

# Practical 03 SG: LD and Haplotype estimation

Kathryn Weissman & Irene Fernández

2022-11-28

## LD (15p.)

1. The file **FOXP2.zip** contains genetic information of individuals of a Japanese population of unrelated individuals. The genotype information concerns SNPs of the Forkhead box protein P2 (FOXP2) gene region, located the long arm of chromosome number 7. This gene plays an important role in the development of speech and language. The **FOXP2.zip** file contains:
  - **FOXP2.dat**: a text file with the genotype data which can be read in with R.
  - **FOXP2.fam**: a PLINK file with data on the individuals (family id, individual id, ids of parents, sex and phenotype).
  - **FOXP2.bed**: a PLINK file with binary genotype data.
  - **FOXP2.bim**: a PLINK file with data on the genetic variants (chromosome, SNP identifier, basepair position along the chromosome and alleles).
2. (1p) Load the **FOXP2.dat** file into the R environment. How many individuals and how many SNPs are there in the database? What percentage of the data is missing?
  - There are 104 individuals in the database.
  - There are 543 SNPs in the database.
  - 0% of the data is missing.

```
file <- "~/Downloads/FOXP2.dat"
data <- fread(file, header = TRUE, index="id")
data[1:10, 1:10]
```

```
##           id rs34684677 rs1839115 rs4727804 rs4727805 rs200888633 rs12534908
## 1: NA18939      T/G      C/T      G/A      T/G      T/G      G/A
## 2: NA18940      G/G      T/T      A/A      G/G      T/G      A/A
## 3: NA18941      G/G      T/T      A/A      G/G      T/G      A/A
## 4: NA18942      G/G      T/T      A/A      G/G      T/T      A/A
## 5: NA18943      G/G      T/T      A/A      G/G      T/T      A/A
## 6: NA18944      T/T      C/C      G/G      T/G      G/G      G/G
## 7: NA18945      G/G      T/T      A/A      G/G      G/G      A/A
## 8: NA18946      T/G      C/T      G/A      G/G      G/G      G/A
## 9: NA18947      T/G      C/T      G/A      G/G      T/G      G/A
## 10: NA18948     G/G      T/T      A/A      G/G      G/G     A/A
##          rs12533049 rs77861356 rs6945561
## 1:      C/T      T/T      C/T
## 2:      T/T      T/T      T/T
```

```

## 3:      T/T      T/T      T/T
## 4:      T/T      T/T      T/T
## 5:      T/T      T/T      T/T
## 6:      C/C      T/T      C/C
## 7:      T/T      T/T      T/T
## 8:      C/T      T/T      C/T
## 9:      C/T      T/T      C/T
## 10:     T/T      A/T      T/T

n <- nrow(data); n # number of samples

## [1] 104

p <- ncol(data) - 1; p # number of SNPs - first column is id

## [1] 543

perc.mis <- 100*sum(is.na(data))/(n*p); perc.mis # percentage of missing data overall

## [1] 0

```

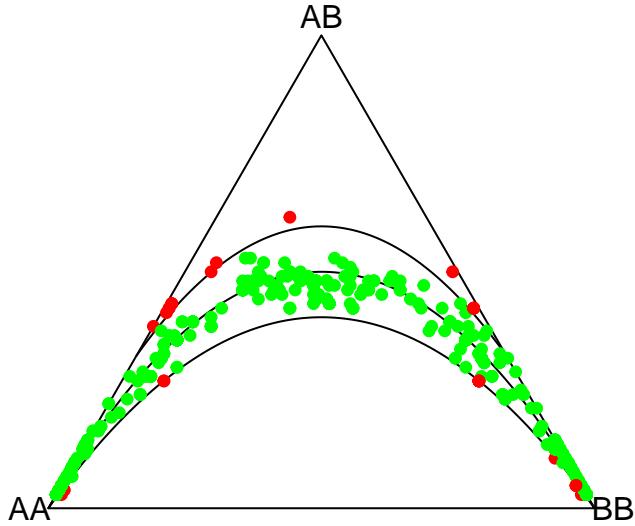
3. (1p) Determine the genotype counts for each SNP, and depict all SNPs simultaneously in a ternary plot, and comment on your result. For how many variants do you reject Hardy-Weinberg equilibrium using an ordinary chi-square test without continuity correction? (hint: you can read the .bim in R in order to determine the alleles of each SNP, and use function `MakeCounts` from the `HardyWeinberg` package to create a matrix of genotype counts).
- The green dots represent SNPs that are in equilibrium, and the red dots represent SNPs that are out of equilibrium. in the data set where the B allele frequency is greater than 0.5.
- Using an ordinary chi-square test we reject Hardy-Weinberg equilibrium for 33 variants.
- *Note:* In order to create the genotype counts matrix, we create a function similar to the one of the anterior practical.

```

genotype_matrix <- matrix(ncol = 3, nrow = 0, dimnames = list(NULL, c("AA", "AB", "BB")))
for(i in colnames(data)[-1]){
  snp <- genotype(data[[i]], reorder = "ascii") # reorder is to use as specific order for the alleles
 .snp.counts <- c(sum(snp==paste(allele.names(snp)[1], allele.names(snp)[1], sep = "/")),
                 sum(snp==paste(allele.names(snp)[1], allele.names(snp)[2], sep = "/")),
                 sum(snp==paste(allele.names(snp)[2], allele.names(snp)[2], sep = "/")))
  genotype_matrix <- rbind(genotype_matrix, .snp.counts)
}
rownames(genotype_matrix) <- colnames(data[-1])

ternaryPlot <- HWTernaryPlot(genotype_matrix)

