

Relate practical: Using genealogies for population genetics

Leo Speidel^{1,2,*}

¹ Genetics Institute, University College London, London, UK

² Francis Crick Institute, London, UK

* Contact: leo.speidel@outlook.com

In this practical, we will infer genealogies for the Simons Genome Diversity Project dataset, downloaded from <https://reichdata.hms.harvard.edu/pub/datasets/sgdp/>. This dataset comprises whole-genome sequencing data of 278 modern humans with sampling locations shown in Fig. 1.

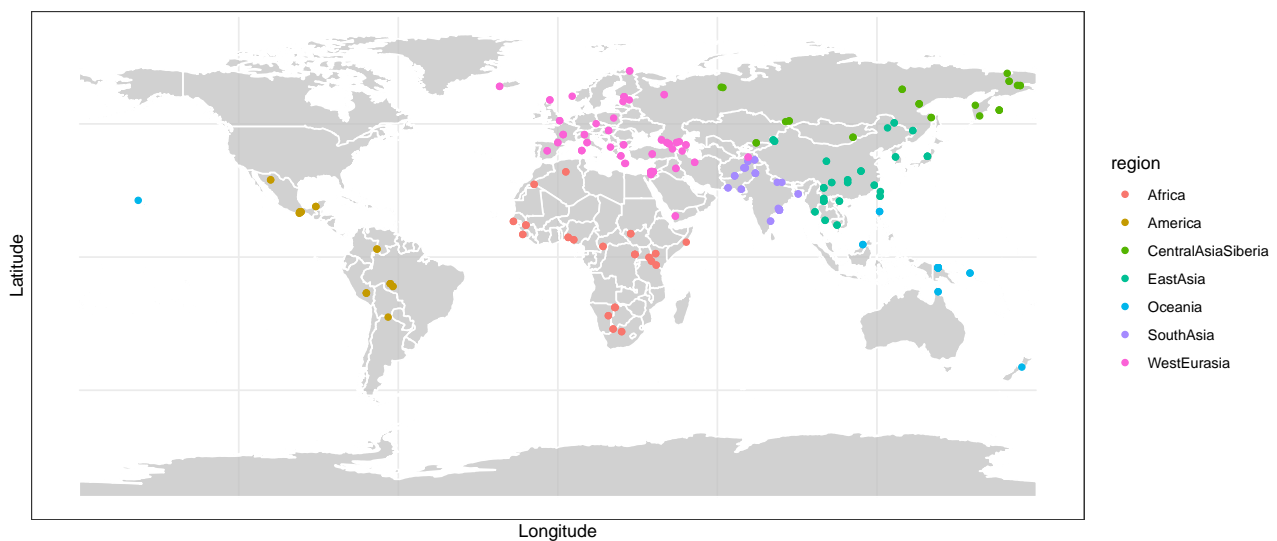


Figure 1: Sampling locations of the 278 modern humans in the Simons Genome Diversity Project. Samples are classified into seven regions shown by colours.

Note 1

The data was downloaded from:

- Phased genotypes: https://sharehost.hms.harvard.edu/genetics/reich_lab/sgdp/phased_data/PS2_multisample_public
- Genomic mask: https://reichdata.hms.harvard.edu/pub/datasets/sgdp/filters/all_samples/
- Human ancestral genome: ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase1/analysis_results/supporting/ancestral_alignments/
- Recombination maps: https://mathgen.stats.ox.ac.uk/impute/1000GP_Phase3.html

We then used `RelateFileFormats --mode ConvertFromVcf` to convert to the haps/sample file format and the `PrepareInputFiles.sh` script to make sure ancestral alleles are denoted by 0 and to filter out regions according to the genomic mask.

Overview

1	Relate	2
1.1	Data requirements and file formats	2
1.2	The arguments	3
2	Running Relate on the Simons Genome Diversity Project dataset	4
2.1	Running Relate	4
3	Effective population sizes and split times	6
3.1	Estimating population sizes given a genealogy	6
3.2	Joint fitting of population size and branch lengths	8
4	Detecting evidence for positive selection	9

1 Relate

Relate estimates the joint genealogies of many thousands of modern individuals genome-wide. These genealogies describe how individuals are related through their most-recent common ancestors back in time and can be seen as the genetic analogue of a family tree for unrelated individuals.

The output of Relate is a sequence of binary trees, each describing the genealogical relationships locally in that part of the genome. Neighbouring genealogical trees differ because of recombination events that change the genetic relationships of individuals.

Note 2

A detailed documentation for **Relate** is available at <https://myersgroup.github.io/relate>.

Details about the method is published in:

L. Speidel, M. Forest, S. Shi, S. R. Myers. A method for genome-wide genealogy estimation for thousands of samples. *Nature Genetics* **51**, 1321–1329 (2019).

