# istoype_ext_analysis.Rmd

```r
library(data.table)
library(dplyr)
```

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

## The following objects are masked from 'package:data.table':
##
##     between, first, last

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

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

```r
library(stringr)
library(ggplot2)
library(NMF)
```

```
## Loading required package: pkgmaker

## Loading required package: registry

##
## Attaching package: 'pkgmaker'

## The following object is masked from 'package:base':
##
##     isNamespaceLoaded

## Loading required package: rngtools

## Loading required package: cluster

## NMF - BioConductor layer [NO: missing Biobase] | Shared memory capabilities [NO: bigmemory] | Cores 5

##   To enable the Bioconductor layer, try: install.extras('
## NMF
## ') [with Bioconductor repository enabled]
##   To enable shared memory capabilities, try: install.extras('
## NMF
## ')
```

```r
library(reshape2)
```

```
##
## Attaching package: 'reshape2'

## The following objects are masked from 'package:data.table':
##
##     dcast, melt
```

```r
library(parallel)
library(RColorBrewer)
library(scales)
```

```
summarise = dplyr::summarise

load("shm_rep12_downsampled.rda")

df = shm %>%
  mutate(replacement = ifelse(as.character(from.aa) != as.character(to.aa), "replacement", "silent"),
         i = isotype) %>%
  group_by(clone, sample, proj, i, replacement) %>%
  summarise(count = n())

df.s = df %>%
  group_by(proj, i) %>%
  summarise(count = n())

print(df.s)
```

```