```



```
results_ChiSq <- HWChisqStats(genotype_matrix, pvalues=TRUE)
sum(results_ChiSq <= .05)
```

```
## [1] 33
```

4. (1p) Using the function LD from the `genetics` package, compute the LD statistic  $D$  for the SNPs rs34684677 and rs2894715 of the database. Is there significant association between the alleles of these two SNPs?

Yes, there is association between the alleles because the p-value is small and  $D'$  is close to 1.

```
LD(genotype(data$rs34684677), genotype(data$rs2894715))
```

```
##
## Pairwise LD
## -----
##          D      D'      Corr
## Estimates: -0.05493703 0.9986536 -0.3144048
## 
##          X^2      P-value     N
## LD Test: 20.56088 5.77645e-06 104
```

5. (2p) Also compute the LD statistic  $D$  for the SNPs rs34684677 and rs998302 of the database. Is there significant association between these two SNPs? Is there any reason why rs998302 could have stronger or weaker correlation than rs2894715?

These two SNPs are not significantly associated because the p-value is large. rs34684677 and rs998302 may have a weaker correlation because the sites are located further from each other or they don't share any functional relationship.

```
LD(genotype(data$rs34684677), genotype(data$rs998302))
```

```
##  
## Pairwise LD  
## -----  
##           D       D'      Corr  
## Estimates: 0.007208888 0.1792444 0.09112725  
##  
##           X^2   P-value   N  
## LD Test: 1.727268 0.1887601 104
```

6. (2p) Given your previous estimate of  $D$  for SNPs rs34684677 and rs2894715, infer the haplotype frequencies. Which haplotype is the most common?

```
SNP1 <- genotype(data$rs34684677)  
SNP2 <- genotype(data$rs2894715)  
summary(SNP1)
```

```
##  
## Number of samples typed: 104 (100%)  
##  
## Allele Frequency: (2 alleles)  
##   Count Proportion  
##   G    174      0.84  
##   T     34      0.16  
##  
##  
## Genotype Frequency:  
##   Count Proportion  
##   G/G     73      0.70  
##   G/T     28      0.27  
##   T/T      3      0.03  
##  
## Heterozygosity (Hu) = 0.2748049  
## Poly. Inf. Content = 0.2360871
```

```
summary(SNP2)
```

```
##  
## Number of samples typed: 104 (100%)  
##  
## Allele Frequency: (2 alleles)  
##   Count Proportion  
##   T    138      0.66  
##   G     70      0.34  
##  
##  
## Genotype Frequency:
```

```

##      Count Proportion
## T/T     46      0.44
## T/G     46      0.44
## G/G     12      0.12
##
## Heterozygosity (Hu) = 0.4487179
## Poly. Inf. Content = 0.3468524

LD_res <- LD(genotype(data$rs34684677), genotype(data$rs2894715))

pA <- summary(SNP1)$allele.freq[,2][1]
pB <- summary(SNP2)$allele.freq[,2][1]
pa <- summary(SNP1)$allele.freq[,2][2]
pb <- summary(SNP2)$allele.freq[,2][2]

pAB <- LD_res$D + pA*pB #D = pAB - pApB
pab <- LD_res$D + pa*pb #D = pab - papb

```

- The haplotype GT have a frequency of 0.5000741.
- The haplotype TG have a frequency of  $7.4065053 \times 10^{-5}$ .
- As we could observe, the most common haplotype is the GT.