For this practical, we will use **Relate binaries**, **bash scripts**, and **R**.

The Relate binaries can be found in my home directory or downloaded from

<https://myersgroup.github.io/relate>. It will be convenient to define a variable storing the location of the Relate binaries, e.g.,

```
#no need to copy this to your home directory
PATH_TO_RELATE="/home/speidel/relate_v1.1.6_x86_64_static/"
```

1.1 Data requirements and file formats

Relate uses the **haps/sample** file format (output file format of SHAPEIT2) as input (see https://myersgroup.github.io/relate/input_data.html#FileFormat).

- You can convert from **vcf** to haps/sample and from **hap/legend/sample** to haps/sample using functions provided by Relate.
(see https://myersgroup.github.io/relate/input_data.html#ConvertToHapsSample)

Note 3: Data requirements

- Please separate your data by chromosome.
- Relate assumes whole-genome sequencing data as input.
- The data needs to be phased. This can be done, for instance, using SHAPEIT.
- Ancestral and derived states at single-nucleotide polymorphisms (SNPs) need to be known.
- A genomic mask should be specified that filters out “bad” regions

We provide a script to prepare your data under

```
${PATH_TO_RELATE}/scripts/PrepareInputFiles/PrepareInputFiles.sh
```

(see https://myersgroup.github.io/relate/input_data.html#Prepare).

1.2 The arguments

Required arguments

<code>--mode</code>	Mode in which to run Relate. Mode "All" will execute all stages of the algorithm and other modes can be used to execute individual stages (see here). This is useful, for instance, when parallelizing Relate.
<code>-m,--mutation_rate</code>	Mutation rate per base per generation.
<code>-N,--effective_size</code>	Haploid effective population size. (For diploid organisms, multiply the effective population size of individuals by 2)
<code>--haps</code>	Filename of haps file.
<code>--sample</code>	Filename of sample file.
<code>--map</code>	Filename of genetic recombination map.
<code>-o,--output</code>	Filename of output files without file extension.

Note 4: Optional arguments

There are other optional arguments for Relate, please consult the documentation for these (https://myersgroup.github.io/relate/getting_started.html#GettingStarted).

2 Running Relate on the Simons Genome Diversity Project dataset

Let's start with running Relate. Please copy the `worksheet_data` directory to your home directory

```
#assuming you are in your home directory
cp -r /home/speidel/worksheet_data .
cd worksheet_data
```

In this directory, you will find three subdirectories, named `data` (containing all input data), `precomp_results` (containing precomputed Relate output files), and `Rscripts` (containing R scripts for plotting results).

2.1 Running Relate

Note 5

Running Relate on all 278 samples and four cores took 52 minutes, so here we will run Relate only on the African samples of this dataset.

We run Relate on chromosome 15 for the 44 African samples of the Simons Genome Diversity Project.

```
#assuming you are under worksheet_data
mkdir results
cd results

lab="Africa"
chr=15
#This takes about 3 minutes
${PATH_TO_RELATE}/bin/Relate \
  --mode All \
  --haps "../data/data_subgroups/SGDP_input_${lab}_chr${chr}.haps.gz" \
  --sample "../data/data_subgroups/SGDP_input_${lab}_chr${chr}.sample.gz" \
  --map "../data/genetic_map_chr${chr}_combined_b37.txt" \
  --dist "../data/data_subgroups/SGDP_input_${lab}_chr${chr}.dist.gz" \
  -m 1.25e-8 \
  -N 30000 \
  -o "SGDP_${lab}_chr${chr}"

gzip SGDP_${lab}_chr15.*
```

The haps/sample files store the genetic variation data (similar to a vcf file) and is outputted by the `PrepareInputFiles.sh` script (see above). The genetic recombination map stores recombination rates along the genome. In addition, we specify a `.dist` file to adjust the distances (in units of BP) between SNPs removing regions of lower quality specified in the mask – this file is also outputted by the `PrepareInputFiles.sh` script and is necessary to adjust mutation rates (we can only observe mutations in regions that pass the filters of the mask).

Exercise 1

Let's plot a few trees. We use the TreeView.sh script provided with Relate: For this, we need an additional file storing assignment of individuals to populations. This file has four columns, named "sample", "population", "group", and "sex". For us, only the second column is of interest. The order in which individuals are listed has to be consistent with the order in the samples file.