## # A tibble: 17 x 3
## # Groups:   proj [?]
##     proj  i      count
##     <chr> <chr>  <int>
##  1 old    ""         1
##  2 old    IGHA1   1430
##  3 old    IGHD       2
##  4 old    IGHE      60
##  5 old    IGHG1   1434
##  6 old    IGHG2      2
##  7 old    IGHG3     73
##  8 old    IGHGP      6
##  9 old    IGHM    2158
## 10 young  IGHA1   2490
## 11 young  IGHD      39
## 12 young  IGHE      73
## 13 young  IGHG1   1823
## 14 young  IGHG2      2
## 15 young  IGHG3    118
## 16 young  IGHGP     13
## 17 young  IGHM    2039
```

```
df = df %>% filter(i != "IGHD", !is.na(i), i != "") %>%
  mutate(isotype.full = i, isotype = str_sub(i, 4, 4))

df$proj = factor(df$proj, levels = c('young', 'old'))
```

```
dt.p = data.table()

for (iso in unique(df$isotype)) {
  tmp = df %>% filter(isotype == iso)
  x = (tmp %>% filter(proj == "old"))$count
  y = (tmp %>% filter(proj != "old"))$count
  kk = ks.test(x, y)
  p = kk$p.value
  dt.p = rbind(dt.p,
               data.table(isotype = iso, p=p))
}
```

```
## Warning in ks.test(x, y): p-value will be approximate in the presence of
## ties

## Warning in ks.test(x, y): cannot compute exact p-value with ties

## Warning in ks.test(x, y): p-value will be approximate in the presence of
## ties

## Warning in ks.test(x, y): p-value will be approximate in the presence of
## ties
```

```r
dt.p$p.adj = p.adjust(dt.p$p, method = "BH")
print(dt.p %>% arrange(p.adj))
```

```
##   isotype         p        p.adj
## 1       M 5.998098e-08 2.399239e-07
## 2       A 3.384933e-01 4.513244e-01
## 3       G 2.639686e-01 4.513244e-01
## 4       E 9.636272e-01 9.636272e-01
```
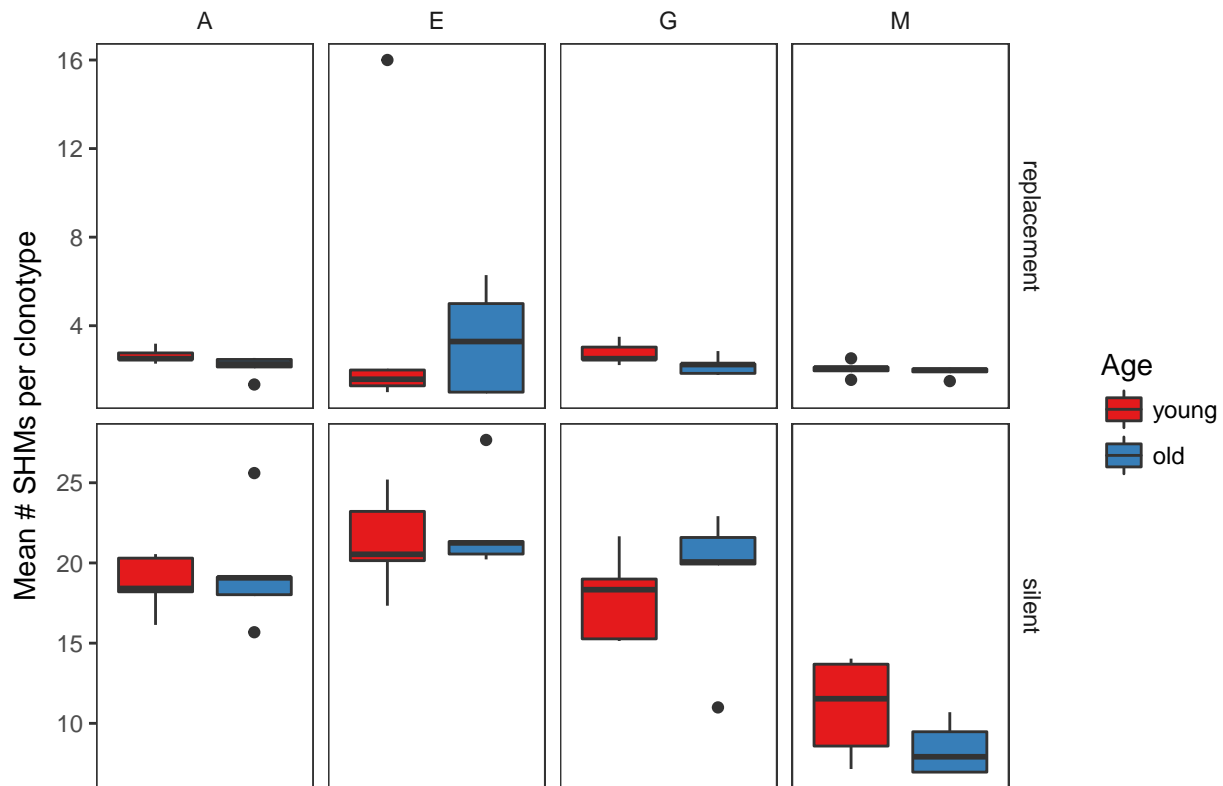
```r
p13=ggplot(df, aes(x = count, fill = proj)) +
  geom_density(alpha = 0.9, color = NA) +
  facet_wrap(~isotype) +
