

Name:

Perform the computations and make the graphics that are asked for in the practical below. You can write your answers on this sheet, and attach graphics at the end. Give each graphic a title, and clearly label x and y axes. Send your solution to ivan.galvan@upc.edu before the 10th of November 2016. You can make use of the R-package **genetics** (and other packages) to compute your answers. The first part of the practical is dedicated to the descriptive analysis of SNP data, whereas the second part is dedicated to the analysis of STR data. The datasets can be downloaded by clicking on their file names given below.

1 SNP dataset

1. The file CHBChr2.rda contains genotype information (10.000 SNPs) of 45 individuals of a Chinese population of unrelated individuals. Load this data into the R environment.

2. (1p) Recode all “NN” genotypes as missing values (NA). What percentage of the data is missing?

.....

3. (1p) For how many SNPs the genotype information is completely missing? Remove these SNPs from the database.

.....

4. (1p) What is, on the average, the percentage of missing information per individual, after fully missing SNPs have been removed?

.....

5. (1p) How many markers are monomorphic?

.....

6. (2p) Write a function to compute the minor allele frequency. Make sure the function also produces sensible answers for markers that consist of missing values only, or markers that are monomorphic. Include the source code of your function here.

.....

.....

.....

.....

-
-
7. (1p) Compute the minor allele frequencies for all markers, and make a histogram of it.
 8. (1p) What percentage of the markers have a maf below 0.05? And below 0.01?
.....
 9. (2p) Compute for each marker its **expected heterozygosity**, where the expected heterozygosity for a bi-allelic markers is defined as $1 - \sum_{i=1}^k p_i^2$, where p_i is the frequency of the i th allele. Compute the average expected heterozygosity over all markers. Make a histogram of the expected heterozygosity.

2 STR dataset

1. The file FrenchStrs.dat contains genotype information (STRs) of individuals from a French population. The first column of the data set contains an identifier the individual. STR data starts at the second column. Load this data into the R environment.
2. (1p) How many individuals and how many STRs contains the database?
.....
3. (1p) The value -9 indicates a missing value. Replace all missing values by NA. What percentage of the total amount of datavalues is missing?
.....
4. (2p) Write a function that determines the number of alleles for a STR. Determine the number of alleles for each STR in the database. Compute basic descriptive statistics of the number of alleles (mean, standard deviation, median, minimum, maximum).
.....
5. (2p) Make a boxplot and a histogram of the number of alleles per STR. What is the most common number of alleles for an STR?
.....
6. (2p) Compute the expected heterozygosity for each STR. Make a histogram of the expected heterozygosity over all STRS. Compute the average expected heterozygosity over all STRs.
.....

7. (2p) Compare the results you obtained for the SNP database with those you obtained for the STR database. What differences do you observe between these two types of genetic markers?

.....