Example:

```
ID POP GROUP SEX
S_Mozabite-1 Mozabite Africa NA
S_Mozabite-2 Mozabite Africa NA
S_Saharawi-1 Saharawi Africa NA
S_Saharawi-2 Saharawi Africa NA
```

We can now run the TreeView.sh script as follows (this will use R and requires ggplot2 and cowplot - these will be installed if missing):

```
${PATH_TO_RELATE}/scripts/TreeView/TreeView.sh \
  --haps "../data/data_subgroups/SGDP_input_Africa_chr${chr}.haps.gz" \
  --sample "../data/data_subgroups/SGDP_input_Africa_chr${chr}.sample.gz" \
  --anc "SGDP_Africa_chr15.anc.gz" \
  --mut "SGDP_Africa_chr15.mut.gz" \
  --poplabels "../data/poplabels/SGDP_Africa_lang.poplabels" \
  --bp_of_interest 20000000 \
  --years_per_gen 28 \
  -o "SGDP_Africa_chr${chr}_BP20000000"
```

This will produce a pdf named SGDP_Africa_chr15_BP20000000.pdf. We can also plot trees for the precomputed trees of West Eurasian SGDP individuals:

```
#This will take approx. 2 minutes
${PATH_TO_RELATE}/scripts/TreeView/TreeView.sh \
  --haps "../data/data_subgroups/SGDP_input_WestEurasia_chr15.haps.gz" \
  --sample "../data/data_subgroups/SGDP_input_WestEurasia_chr15.sample.gz" \
  --anc "../precomp_results/WestEurasia/SGDP_WestEurasia_ne_chr15.anc.gz" \
  --mut "../precomp_results/WestEurasia/SGDP_WestEurasia_ne_chr15.mut.gz" \
  --poplabels "../data/poplabels/SGDP_WestEurasia.poplabels" \
  --bp_of_interest 48426484 \
  --years_per_gen 28 \
  -o "SGDP_WestEurasia_chr${chr}_BP48426484"
```

Can you see anything unusual about this tree. You can also rerun the above replacing TreeView.sh by TreeViewMutation.sh; this will highlight branches carrying the mutation at chromosome 15, BP 48426484.

Feel free to plot a few other trees!

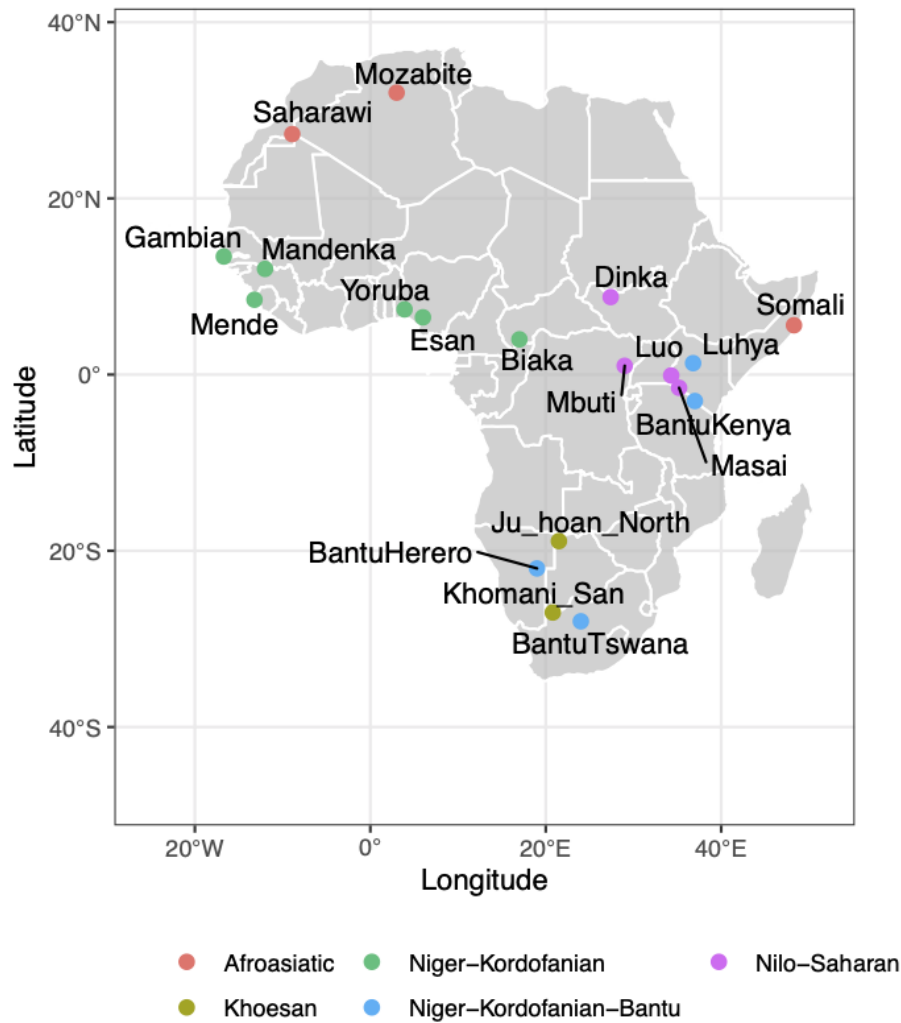


Figure 2: Map showing SGDP groups in Africa, coloured by language family.

