

Covid-19 SNP analysis

Code ▼

Xiaomi Liu

Hide

```
library(dplyr)
```

Warning: package 'dplyr' was built under R version 4.3.2

Hide

```
library(ggplot2)  
library(adeigenet)
```

Warning: package 'adeigenet' was built under R version 4.3.3Warning: package 'ade4' was built under R version 4.3.3

Hide

```
library(LEA)  
library(dartR)
```

Warning: package 'dartR' was built under R version 4.3.3Warning: package 'dartR.data' was built under R version 4.3.3

Hide

```
library(gtools)
```

Warning: package 'gtools' was built under R version 4.3.3

Hide

```
library(tidyverse)
```

Warning: package 'tidyverse' was built under R version 4.3.3Warning: package 'tibble' was built under R version 4.3.2Warning: package 'tidyr' was built under R version 4.3.3Warning: package 'readr' was built under R version 4.3.3Warning: package 'purrr' was built under R version 4.3.3Warning: package 'lubridate' was built under R version 4.3.3

Hide

```
library(RColorBrewer)
```

Read in data

Hide

```
map <- read.table("../cov19_snp/cov19.map", header = F)
ped <- read.table("../cov19_snp/cov19.ped", header = F)
info <- read.table("../../sample_info.txt", header = T)
```

Modify data

Hide

```
#rename chromosome
map$V1 <- "genome"
#write to output
write.table(map,file = "../cov19_snp/covid19edit.map", quote = FALSE, sep = "\t", row.names = FALSE, col.names = FALSE)
#assign group
ped$V1[1] <- "reference"
ped$V1[2:1931] <- info$state
#set V1 as factor
ped$V1 <- as.factor(ped$V1)
levels(ped$V1)
```

```
[1] "AK"      "AL"      "AR"      "AZ"      "CA"      "CO"      "CT"
[8] "DE"      "FL"      "GA"      "HI"      "IA"      "ID"      "IL"
[15] "IN"      "KS"      "KY"      "LA"      "MA"      "MD"      "ME"
[22] "MI"      "MN"      "MO"      "MS"      "MT"      "NC"      "NE"
[29] "NH"      "NJ"      "NM"      "NV"      "NY"      "OH"      "OK"
[36] "OR"      "PA"      "PR"      "reference" "RI"      "SC"      "SD"
[43] "TN"      "TX"      "UT"      "VA"      "VT"      "WA"      "WI"
[50] "WV"      "WY"
```

Hide

```
#reorder levels in Column 1
ped$V1 <- factor(ped$V1, levels = c("reference", unique(info$state) ))
ped <- ped[order(ped$V1),]
#write to output
write.table(ped,file = "../cov19_snp/covid19edit.ped", quote = FALSE, sep = "\t", row.names = FALSE, col.names = FALSE)
#generate geno file
ped2geno("../cov19_snp/covid19edit.ped", output.file = "../cov19_snp/covid19edit.geno")
```

- number of detected individuals: 1931
- number of detected loci: 2599

```
[1] "../cov19_snp/covid19edit.geno"
```

Hide

```
#use PLINK to reformat our data into raw output  
system("bash -c '../plink/plink.exe --version'")
```

```
PLINK v1.90b6.20 64-bit (21 Sep 2020)  
[1] 0
```

[Hide](#)

```
system("bash -c '../plink/plink.exe --file ../cov19_snp/covid19edit --chr-set -1 --allow-extra-c  
hr --recodeA --out ../cov19_snp/covid19Raw'")
```

0%
1%
2%
3%
4%
5%
6%
7%
8%
9%
10%
11%
12%
13%
14%
15%
16%
17%
18%
19%
20%
21%
22%
23%
24%
25%
26%
27%
28%
29%
30%
31%
32%
33%
34%
35%
36%