7. (2p) Compute the LD statistics  $R^2$  for all the marker pairs in this data base, using the LD function of the packages **genetics**. Be prepared that this make take a few minutes. Also compute an alternative estimate of  $R^2$  obtained by using the PLINK program. For this purpose you should:

- Download and install PLINK 1.90 from <https://www.cog-genomics.org/plink2/>
- Take care to store the files **FOXP2.bim**, **FOXP2.fam** and **FOXP2.bed** in a directory where PLINK can find them.
- Compute LD estimates with PLINK using `plink --bfile FOXP2 --r2 --matrix --out FOXP2` This creates a file with extension **FOXP2.ld** that contains a matrix with all  $R^2$  statistics. Read this file into the R environment. Make a scatter plot for R's LD estimates against PLINK's LD estimates. Are they identical or do they at least correlate? What's the difference between these two estimators? Which estimator would you prefer and why?
- As the plot shows, the  $R^2$  statistics are very similar with a correlation of 0.99. If the results do not differ, we prefer PLINK because is faster in the calculation.

```

# Convert one SNP at a time to genotype and save it as a column in dataframe.
# RES <- data.frame(genotype(data[[2]],sep="/"))
# for(i in 3:ncol(data)) {
#   snp <- genotype(data[[i]],sep="/")
#   RES <- cbind(RES,snp)
# }

#output <- LD(RES)

#R2 <- output$"R^2"
#write.table(R2, file="FOXP2_R2_r.txt", row.names=TRUE, col.names=TRUE)

```

```

R2 <- read.table("FOXP2_R2_r.txt", row.names = NULL)
#head(R2)
R2 <- R2[, -1]
dim(R2)

## [1] 543 543

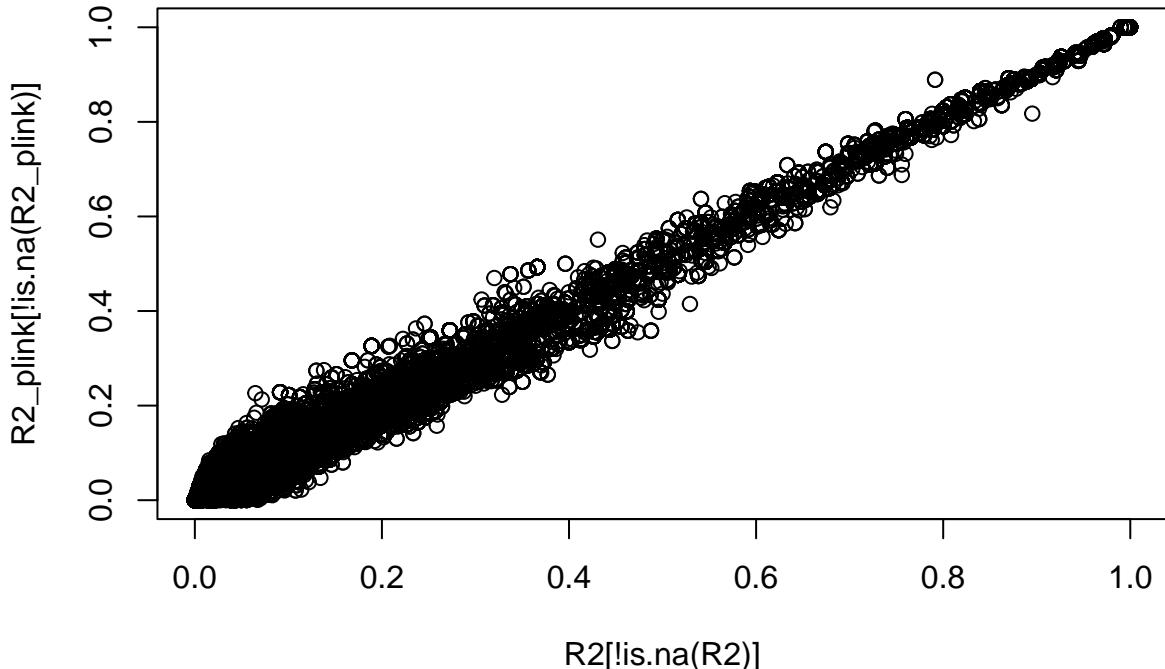
R2_plink <- read.table("FOXP2.ld")
#head(R2_plink)
dim(R2_plink)

## [1] 543 543

for(i in 1:ncol(R2_plink)) {
  for(j in 1:i) {
    R2_plink[i,j] = NA
  }
}

plot(R2[!is.na(R2)], R2_plink[!is.na(R2_plink)])