  scale_fill_brewer("Age", palette = "Set1") +
  xlab("SHMs per clonotype") + ylab("") +
  theme_bw() +
  theme(panel.grid.major = element_blank(), panel.grid.minor = element_blank(),
        strip.background = element_blank())
ggsave("figures/p13.pdf", p13)
```

```
## Saving 6.5 x 4.5 in image
```

```r
df.1 = df %>%
  group_by(sample, proj, replacement, isotype) %>%
  summarise(shms = mean(count))

p10=ggplot(df.1, aes(x=proj, fill = proj, y = shms)) +
  geom_boxplot() +
  facet_grid(replacement~isotype, scales = "free") +
  scale_fill_brewer("Age", palette = "Set1") +
  xlab("") + ylab("Mean # SHMs per clonotype") +
  theme_bw() +
  theme(panel.grid.major = element_blank(), panel.grid.minor = element_blank(),
        axis.text.x = element_blank(), axis.ticks.x = element_blank(),
        strip.background = element_blank())
p10
```

```
ggsave("figures/p10.pdf", p10)
```

## Saving 6.5 x 4.5 in image

```
a = aov(shms ~ replacement + isotype + proj, df.1)
summary(a)
```

```
##              Df Sum Sq Mean Sq F value   Pr(>F)
## replacement   1   4254    4254  368.25  < 2e-16 ***
## isotype       3    518     173   14.95 1.06e-07 ***
## proj          1      1       1    0.10    0.753
## Residuals    74    855      12
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
TukeyHSD(a, "proj")
```

```
##   Tukey multiple comparisons of means
##     95% family-wise confidence level
##
## Fit: aov(formula = shms ~ replacement + isotype + proj, data = df.1)
##
## $proj
##                 diff       lwr      upr     p adj
## old-young -0.2399598 -1.754361 1.274442 0.7531021
```
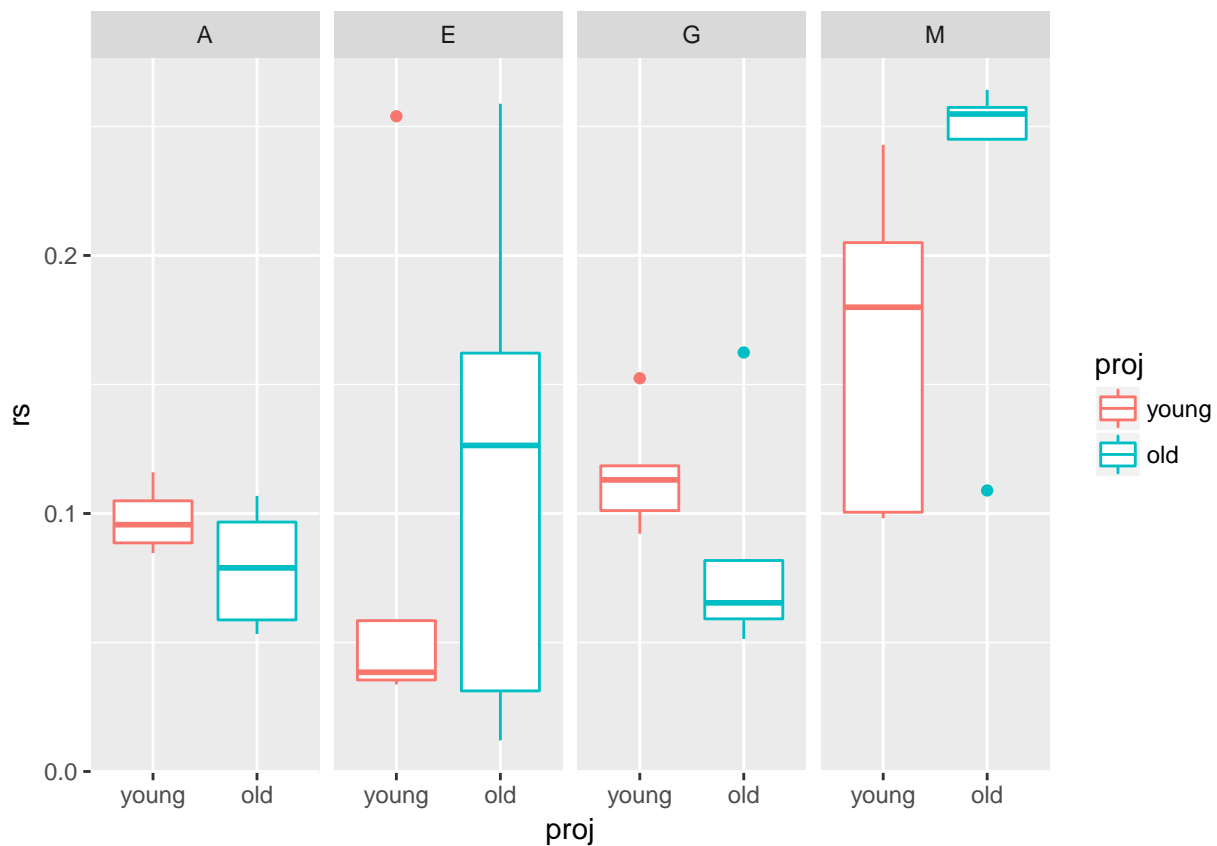
```
TukeyHSD(a, "isotype")
```

```
