# Somatic hypermutations signatures

*Anna Obraztsova*

*4/20/2017*

```r
library(dplyr)
```

```
##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
##
##     filter, lag

## The following objects are masked from 'package:base':
##
##     intersect, setdiff, setequal, union
```

```r
library(ggplot2)
library(reshape2)
library(stringr)
library(tidyr)
```

```
##
## Attaching package: 'tidyr'

## The following object is masked from 'package:reshape2':
##
##     smiths
```

```r
load('shm.Rda')
shm$region <- factor(shm$region, c("FR1","CDR1","FR2","CDR2","FR3","CDR3"))
shm$clonotypes <- 1
shm$name <- paste(shm$proj, shm$sample, shm$cells, sep='_')

contig.melted <- shm %>% dplyr::select(proj, sample, cells, name, contignt) %>%
  unique() %>%
  mutate(nt = str_split(contignt, '')) %>%
  unnest(nt = nt)

contig.total <- contig.melted %>% dplyr::group_by(proj, sample, cells, name) %>%
  dplyr::summarise(total.nt = n())

contig.freq <- contig.melted %>% dplyr::group_by(proj, sample, cells, nt, name) %>%
  dplyr::summarise(count = n()) %>%
  merge(contig.total) %>%
  mutate(nt.freq = count/total.nt)

total.mutations <- shm %>% dplyr::group_by(proj, sample, cells, name) %>%
  dplyr::summarise(total.mutations = n())

f <- shm %>% dplyr::group_by(proj, sample, cells, name, from, to) %>%
  dplyr::summarise(total.clonotypes = sum(clonotypes)) %>%
  merge(dplyr::select(contig.freq, from=nt, proj, sample, cells, nt.freq, name)) %>%
  merge(total.mutations) %>%
```

```r
    mutate(weight = total.clonotypes/(total.mutations*nt.freq))

f1 <- f %>% dplyr::group_by(proj, sample, cells, name) %>%
  dplyr::summarise(weight.total = sum(weight))

f <- f %>% merge(f1) %>%
  mutate(freq = weight / weight.total) %>%
  dplyr::select(proj, sample, cells, from, to, freq, name)

st <- read.table('freq_steele.txt', header=T) %>%
  mutate(proj = 'steele', sample = 'steele', cells='N', name = 'steele',
         freq = freq/100)

f <- rbind(f, st)

ggplot(f, aes(x = from, y = to)) +
  geom_tile(aes(fill = freq)) +
  geom_text(aes(label = round(freq, 3)), cex=3) +