```



```

cor(R2[!is.na(R2)], R2_plink[!is.na(R2_plink)])

## [1] 0.9948306

```

8. (2p) Compute a distance matrix with the distance in base pairs between all possible pairs of SNPs, using the basepair position of each SNP given in the .bim file. Make a plot of R's  $R^2$  statistics against the distance (expressed as the number of basepairs) between the markers. Comment on your results.

```
dist_data <- fread("~/Downloads/FOXP2.bim")
head(dist_data)

##   V1          V2  V3          V4  V5  V6
## 1: 7  rs34684677 0 114400288  T  G
## 2: 7  rs1839115 0 114400872  C  T
## 3: 7  rs4727804 0 114400999  G  A
## 4: 7  rs4727805 0 114401123  T  G
## 5: 7 rs200888633 0 114401388  T  G
## 6: 7 rs12534908 0 114402258  G  A

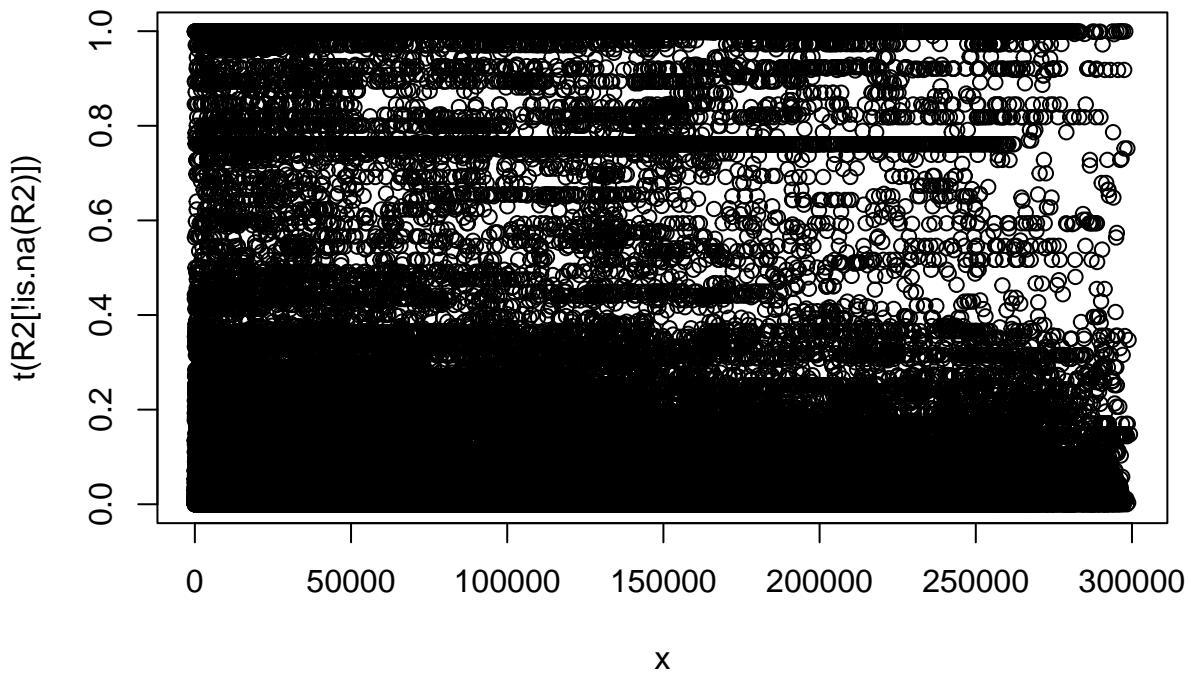
dist_data <- dist_data[match(colnames(data)[-1], dist_data$V2)]
dist_data <- dist_data[, 4]
x <- dist(dist_data)
head(x)

## [1] 584 711 835 1100 1970 2045

length(x)

## [1] 147153

plot(x, t(R2[!is.na(R2)]))
```



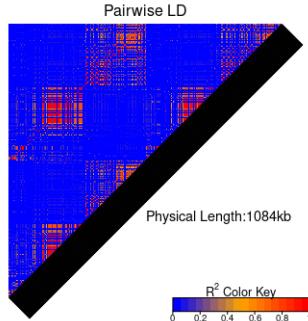
9. (2p) Make an LD heatmap of the markers in this database, using the  $R^2$  statistic with the LD function. Make another heatmap obtained by filtering out all variants with a MAF below 0.35, and redoing the computations to obtain the  $R^2$  statistics in R. Can you explain any differences observed between the two heatmaps?