5/16

```
.ped scan complete (for binary autoconversion).
Performing single-pass .bed write (2599 variants, 1931 samples).
0%001%002%003%004%005%006%007%008%009%010%011%012%013%014%015%016%017%018%019%
09%020%021%022%023%024%025%026%027%028%029%030%031%032%033%034%035%036%037%038%039%040%041%042%043%044%045%046%047%048%049%050%051%052%053%054%055%056%057%058%059%060%061%062%063%064%065%066%067%068%069%070%071%072%073%074%075%076%077%078%079%080%081%082%083%084%085%086%087%088%089%090%091%092%093%094%
--file: ../cov19_snp/covid19Raw-temporary.bed +
../cov19_snp/covid19Raw-temporary.bim + ../cov19_snp/covid19Raw-temporary.fam
written.
2599 variants loaded from .bim file.
1931 samples (0 males, 0 females, 1931 ambiguous) loaded from .fam.
Ambiguous sex IDs written to ../cov19_snp/covid19Raw.nosex .
Using 1 thread (no multithreaded calculations invoked).
Before main variant filters, 1931 founders and 0 nonfounders present.
Calculating allele frequencies... 0%001%002%003%004%005%006%007%008%009%010%011%012%013%014%015%016%017%018%019%020%021%022%023%024%025%026%027%028%029%030%031%032%033%034%035%036%037%038%039%040%041%042%043%044%045%046%047%048%049%050%051%052%053%054%055%056%057%058%059%060%061%062%063%064%065%066%067%068%069%070%071%072%073%074%075%076%077%078%079%080%081%082%083%084%085%086%087%088%089%090%091%092%093%094%095%096%097%098%099%100% done.
Total genotyping rate is 0.991728.
2599 variants and 1931 samples pass filters and QC.
Note: No phenotypes present.
--recode A to ../cov19_snp/covid19Raw.raw ... 0%001%002%003%004%005%006%007%008%009%010%011%
```

```
1%22212%22213%22214%22215%22216%22217%22218%22219%22220%22221%22222%22223%22224%22225%22226%22227%22228%22229%22230%22231%22232%22233%22234%22235%22236%22237%22238%22239%22240%22241%22242%22243%22244%22245%22246%22247%22248%22249%22250%22251%22252%22253%22254%22255%22256%22257%22258%22259%22260%22261%22262%22263%22264%22265%22266%22267%22268%22269%22270%22271%22272%22273%22274%22275%22276%22277%22278%22279%22280%22281%22282%22283%22284%22285%22286%22287%22288%22289%22290%22291%22292%22293%22294%22295%22296%22297%22298%22299%222done.
[1] 0
```

Read in edited SNP data

Hide

```
# now import data into adegenet using the function read.PLINK()
# Windows users will want to add the argument parallel = FALSE
covidSNPs <- read.PLINK("../cov19_snp/covid19Raw.raw", parallel = F)
```

Reading PLINK raw format into a genlight object...

Reading loci information...

Reading and converting genotypes...

..

Building final object...

...done.

Hide

```
# check to see if it worked by just executing the object name
names(covidSNPs)
```

```
[1] "gen"      "n.loc"    "ind.names" "loc.names" "loc.all"   "chromosome" "position"
[8] "ploidy"   "pop"      "strata"    "hierarchy" "other"
```

Hide

```
covidSNPs$ploidy <- as.integer(rep(1, 1931))
# save
#save.image(file = "../SNP_analysis.rdata")
```

Sanity check

Hide

```
#use nInd() to check the number of individuals
nInd(covidSNPs)
```

```
[1] 1931
```

Hide

```
# use nLoc() to check the number of SNPs (loci)
nLoc(covidSNPs)
```

```
[1] 2599
```

Hide

```
#pop() returns a factor indicating the population of each individual
#use table() to produce a table of the output of pop()
table(pop(covidSNPs))
```

AK	AL	AR	AZ	CA	CO	CT	DE	FL
50	16	50	50	50	50	50	50	50
GA	HI	IA	ID	IL	IN	KS	KY	LA
50	10	50	50	50	34	3	5	50
MA	MD	ME	MI	MN	MO	MS	MT	NC
50	50	50	50	50	50	44	4	50
NE	NH	NJ	NM	NV	NY	OH	OK	OR
50	5	50	50	50	50	21	13	50
PA	PR	reference	RI	SC	SD	TN	TX	UT
50	35	1	7	50	16	11	50	50
VA	VT	WA	WI	WV	WY			
50	4	50	50	2	50			