  facet_wrap(~name, ncol=3) +
  scale_fill_gradient2(low = "#2F6B89", mid = "#67CB87", high = "#F4F27B", midpoint = 0.1, na.value = "
```
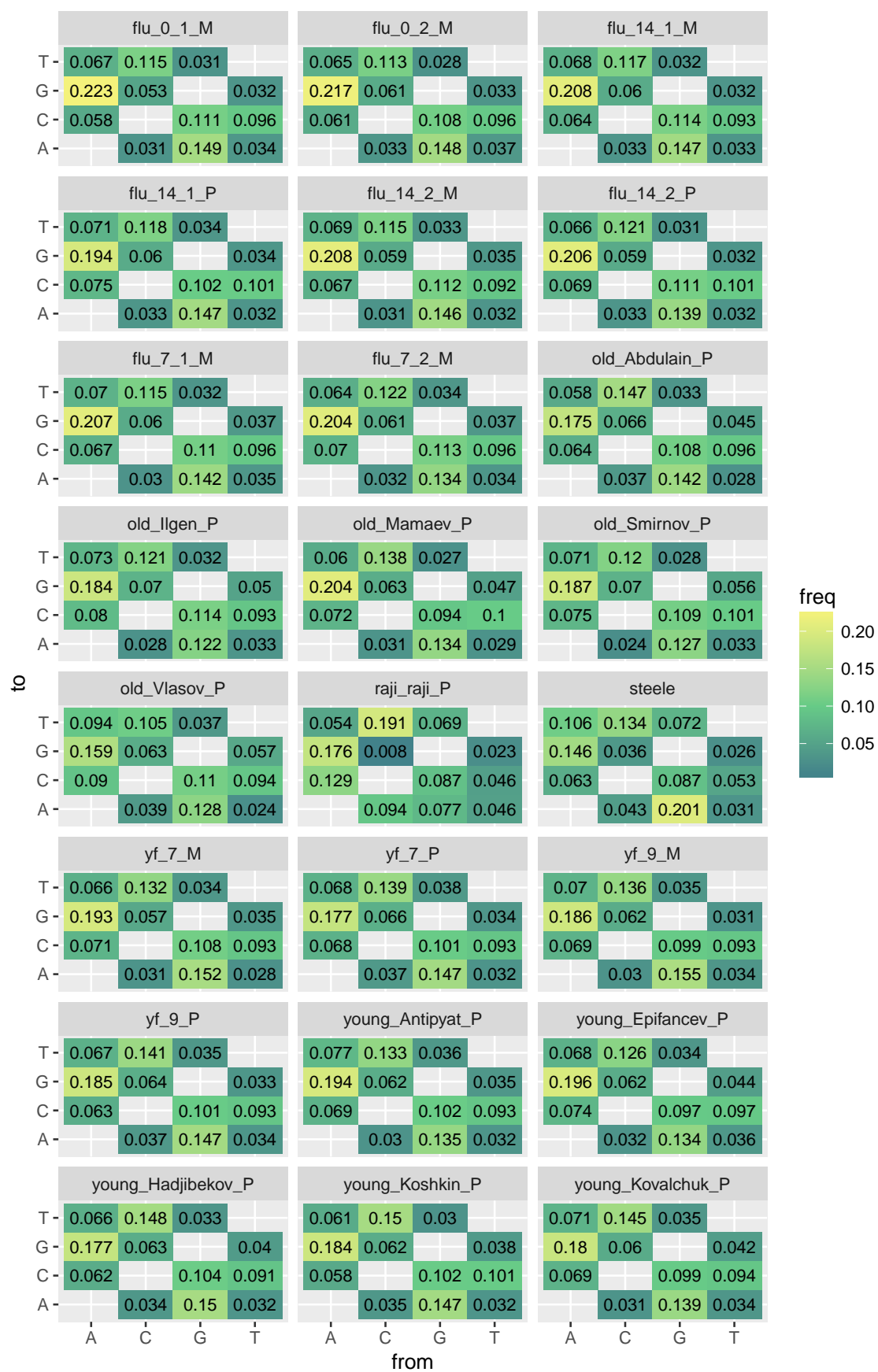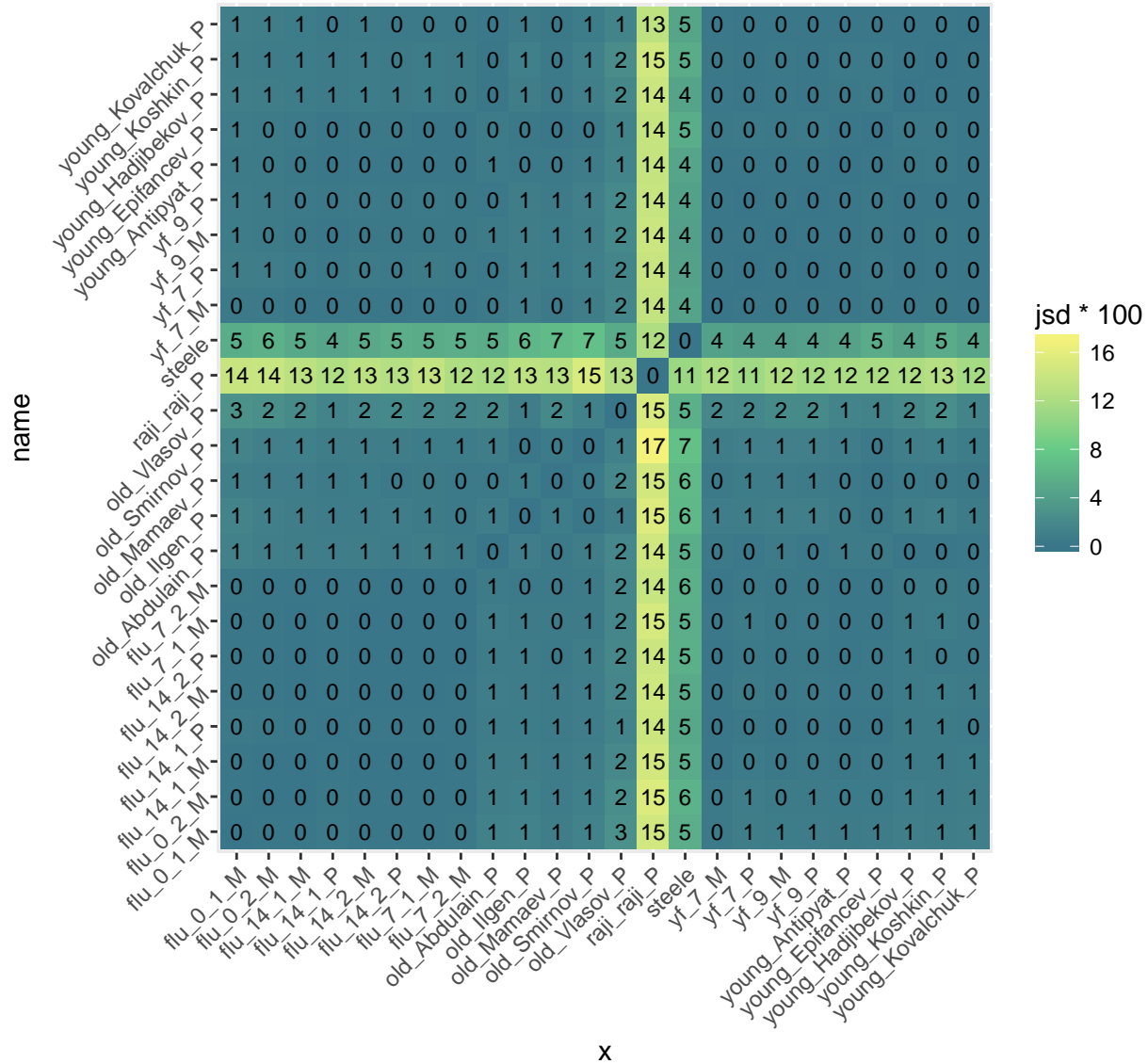
Calculate Jensen-Shannon divergence

```r
jsd <- function(v1, v2){
  m <- 0.5 * (v1 + v1)
  0.5 * (sum(v1 * log2(v1 / m)) + sum(v2 * log2(v2 / m)))
}


f <- f[ order(f$from, f$to), ]


jsd.df <- data.frame()


for (i in unique(f$name)){
  x <- filter(f, name == i)$freq

  .jsd.df <- f %>% dplyr::group_by(name) %>%
  dplyr::summarise(jsd = jsd(x, freq))

  .jsd.df$x <- i

  jsd.df <- rbind(jsd.df, .jsd.df)
}

ggplot(jsd.df, aes(x, name)) +
  geom_tile(aes(fill = jsd*100)) +
  geom_text(aes(label = round(jsd*100)), cex=3) +
  scale_fill_gradient2(low = "#2F6B89", mid = "#67CB87", high = "#F4F27B", midpoint = 8,
                       na.value = "white") +
    theme(axis.text.x = element_text(angle = 45, hjust = 1),
          axis.text.y = element_text(angle = 45, hjust = 1))
```