- In the first LDheatmap there is more uncorrelation in general, this could be because we are keeping those SNPs with lower maf and these SNPs not provide much variation to have a correlation with other SNPs.

```

RES <- data.frame(genotype(data[[2]], sep="/"))
for(i in 3:ncol(data)) {
  snp <- genotype(data[[i]], sep="/")
  RES <- cbind(RES, snp)
}
rgb.palette <- colorRampPalette(rev(c("blue", "orange", "red")), space = "rgb")
#LDheatmap(gdat=RES,color=rgb.palette(18), LDmeasure = 'r')
include_graphics("LDheatmap.png", dpi = 200)

```

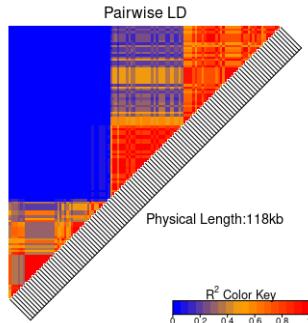


```

maf <- function(x){
  x <- genotype(x,sep="/")
  out <- summary(x)
  af1 <- min(out$allele.freq[,2],na.rm=TRUE)
  af1[af1==1] <- 0
  return(af1)
}
mafs <- apply(data[,-1], 2, maf)
sub_snp <- mafs[mafs > 0.35]

data_maf <- as.data.frame(data)
data_maf <- data_maf[, names(sub_snp)]
RES_maf <- data.frame(genotype(data_maf[[1]],sep="/"))
for(i in 2:ncol(data_maf)) {
  snp <- genotype(data_maf[[i]],sep="/")
  RES_maf <- cbind(RES_maf,snp)
}
rgb.palette <- colorRampPalette(rev(c("blue", "orange", "red")), space = "rgb")
#LDheatmap(gdat=RES_maf,color=rgb.palette(18))
include_graphics("LDheatmap_maf.png", dpi = 200)

```



10. (1p) Can you distinguish blocks of correlated markers in the area of the FOXP2 gene? How many blocks do you think that *at least* seem to exist?
  - In the second plot we could distinguish 3 main blocks with high correlation, two of them also being correlated but in a weaker manner.
11. (1p) Simulate independent SNPs under the assumption of Hardy-Weinberg equilibrium, using R's `sample` instruction `sample(c("AA","AB","BB"),n,replace=TRUE,prob=c(p*p,2*p*q,q*q))`. Simulate as many SNPs as you have in your database, and take care to match each SNP in your database

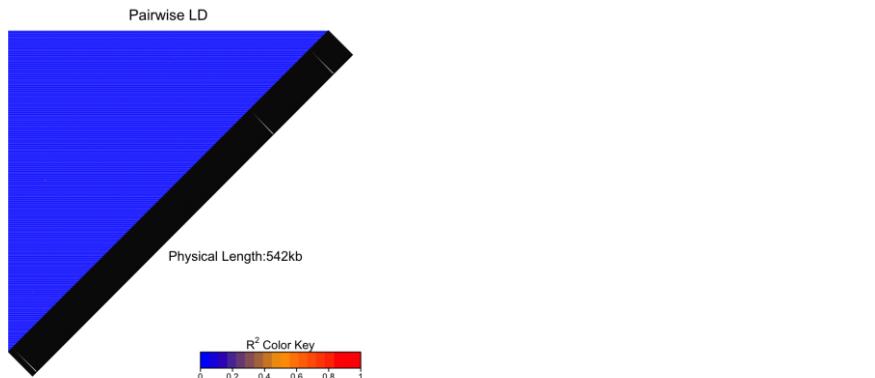
with a simulated SNP that has the same sample size and allele frequency. Make an LD heatmap of the simulated SNPs, using  $R^2$  as your statistic. Compare the results with the LD heatmap of the FOXP2 region. What do you observe? State your conclusions.

```

f_sample_data <- function(x) {
  x = genotype(x, sep = "/")
  out = summary(x)
  p = out$allele.freq[1,2]
  q = out$allele.freq[2,2]
  A = out$allele.names[1]
  B = out$allele.names[2]
  AA = paste(A, A, sep = "/")
  AB = paste(A, B, sep = "/")
  BB = paste(B, B, sep = "/")
  return(sample(c(AA,AB,BB), nrow(data), replace=TRUE, prob=c(p*p,2*p*q,q*q)))
}
sample_data = data.frame(apply(as.data.frame(data[,-1]), 2, f_sample_data))
#head(sample_data)

RES_sample <- data.frame(genotype(sample_data[,1],sep="/"))
for(i in 2:ncol(sample_data)) {
 .snp <- genotype(sample_data[,i],sep="/")
  RES_sample <- cbind(RES_sample,.snp)
}
#LDheatmap(gdat=RES_sample,color=rgb.palette(18), LDmeasure = 'r')
include_graphics("LDheatmap_sample.png", dpi = 200)