Hide

```
#remove groups with only 1 individual
covidSNPs <- subset(covidSNPs, covidSNPs$pop != "reference" & covidSNPs$pop != "FL")
table(pop(covidSNPs))
```

```
AK AL AR AZ CA CO CT DE GA HI IA ID IL IN KS KY LA MA MD ME MI MN MO MS MT NC NE NH NJ NM NV NY
50 16 50 50 50 50 50 50 50 10 50 50 50 34 3 5 50 50 50 50 50 50 44 4 50 50 5 50 50 50 50
OH OK OR PA PR RI SC SD TN TX UT VA VT WA WI WV WY
21 13 50 50 35 7 50 16 11 50 50 50 4 50 50 2 50
```

DAPC

DAPC analysis with state as group

Hide


```
#change genlight object to a matrix
covidSNPsMat <- as.matrix(covidSNPs)
# replace any missing genotypes with the mean
covidSNPsImpute <- na.replace(covidSNPsMat, 0, na.rm = TRUE)
#use pop() to get group membership information for all individuals
group <- pop(covidSNPs)
#now use xvalDapc() with the group membership info from above, and
#a maximum of 100 PCs and 100 repetitions
xval<-xvalDapc(covidSNPsImpute, grp = group, n.pca.max=100, n.rep=100, parallel = "multicore", n
cpus = 8)
#look at the items two to six in the output
xval[2:6]
```

```
$`Median and Confidence Interval for Random Chance`
      2.5%      50%      97.5%
0.01436781 0.02052345 0.02708827

$`Mean Successful Assignment by Number of PCs of PCA`
      10      20      30      40      50      60      70      80
0.06102281 0.10204658 0.12926474 0.13867179 0.15536037 0.16511666 0.17612022 0.18371017
      90
0.19302710

$`Number of PCs Achieving Highest Mean Success`
[1] "90"

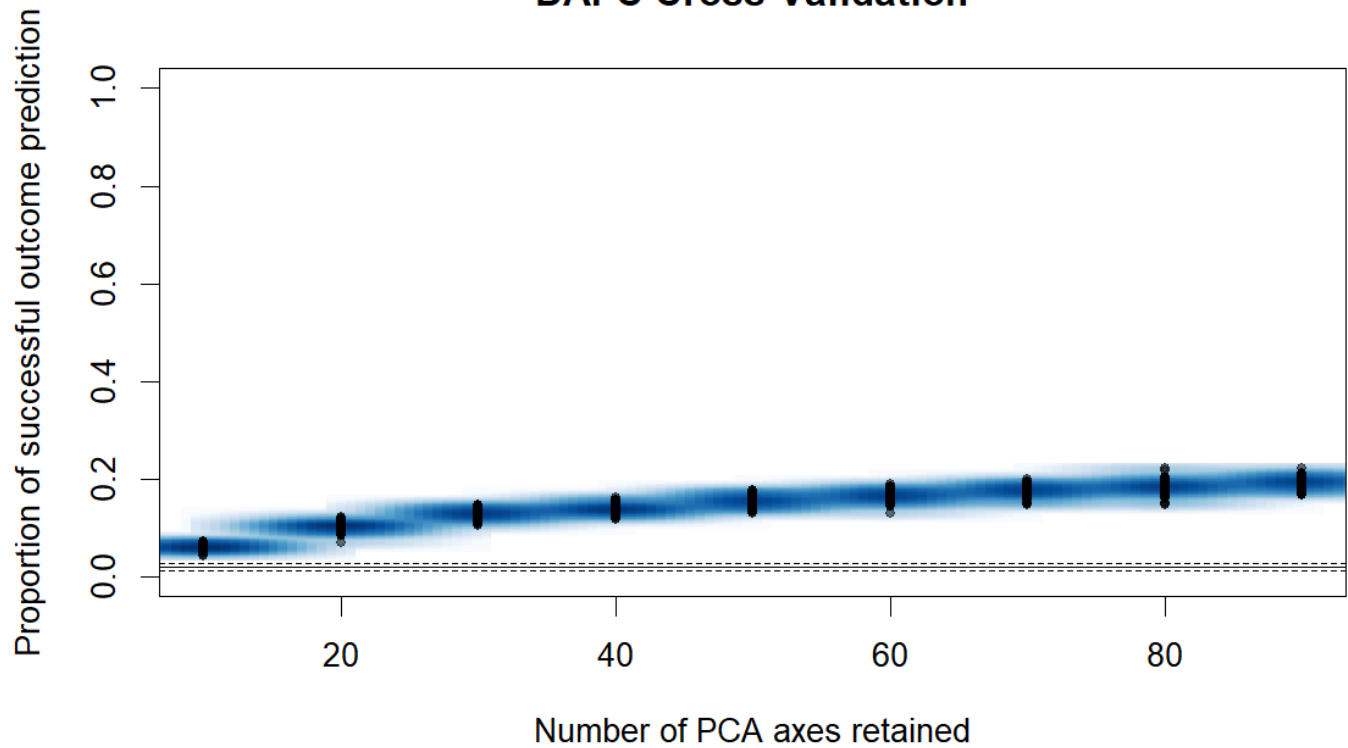
$`Root Mean Squared Error by Number of PCs of PCA`
      10      20      30      40      50      60      70      80      90
0.9390014 0.8979951 0.8707858 0.8613778 0.8447108 0.8349504 0.8239554 0.8163831 0.8070508

$`Number of PCs Achieving Lowest MSE`
[1] "90"
```

