taking the difference (calculate q values)
```

```
log_odds <- ATP_forward - GTP_forward
```

```
names(log_odds) <- 1:length(unclass_proteins)
```

```
#plot q values
```

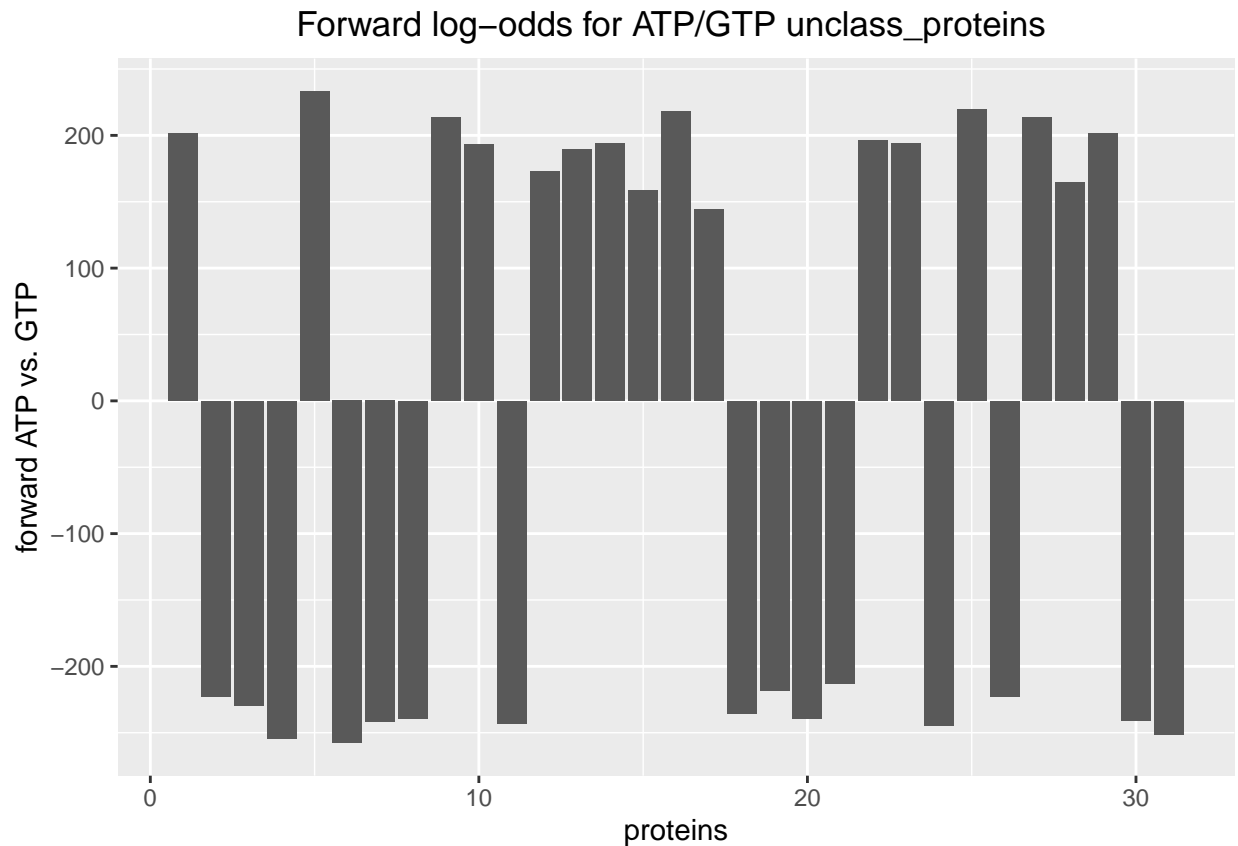
```
df_log_odds <- data.frame(log_odds = log_odds, num_proteins = 1:length(log_odds))
```