```



## Haplotype estimation (10p.)

- Apolipoprotein E (APOE) is a protein involved in Alzheimer's disease. The corresponding gene *APOE* has been mapped to chromosome 19. The file *APOE.dat* contains genotype information of unrelated individuals for a set of SNPs in this gene. Load this data into the R environment. *APOE.zip* contains the corresponding *.bim*, *.fam* and *.bed* files. You can use the *.bim* file to obtain information about the alleles of each polymorphism.

```

rm(list=ls())
file <- "~/Downloads/APOE.dat"
data <- as.data.frame(fread(file, header = TRUE))
data[1:10, 1:10]

##          id rs112748686 rs374311741 rs569874826 rs538789834 rs144831227
## 1  NA20502      C/C      C/C      C/C      G/G      C/C
## 2  NA20503      C/C      C/C      C/C      G/G      C/C
## 3  NA20504      C/C      C/C      C/C      G/G      C/C
## 4  NA20505      C/C      C/C      C/C      G/G      C/C
## 5  NA20506      C/C      C/C      C/C      G/G      C/C
## 6  NA20507      C/C      C/C      C/C      G/G      C/C
## 7  NA20508      C/C      C/C      C/C      G/G      C/C
## 8  NA20509      C/C      C/C      C/C      G/G      C/C
## 9  NA20510      C/C      C/C      C/C      G/G      C/C
## 10 NA20511      C/C      C/C      C/C      G/G      C/C
##          rs148595630 rs534353668 rs551394116 rs892593
## 1      G/G      C/C      G/G      G/G
## 2      G/G      C/C      G/G      G/G
## 3      G/G      C/C      G/G      G/G
## 4      G/G      C/C      G/G      G/G
## 5      G/G      C/C      G/G      G/G
## 6      G/G      C/C      G/G      C/C
## 7      G/G      C/C      G/G      C/G
## 8      A/G      C/C      G/G      G/G
## 9      G/G      C/C      G/G      C/G
## 10     G/G      C/C      G/G      G/G

data_id <- data$id
data <- data[,-1]

```

2. (1p) How many individuals and how many SNPs are there in the database? What percentage of the data is missing?

- There are 107 individuals in the database.
- There are 161 SNPs in the database.
- 0% of the data is missing.

```

n <- nrow(data); n # number of samples

## [1] 107

p <- ncol(data); p # number of SNPs - first column is id

## [1] 162

perc.mis <- 100*sum(is.na(data))/(n*p); perc.mis # percentage of missing data overall

## [1] 0

```

3. (1p) Assuming all SNPs are bi-allelic, how many haplotypes can theoretically be found for this data set?

- Theoretically  $5.846007e+48$  haplotypes can be found for this dataset.

$2^p$

```
## [1] 5.846007e+48
```

4. (2p) Estimate haplotype frequencies using the `haplo.stats` package (set the minimum posterior probability to 0.001). How many haplotypes do you find? List the estimated probabilities in decreasing order. Which haplotype number is the most common?

- We find 31 haplotypes.
- The most common haplotype is the number 27.

```
Geno <- cbind(substr(data[,1],1,1),substr(data[,1],3,3))
for(i in 2:ncol(data)) {
  Geno <- cbind(Geno,substr(data[,i],1,1),substr(data[,i],3,3))
}
Haplo.Res <- haplo.em(Geno,locus.label=colnames(data),control=haplo.em.control(min.posterior=0.001))
#Haplo.Res
Haplo.Res$nreps
```

```
## indx.subj
##   1   2   3   4   5   6   7   8   9   10  11  12  13  14  15  16  17  18  19  20
##   1   1   2   1   1   1   1   1   1   1   1   1   1   1   1   1   1   2   1   1   1
##  21  22  23  24  25  26  27  28  29  30  31  32  33  34  35  36  37  38  39  40
##   2   2   1   1   1   2   1   2   1   1   1   1   2   1   1   1   1   1   1   1   1   2
##  41  42  43  44  45  46  47  48  49  50  51  52  53  54  55  56  57  58  59  60
##   1   1   1   2   1   1   1   1   1   1   1   2   1   1   1   1   1   1   1   1   3   2
##  61  62  63  64  65  66  67  68  69  70  71  72  73  74  75  76  77  78  79  80
##   1   2   1   1   1   1   1   1   1   1   1   1   1   1   1   1   1   1   2   1   2
##  81  82  83  84  85  86  87  88  89  90  91  92  93  94  95  96  97  98  99  100
##   1   2   1   1   1   1   1   2   1   1   1   1   1   1   1   1   1   1   1   1   2   2
## 101 102 103 104 105 106 107
##   1   1   1   1   1   1   1
```

```
sort(Haplo.Res$hap.prob, decreasing = TRUE)
```

```
## [1] 0.3995237131 0.1308411215 0.0744796311 0.0683949710 0.0501815979
## [6] 0.0467289720 0.0358505426 0.0351610760 0.0225442605 0.0204954576
## [11] 0.0186915888 0.0161159156 0.0086862045 0.0073499717 0.0046728972
## [16] 0.0046728972 0.0046728972 0.0046728972 0.0046728972 0.0046728972
## [21] 0.0046728972 0.0046728972 0.0046728972 0.0046728972 0.0040259704
## [26] 0.0034124961 0.0033023394 0.0028690284 0.0021386625 0.0016572819
## [31] 0.0008202255
```

```
which(Haplo.Res$hap.prob == max(Haplo.Res$hap.prob))
```

```
## [1] 27
```