[Hide](#)

```
covidDAPC <- xval[7]$DAPC
#plot
par(mfrow = c(2,2))
```

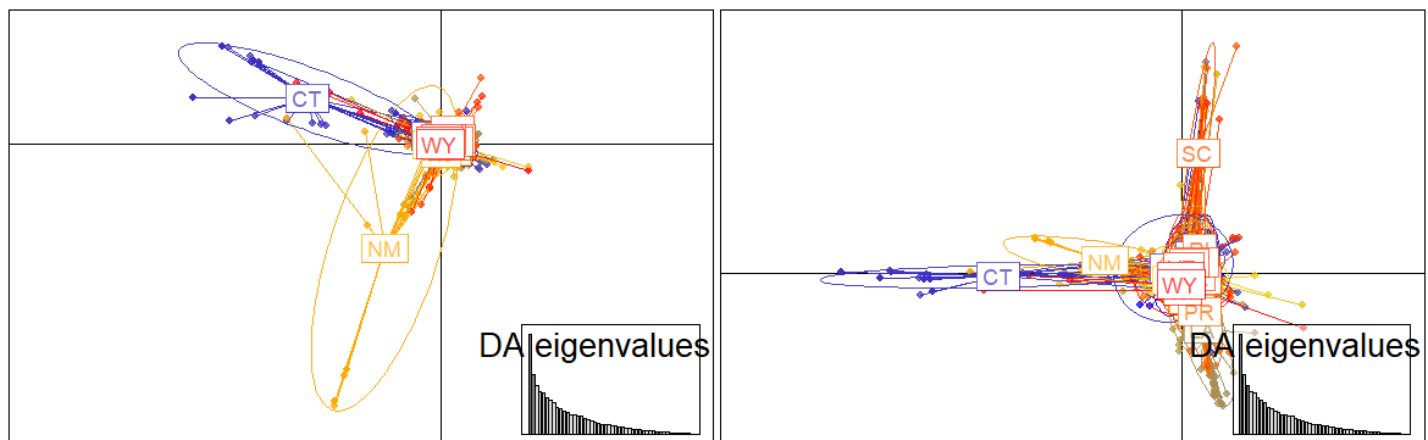
DAPC Cross-Validation

[Hide](#)

```
scatter(covidDAPC, xax = 1, yax = 2) #use discriminant functions 1 for the x-axis & 2 for the y-axis  
scatter(covidDAPC, xax = 1, yax = 3) #use discriminant functions 1 for the x-axis & 3 for the y-axis
```