##   Tukey multiple comparisons of means
##     95% family-wise confidence level
##
```

```
## Fit: aov(formula = shms ~ replacement + isotype + proj, data = df.1)
##
## $isotype
##           diff        lwr        upr      p adj
## E-A  2.0394954 -0.785626  4.8646168 0.2380135
## G-A -0.2741454 -3.099267  2.5509760 0.9941308
## M-A -4.9141574 -7.739279 -2.0890360 0.0001101
## G-E -2.3136408 -5.138762  0.5114806 0.1463769
## M-E -6.9536528 -9.778774 -4.1285314 0.0000001
## M-G -4.6400120 -7.465133 -1.8148907 0.0002779
```

```r
df.2 = df %>%
  group_by(sample, proj, isotype) %>%
  summarise(rs = sum(count[which(replacement == "replacement")]) / sum(count[which(replacement != "repl
```

```r
ggplot(df.2, aes(x=proj, color = proj, y = rs)) +
  geom_boxplot() +
  facet_grid(.~isotype, scales = "free")
```
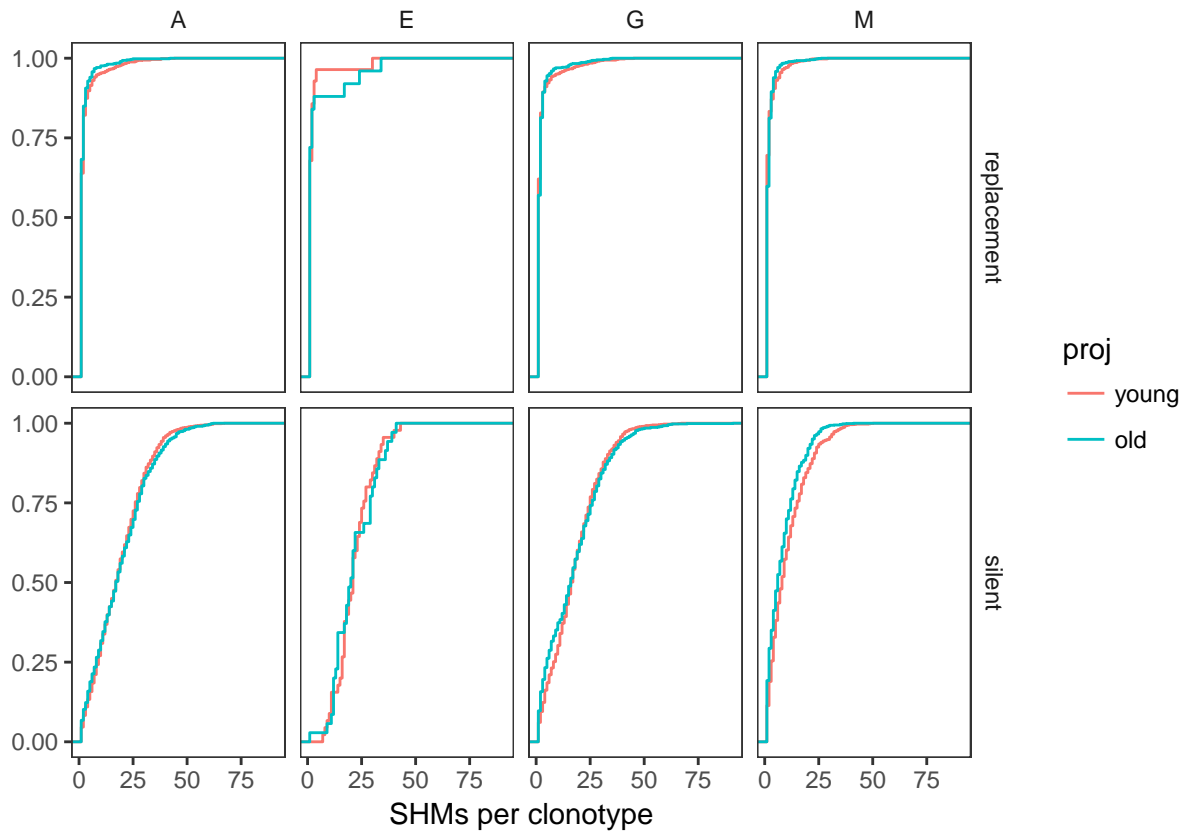


```r
a = aov(rs ~ isotype + proj, df.2)
summary(a)
```

```
##             Df  Sum Sq  Mean Sq F value  Pr(>F)
## isotype      3 0.07488 0.024959   6.313 0.00155 **
## proj         1 0.00123 0.001234   0.312 0.58001
## Residuals   35 0.13837 0.003954
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
TukeyHSD(a, "proj")
```

```
##   Tukey multiple comparisons of means
##     95% family-wise confidence level
##
## Fit: aov(formula = rs ~ isotype + proj, data = df.2)
##
## $proj
##                 diff        lwr        upr      p adj
## old-young 0.01110631 -0.0292593 0.05147192 0.5800118
```

```r
ggplot(df, aes(x=count, color = proj)) +
  stat_ecdf() +
  facet_grid(replacement~isotype) + #, scales = "free") +
  #scale_fill_brewer("Age", palette = "Set1") +
  xlab("SHMs per clonotype") + ylab("") +
  theme_bw() +