5. (2p) Is the haplotypic constitution of any of the individuals in the database ambiguous or uncertain? For how many? What is the most likely haplotypic constitution of individual NA20763? (identify the constitution by the corresponding haplotype numbers).

- There are 19 ambiguous constitutions.
- The most likely haplotypic constitution of individual NA20763 is 8.

```
sum(Haplo.Res$nreps > 1)

## [1] 19

Haplo.Res$hap1code[which(data_id == "NA20763")]

## [1] 8
```

6. (1p) Suppose we would delete polymorphism rs374311741 from the database prior to haplotype estimation. Would this affect the results obtained? Justify your answer.

- Deleting one polymorphism generates very minimal variations ( $\text{var} = 0.005$ ) in the estimated probabilities and there is still 31 haplotypes.

```
data_sub <- data[, !(colnames(data) %in% c("rs374311741"))]
Geno_sub <- cbind(substr(data_sub[,1],1,1),substr(data_sub[,1],3,3))
for(i in 2:ncol(data_sub)) {
  Geno_sub <- cbind(Geno_sub,substr(data_sub[,i],1,1),substr(data_sub[,i],3,3))
}
Haplo.Res_sub <- haplo.em(Geno_sub,locus.label=colnames(data_sub),control=haplo.em.control(min.posterior=0.005))
#Haplo.Res_sub
Haplo.Res_sub$nreps
```

```
## indx.subj
##   1   2   3   4   5   6   7   8   9   10  11  12  13  14  15  16  17  18  19  20
##   1   1   2   1   1   1   1   1   1   1   1   1   1   1   1   1   1   2   1   1   1
##  21  22  23  24  25  26  27  28  29  30  31  32  33  34  35  36  37  38  39  40
##   2   2   1   1   1   2   1   2   1   1   1   1   2   1   1   1   1   1   1   1   1   2
##  41  42  43  44  45  46  47  48  49  50  51  52  53  54  55  56  57  58  59  60
##   1   1   1   2   1   1   1   1   1   1   1   2   1   1   1   1   1   1   1   1   3   2
##  61  62  63  64  65  66  67  68  69  70  71  72  73  74  75  76  77  78  79  80
##   1   2   1   1   1   1   1   1   1   1   1   1   1   1   1   1   1   1   1   2   1   2
##  81  82  83  84  85  86  87  88  89  90  91  92  93  94  95  96  97  98  99  100
##   1   2   1   1   1   1   1   2   1   1   1   1   1   1   1   1   1   1   1   2   1   2
## 101 102 103 104 105 106 107
##   1   1   1   1   1   1   1
```

```
var(Haplo.Res_sub$hap.prob, Haplo.Res$hap.prob)

## [1] 0.005461913
```

7. (1p) Remove all genetic variants that have a minor allele frequency below 0.10 from the database, and re-run `haplo.em`. How does this affect the number of haplotypes?