[Hide](#)

```
compoplot(covidDAPC)
```



Save plot

[Hide](#)

```
pdf("./plot/DAPC.pdf", width = 8, height = 5)
par(mfrow = c(2,2))
scatter(covidDAPC, xax = 1, yax = 2)
scatter(covidDAPC, xax = 1, yax = 3)
```

[Hide](#)

```
dev.off()
```

```
null device
      1
```

[Hide](#)

```
pdf("./plot/DAPC_composition.pdf", width = 10, height = 4)
compoplot(covidDAPC)
```

[Hide](#)

```
dev.off()
```

```
null device
      1
```

Clustering with no prior group

[Hide](#)

```
#run find.clusters() on the raw data
#use max.n.clust = 100
cluster <- find.clusters(covidSNPs, max.n.clust = 300, n.pca = 200)
```

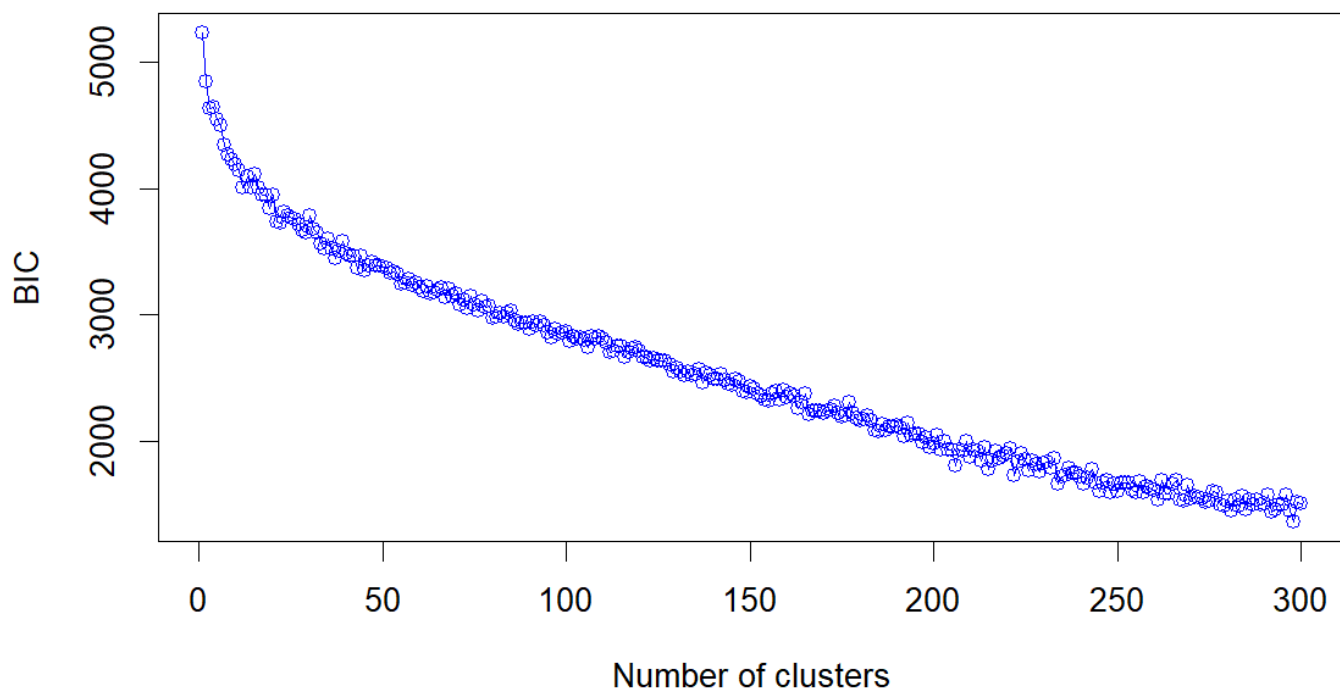
Warning: did not converge in 100000 iterations

Choose the number of clusters (≥ 2):

[Hide](#)