Heatmap of pairwise jsd * 100 values (legend: 0, 4, 8, 12, 16).

| name \ x | flu_0_1_M | flu_0_2_M | flu_14_1_M | flu_14_1_P | flu_14_2_M | flu_14_2_P | flu_7_1_M | flu_7_2_P | old_Abdulain_P | old_Ilgen_P | old_Mamaev_P | old_Smirnov_P | old_Vlasov_P | raji_raji_P | steele | yf_7_M | yf_7_P | yf_9_M | yf_9_P | young_Antipyat_P | young_Epifancev_P | young_Hadjibekov_P | young_Koshkin_P | young_Kovalchuk_P |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| young_Kovalchuk_P | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 13 | 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| young_Koshkin_P | 1 | 1 | 1 | 1 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 2 | 15 | 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| young_Hadjibekov_P | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 1 | 0 | 1 | 2 | 14 | 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| young_Epifancev_P | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 14 | 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| young_Antipyat_P | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 1 | 14 | 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| yf_9_P | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 2 | 14 | 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| yf_9_M | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 2 | 14 | 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| yf_7_P | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 1 | 2 | 14 | 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| yf_7_M | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 2 | 14 | 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| steele | 5 | 6 | 5 | 4 | 5 | 5 | 5 | 5 | 5 | 6 | 7 | 7 | 5 | 12 | 0 | 4 | 4 | 4 | 4 | 4 | 5 | 4 | 5 | 4 |
| raji_raji_P | 14 | 14 | 13 | 12 | 13 | 13 | 13 | 12 | 12 | 13 | 13 | 15 | 13 | 0 | 11 | 12 | 11 | 12 | 12 | 12 | 12 | 12 | 13 | 12 |
| old_Vlasov_P | 3 | 2 | 2 | 1 | 2 | 2 | 2 | 2 | 1 | 2 | 1 | 0 | 15 | 5 | 2 | 2 | 2 | 1 | 1 | 2 | 2 | 1 |  |  |
| old_Smirnov_P | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 17 | 7 | 1 | 1 | 1 | 1 | 1 | 0 | 1 | 1 | 1 |
| old_Mamaev_P | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 2 | 15 | 6 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| old_Ilgen_P | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 15 | 6 | 1 | 1 | 1 | 1 | 0 | 1 | 1 | 1 |
| old_Abdulain_M | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 2 | 14 | 5 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 |
| flu_7_2_M | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 2 | 14 | 6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| flu_7_1_P | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 2 | 15 | 5 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 |
| flu_14_2_M | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 2 | 14 | 5 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| flu_14_2_P | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 2 | 14 | 5 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 |  |
| flu_14_1_P | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 14 | 5 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 |
| flu_14_1_M | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 2 | 15 | 5 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 |
| flu_0_2_M | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 2 | 15 | 6 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 1 | 1 |
| flu_0_1_M | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 3 | 15 | 5 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |

```r
f2 <- f %>% dplyr::group_by(proj, sample, cells, name, from) %>%
  dplyr::summarise(freq = sum(freq)) %>%
  dcast(proj + sample + cells +name ~ from, value.var='freq')

f2$A.AT.ratio = f2$A/(f2$A+f2[['T']])
f2$G.GC.ratio = f2$G/(f2$G+f2$C)
f2 <- f2 %>%
  dplyr::select(name, A.AT.ratio, G.GC.ratio) %>%
  melt(id.vars=c('name'))

ggplot(f2, aes(variable, name)) +
  geom_tile(aes(fill = value)) +
  geom_text(aes(label = round(value, 3)), cex=3) +
  scale_fill_gradient2(low = "#2F6B89", mid = "#67CB87", high = "#F4F27B", midpoint = 0.6,
                       na.value = "white")
```

| name | A.AT.ratio | G.GC.ratio |
|---|---|---|
| young_Kovalchuk_P | 0.653 | 0.535 |
| young_Koshkin_P | 0.639 | 0.53 |
| young_Hadjibekov_P | 0.651 | 0.539 |
| young_Epifancev_P | 0.657 | 0.547 |
| young_Antipyat_P | 0.681 | 0.548 |
| yf_9_P | 0.663 | 0.54 |
| yf_9_M | 0.673 | 0.559 |
| yf_7_P | 0.663 | 0.541 |
| yf_7_M | 0.68 | 0.572 |
| steele | 0.741 | 0.628 |
| raji_raji_P | 0.757 | 0.443 |
| old_Vlasov_P | 0.663 | 0.571 |
| old_Smirnov_P | 0.638 | 0.553 |
| old_Mamaev_P | 0.656 | 0.522 |
| old_Ilgen_P | 0.657 | 0.55 |
| old_Abdulain_P | 0.637 | 0.53 |
| flu_7_2_M | 0.671 | 0.566 |
| flu_7_1_M | 0.671 | 0.58 |
| flu_14_2_P | 0.674 | 0.569 |
| flu_14_2_M | 0.683 | 0.586 |
| flu_14_1_P | 0.671 | 0.573 |
| flu_14_1_M | 0.683 | 0.583 |
| flu_0_2_M | 0.674 | 0.578 |
| flu_0_1_M | 0.681 | 0.594 |