- Removing the variants with MAF < 0.10, we only remain with 21 variants. And the number of haplotypes decreases to 8.

```

maf <- function(x){
  x <- genotype(x, sep="/")
  out <- summary(x)
  af1 <- min(out$allele.freq[,2], na.rm=TRUE)
  af1[af1==1] <- 0
  af1
}

mafs <- apply(data, 2, maf)
data_maf <- data[, mafs >= 0.10]
dim(data_maf)

## [1] 107 21

Geno_maf <- cbind(substr(data_maf[,1],1,1), substr(data_maf[,1],3,3))
for(i in 2:ncol(data_maf)) {
  Geno_maf <- cbind(Geno_maf, substr(data_maf[,i],1,1), substr(data_maf[,i],3,3))
}

Haplo.Res_maf <- haplo.em(Geno_maf, locus.label=colnames(data_maf), control=haplo.em.control(min.posterior=0.01))
#Haplo.Res_maf
Haplo.Res_maf$nreps

## indx.subj
##   1   2   3   4   5   6   7   8   9   10  11  12  13  14  15  16  17  18  19  20
##   1   1   1   1   1   1   1   2   1   1   1   1   1   1   1   1   1   1   1   1   1
##  21  22  23  24  25  26  27  28  29  30  31  32  33  34  35  36  37  38  39  40
##   1   1   1   1   1   1   1   1   1   1   1   1   1   1   1   1   1   1   1   1   1
##  41  42  43  44  45  46  47  48  49  50  51  52  53  54  55  56  57  58  59  60
##   1   1   1   1   1   1   1   1   2   1   1   1   1   1   1   1   1   1   1   1   2
##  61  62  63  64  65  66  67  68  69  70  71  72  73  74  75  76  77  78  79  80
##   1   1   1   1   1   1   1   1   1   2   1   1   1   1   1   1   1   1   1   1   1
##  81  82  83  84  85  86  87  88  89  90  91  92  93  94  95  96  97  98  99  100
##   1   1   1   1   1   1   1   1   2   1   1   1   1   1   1   2   1   1   1   1   1
## 101 102 103 104 105 106 107
##   1   1   1   1   1   1   1

Haplo.Res_maf$hap.prob

## [1] 0.130841121 0.031850225 0.005532952 0.074766355 0.004672897 0.113009588
## [7] 0.018691589 0.620635272

```

8. (2p) We could consider the newly created haplotypes in our last run of `haplo.em` as the alleles of a new superlocus. Which is, under the assumption of Hardy-Weinberg equilibrium, the most likely genotype at this new locus? What is the probability of this genotype? Which genotype is the second most likely, and what is its probability?

- The most likely genotype is the one of the haplotype 36 with probability of 0.38.
- The second most likely genotype is the one of the haplotype 8 with probability of 0.16.

```
Haplo.Res_maf$haplotype
```

```
##   rs892593 rs892594 rs2722659 rs2722660 rs2571147 rs2571148 rs34762924
## 1      C      A      C      C      G      T      C
## 2      C      G      C      C      G      T      C
## 3      C      G      C      C      G      T      C
## 4      C      G      C      C      G      T      C
## 5      G      G      T      T      A      C      A
## 6      G      G      T      T      A      C      A
## 7      G      G      T      T      A      C      A
## 8      G      G      T      T      A      C      A
##   rs2571149 rs2722661 rs2571150 rs147663893 rs8102685 rs2571151 rs35570438
## 1      T      A      G      C      T      G      A
## 2      T      A      G      C      T      G      T
## 3      T      A      G      T      T      G      T
## 4      T      A      G      T      T      G      T
## 5      C      G      T      C      C      T      T
## 6      C      G      T      C      C      T      T
## 7      C      G      T      T      C      G      T
## 8      C      G      T      T      C      T      T
##   rs2571152 rs2571153 rs2722662 rs35391606 rs2437014 rs2437013 rs2722664
## 1      T      A      C      A      C      T      A
## 2      T      A      C      A      C      T      A
## 3      T      A      C      A      C      T      A
## 4      T      C      C      A      C      T      A
## 5      G      C      T      A      T      A      G
## 6      G      C      T      C      T      A      G
## 7      T      A      C      A      C      T      G
## 8      G      C      T      C      T      A      G
```

```
haplotypes <- Haplo.Res_maf$haplotype
probabilities <- Haplo.Res_maf$hap.prob
```

```
haplotype_pairs <- get.hapPair(haplotypes, probabilities, base.index=1)
```

```
genotype_prob <- haplotype_pairs$p.g
genotype1 <- which(genotype_prob == max(genotype_prob))
genotype2 <- which(genotype_prob == max(genotype_prob[-genotype1]))
```

```
haplotype_pairs$x.haplo[genotype1,]
```

```
## [1] 0 0 0 0 0 0 2
```

```
genotype_prob[genotype1]
```

```
## [1] 0.3851881
```

```
haplotype_pairs$x.haplo[genotype2,]
```

```
## [1] 0 0 0 0 0 0 1
```

```
genotype_prob[genotype2]
```

```
## [1] 0.1624092
```