100

**Value of BIC
versus number of clusters**



No optimal K found, skip DAPC without prior group

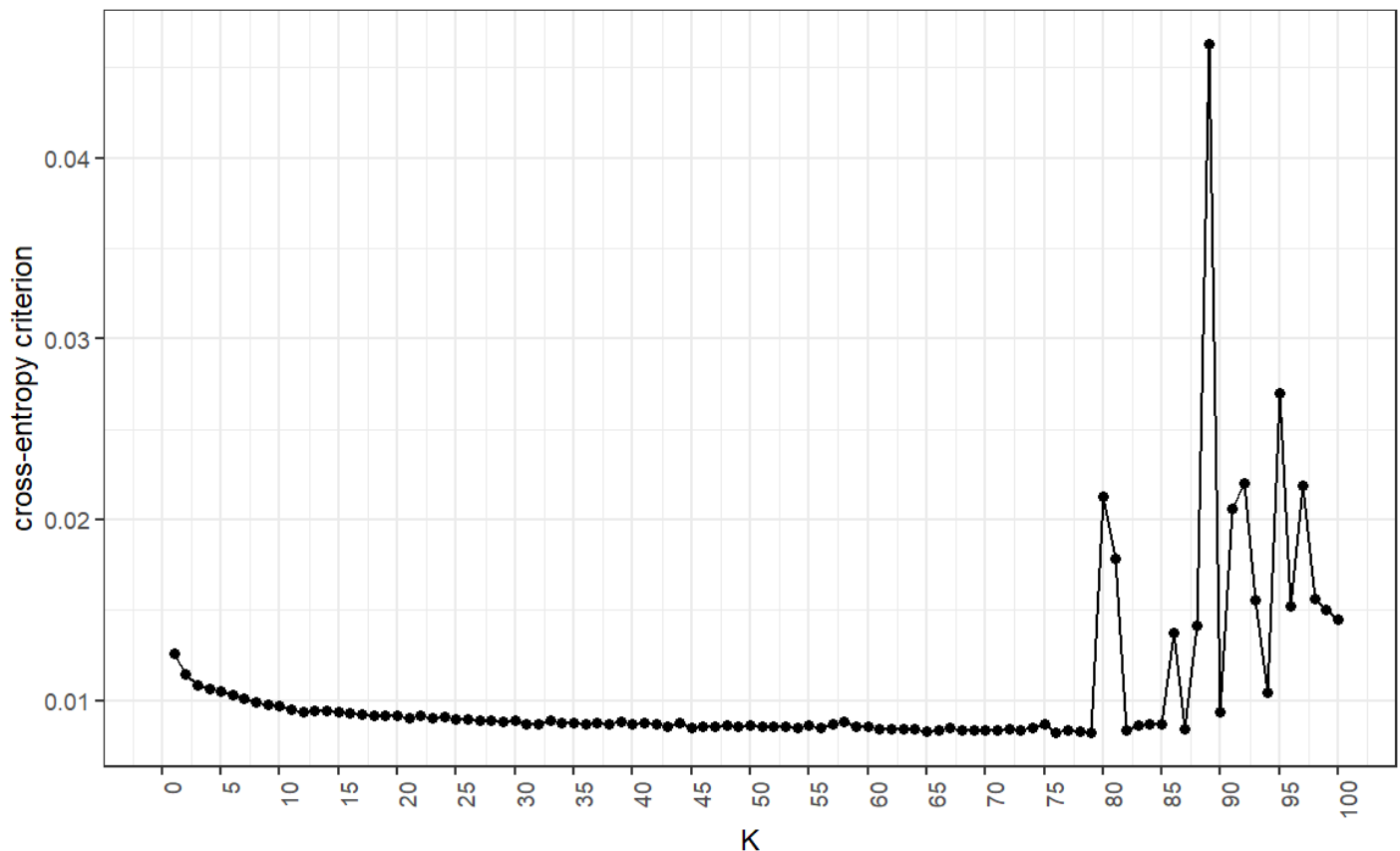
Admixture modeling

[Hide](#)

```
#save.image(file = "./SNP_analysis.rdata")
```

Hide

```
# plot cross-entropy criterion for each K
ces_sum <- data.frame()
ce_sum <- summary(aa_res)$crossEntropy %>%
  as.data.frame()
colnames(ce_sum) <- 1:100
ce_sum <- t(ce_sum) %>%
  as.data.frame()
ce_sum$K <- 1:100
ggplot(ce_sum) +
  geom_point(aes(x = K, y = mean)) +
  geom_path(aes(x = K, y = mean)) +
  scale_x_continuous(limits=c(0,100), breaks=seq(0,100, by = 5)) +
  labs(x = "K", y = "cross-entropy criterion")+
  theme_bw() +
  theme(axis.text.x = element_text(angle = 90, hjust = 0.95))
ggsave(filename = "./plot/cross_entropy.pdf", width = 8, height = 3, units = "in")
```



Hide

```
# choose k =5
# Pick the K where cross-entropy stops dropping sharply
cross.entropy(aa_res,5)
```

```
K = 5
run 1 0.01053063
```

Hide

```
q_k5 <- Q(aa_res, 5, 1) %>% as.data.frame()
# assign group
sample_size <- info %>% group_by(state) %>% summarise("size" = n())
sample_size <- arrange(sample_size, desc(size))
q_k5$state <- ped$V1
# set state as factor
q_k5$state <- as.factor(q_k5$state)
levels(q_k5$state)
```

```
[1] "reference" "AK"      "AL"      "AR"      "AZ"      "CA"      "CO"
[8] "CT"        "DE"      "FL"      "GA"      "HI"      "IA"      "ID"
[15] "IL"        "IN"      "KS"      "KY"      "LA"      "MA"      "MD"
[22] "ME"        "MI"      "MN"      "MO"      "MS"      "MT"      "NC"
[29] "NE"        "NH"      "NJ"      "NM"      "NV"      "NY"      "OH"
[36] "OK"        "OR"      "PA"      "PR"      "RI"      "SC"      "SD"
[43] "TN"        "TX"      "UT"      "VA"      "VT"      "WA"      "WI"
[50] "WV"        "WY"
```