```
ggplot(df_log_odds, aes(x = num_proteins, y = log_odds))+
```

```
  geom_bar(stat = "identity")+
```

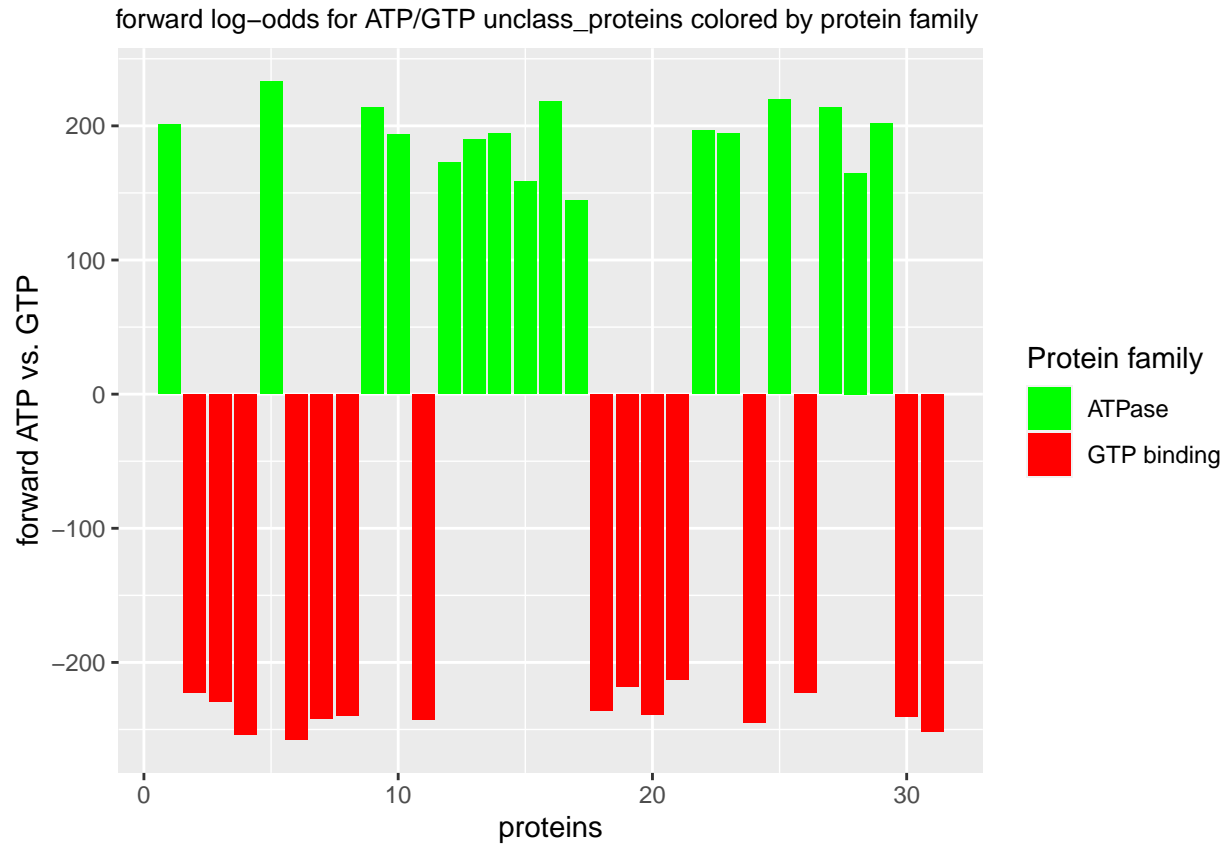


```
labs(x="proteins", y = "forward ATP vs. GTP", title = ("Forward log-odds for ATP/GTP unclass_proteins"),
theme(plot.title = element_text(hjust = 0.5))
```



One can clearly decide for each protein to which family it belongs to based on the log odds ratio ATP/GTP sign. If its positive then the proteins belongs to the ATPase protein family, if its negative to the GTP binding protein family. This can be visualized:

```
df_log_odds$color <- ifelse(df_log_odds$log_odds > 0, "ATPase", "GTP binding")
ggplot(df_log_odds, aes(x = num_proteins, y = log_odds, fill = color)) +
  geom_bar(stat = "identity") +
  labs(x = "proteins", y = "forward ATP vs. GTP",
       title = "forward log-odds for ATP/GTP unclass_proteins colored by protein family ") +
  theme(plot.title = element_text(hjust = 0.5, size = 10)) +
  scale_fill_manual(values = c("ATPase" = "green", "GTP binding" = "red"),
                   name = "Protein family",
                   labels = c("ATPase", "GTP binding"))
```



So proteins number 1,5,9,10,12,13,14,15,16,17,22,23,25,27,28,29 belong to the ATPase family and proteins 2,3,4,6,7,8,11,18,19,20,21,24,26,30,31 to the GTP binding family.