value
0.75
0.70
0.65
0.60
0.55
0.50
0.45

variable

```r
library(seqLogo)
```

```
## Loading required package: grid
```

```r
shm$type <- paste(shm$from, shm$to, sep='>')

for (i in unique(shm$type)){
  print(i)
  set <- filter(shm, type == i, nchar(context) == 7, str_detect(context, 'N') == F) %>%
    separate(context, c('n.1', 'n.2', 'n.3', 'n.4', 'n.5', 'n.6', 'n.7'), c(1,2,3,4,5,6)) %>%
    dplyr::select(starts_with('n.'))

  prob <- data.frame()

  for (j in c('A','C','G','T')){
    .prob <- apply(set, 2, function(x) length(x[x==j])/length(x))
    prob <- rbind(prob, .prob)
  }
  rownames(prob) <- c('A','C','G','T')
  colnames(prob) <- paste0('p', 1:7)
  seqLogo(makePWM(prob), ic.scale=F,  xfontsize=10, yfontsize=10)
}
```
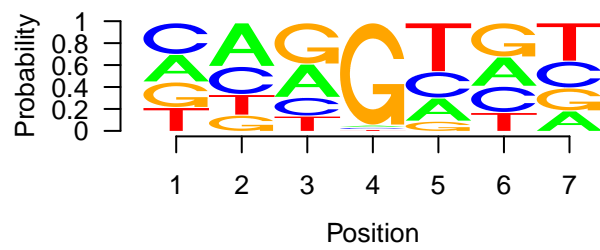
```
## [1] "T>C"
```
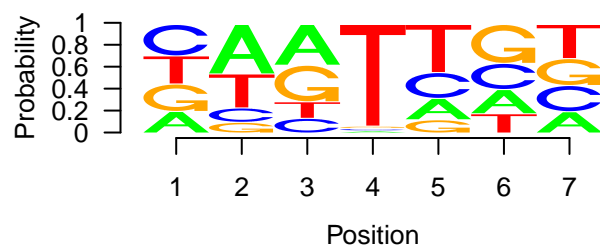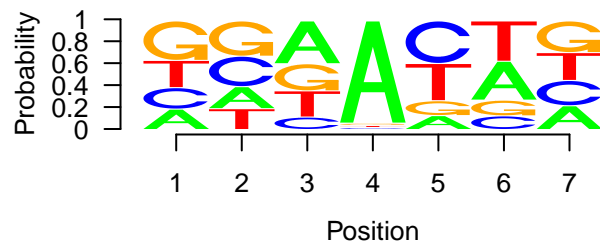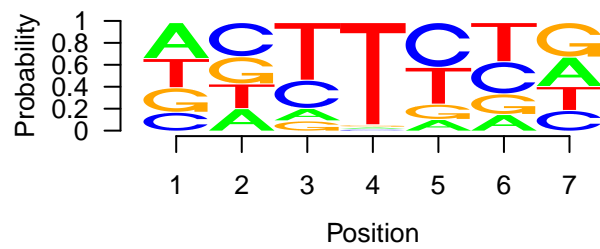
## [1] "A>G"



## [1] "G>C"



## [1] "C>G"



## [1] "C>T"



## [1] "G>A"

7

## [1] "A>T"



## [1] "C>A"



## [1] "A>C"



## [1] "T>A"



## [1] "G>T"

8

## [1] "T>G"