Hide

```
# reorder based on state
q_k5$state <- factor(q_k5$state, levels = c("reference", as.character(sample_size$state)))
q_k5 <- q_k5[order(q_k5$state),]
# remove last column
q_k5$state <- NULL
# store x break info
sample_size <- rbind(c("WUHAN", 1), sample_size)
sample_size$size <- as.numeric(sample_size$size)
sample_size$xbreak <- NA
for(i in 1:nrow(sample_size)){
  sample_size$xbreak[i] <- sum(sample_size$size[1:i])
}

#plot
n <- 5
qual_col_pals = brewer.pal.info
col_vector = unlist(mapply(brewer.pal, qual_col_pals$maxcolors, rownames(qual_col_pals)))
colors <- sample(col_vector, n)

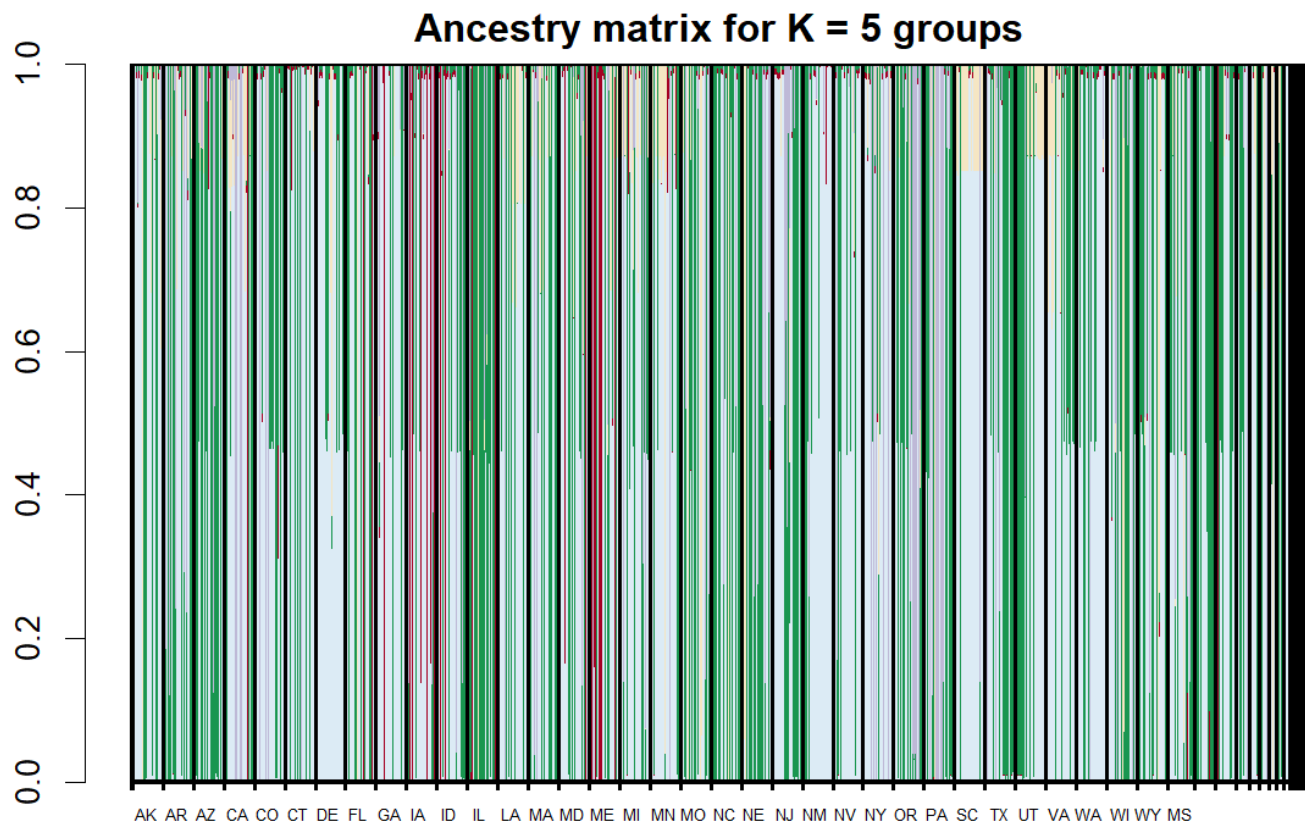
gap <- (0.845-0.04)/33
par(mar = c(2, 2, 2, 2))
barplot(t(q_k5), border = NA, space = 0, col = colors,
        ylab = "Ancestry proportions", xlab = "Individuals",
        main = "Ancestry matrix for K = 5 groups", xaxt = 'n')
# add boxes around each sampling site
abline(v = c(0, sample_size$xbreak), lwd = 2)
```

Hide

```
segments(0, 0, 1931, 0, lwd = 3)
segments(0, 1, 1931, 1, lwd = 3)
```

Hide

```
for (i in 2:35){
  mtext(as.character(sample_size$state[i]), 1, adj = 0.04 + (i-2)*gap, cex = 0.5)
}
mtext("MS", 1, adj = 0.87, cex = 0.5)
```



Hide

```
pdf("./plot/Admixture.pdf", width = 18, height = 5)
gap <- (0.845-0.04)/33
par(mar = c(2, 2, 2, 2))
barplot(t(q_k5), border = NA, space = 0, col = colors,
        ylab = "Ancestry proportions", xlab = "Individuals",
        main = "Ancestry matrix for K = 5 groups", xaxt = 'n')
abline(v = c(0, sample_size$xbreak), lwd = 2)
segments(0, 0, 1931, 0, lwd = 3)
segments(0, 1, 1931, 1, lwd = 3)
for (i in 2:35){
  mtext(as.character(sample_size$state[i]), 1, adj = 0.04 + (i-2)*gap, cex = 0.5)
}
mtext("MS", 1, adj = 0.87, cex = 0.5)
dev.off()
```

```
null device
1
```

Hide

```
save.image(file = "./SNP_analysis.rdata")
```