  theme(panel.grid.major = element_blank(), panel.grid.minor = element_blank(),
        strip.background = element_blank())
```



```r
dt.p = data.table()

for (iso in unique(df$isotype)) {
for (rr in unique(df$replacement)) {
  tmp = df %>% filter(isotype == iso, replacement == rr)
  x = (tmp %>% filter(proj == "old"))$count
  y = (tmp %>% filter(proj != "old"))$count
  kk = ks.test(x, y)
```

```
  p = kk$p.value
  dt.p = rbind(dt.p,
               data.table(isotype = iso, replacement = rr, p=p))
}
}
```

```
## Warning in ks.test(x, y): p-value will be approximate in the presence of
## ties

## Warning in ks.test(x, y): p-value will be approximate in the presence of
## ties

## Warning in ks.test(x, y): cannot compute exact p-value with ties

## Warning in ks.test(x, y): cannot compute exact p-value with ties

## Warning in ks.test(x, y): p-value will be approximate in the presence of
## ties

## Warning in ks.test(x, y): p-value will be approximate in the presence of
## ties

## Warning in ks.test(x, y): p-value will be approximate in the presence of
## ties

## Warning in ks.test(x, y): p-value will be approximate in the presence of
## ties
```

```
dt.p$p.adj = p.adjust(dt.p$p, method = "BH")
print(dt.p %>% arrange(p.adj))
```

```
##   isotype replacement            p         p.adj
## 1       M      silent 1.156105e-05 9.248838e-05
## 2       M replacement 2.043525e-04 8.174100e-04
## 3       G      silent 1.660173e-03 4.427127e-03
## 4       G replacement 3.232664e-01 6.465328e-01
## 5       A replacement 4.837888e-01 7.740621e-01
## 6       E      silent 6.567257e-01 7.832819e-01
## 7       A      silent 6.853717e-01 7.832819e-01
## 8       E replacement 9.999841e-01 9.999841e-01
```