3 Effective population sizes and split times

Note 6

We will use the genealogies for Africans obtained in Section 2.1 as an example here. You can replicate the analysis for other subgroups, or for all samples; each step will take longer for larger sample sizes (for all samples, computation time in Section 3.1 is around 15 min and for Section 3.2 around 70 min).

3.1 Estimating population sizes given a genealogy

We infer effective population sizes and split times from estimated genealogies of Africans obtained in Section 2.1. Extracting effective population sizes from genealogies is done as follows

```
#assuming you are in worksheet_data
mkdir popsizes
cd popsizes

lab="Africa"
chr=15
${PATH_TO_RELATE}/bin/RelateCoalescentRate \
  --mode EstimatePopulationSize \
  -i "../results/SGDP_${lab}_chr${chr}" \
  -o "SGDP_${lab}_chr${chr}" \
  --bins 3,7,0.25 \
  --poplabels "../data/poplabels/SGDP_${lab}.poplabels"
```

This will produce a file named `SGDP_Africa_chr15.coal` and `SGDP_Africa_chr15.bin`. The `SGDP_Africa_chr15.coal` stores cross-coalescence rates for all pairs of group labels in `SGDP_Africa.poplabels` of which there are 19. This might be too many to visualise easily, so we will aggregate some of these groups based on language family. The `.bin` file can be used in combination with different `poplabels` files to extract coalescence rates, e.g.

```
${PATH_TO_RELATE}/bin/RelateCoalescentRate \
  --mode FinalizePopulationSize \
  -o "SGDP_Africa_chr${chr}" \
  --poplabels "../data/poplabels/SGDP_Africa_lang.poplabels"
```

Exercise 2

- (i) Let's plot the estimated effective population sizes stored in `SGDP_Africa_chr15.coal` in R. You can use `relater`, which is an R package that can parse and manipulate Relate output files (preinstalled, available at <https://github.com/leospeidel/relater>):

```
library(relater)
coal <- read.coal("SGDP_Africa_chr15.coal")

#diploid effective population size is the 0.5* inverse coalescence rate
coal$popsize <- 0.5/coal$haploid.coalescence.rate
#multiply epochs times by 28 to scale to years (assuming 28 years per generation)
coal$epoch.start <- 28 * coal$epoch.start
head(coal)
```

Can you estimate the split time of Khoesan and other groups? How are the population sizes of groups speaking Afroasiatic languages different from other groups?

- (ii) Optional: You can also plot the effective population sizes for all samples of the Simons Genome Diversity Project; for this, the relevant precomputed file is located at `../precomp_results/SGDP_all/SGDP_v1_annot_ne.pairwise.bin` and `*.coal`. You may want to aggregate some groups again, e.g., using `../data/poplabels/SGDP_region.poplabels`.

3.2 Joint fitting of population size and branch lengths

So far, we used trees that assumed a pre-specified constant effective population size through time. Next, we use the `EstimatePopulationSize.sh` script to jointly fit effective population sizes and branch lengths.

```
chr=15
lab=Africa

# This will take around 4min
${PATH_TO_RELATE}/scripts/EstimatePopulationSize/EstimatePopulationSize.sh \
-i "../results/SGDP_${lab}_chr${chr}" \
-o "SGDP_${lab}_ne_chr${chr}" \
--poplabels "../data/poplabels/SGDP_${lab}_lang.poplabels" \
-m 1.25e-8 \
--years_per_gen 28 \
--bins 3,7,0.25 \
--num_iter 2 \
--threads 4
```

This will output `SGDP_Africa_ne_chr15.anc.gz`, `SGDP_Africa_ne_chr15.mut.gz`, and `SGDP_Africa_ne_chr15.pairwise.coal`. (The file `SGDP_Africa_ne_chr15.coal` stores coalescence rates grouping all samples into one group and can be ignored here.)

Note 7

To run this on more than one chromosome, the `--chr` argument can be useful (see documentation <https://myersgroup.github.io/relate/modules.html#PopulationSizeScript>.)

Exercise 3

- (i) Similarly to Exercise 2, let's plot the effective population sizes for `SGDP_Africa_ne_chr15.pairwise.*`. How is this plot different to the one generated in Exercise 2 (which assumed a constant population size through time)?
(If you like, you can use the script under `Rscripts/`.)
- (ii) Now, let's get coalescence rates between all pairs of African individuals from `SGDP_Africa_ne_chr15.pairwise` using the poplabels file `../data/poplabels/SGDP_Africa_ind.poplabels`. We can plot coalescence rates in a heatmap for different epochs to visualise how structure has evolved through time.
(If you like, you can use the script under `Rscripts/`.)

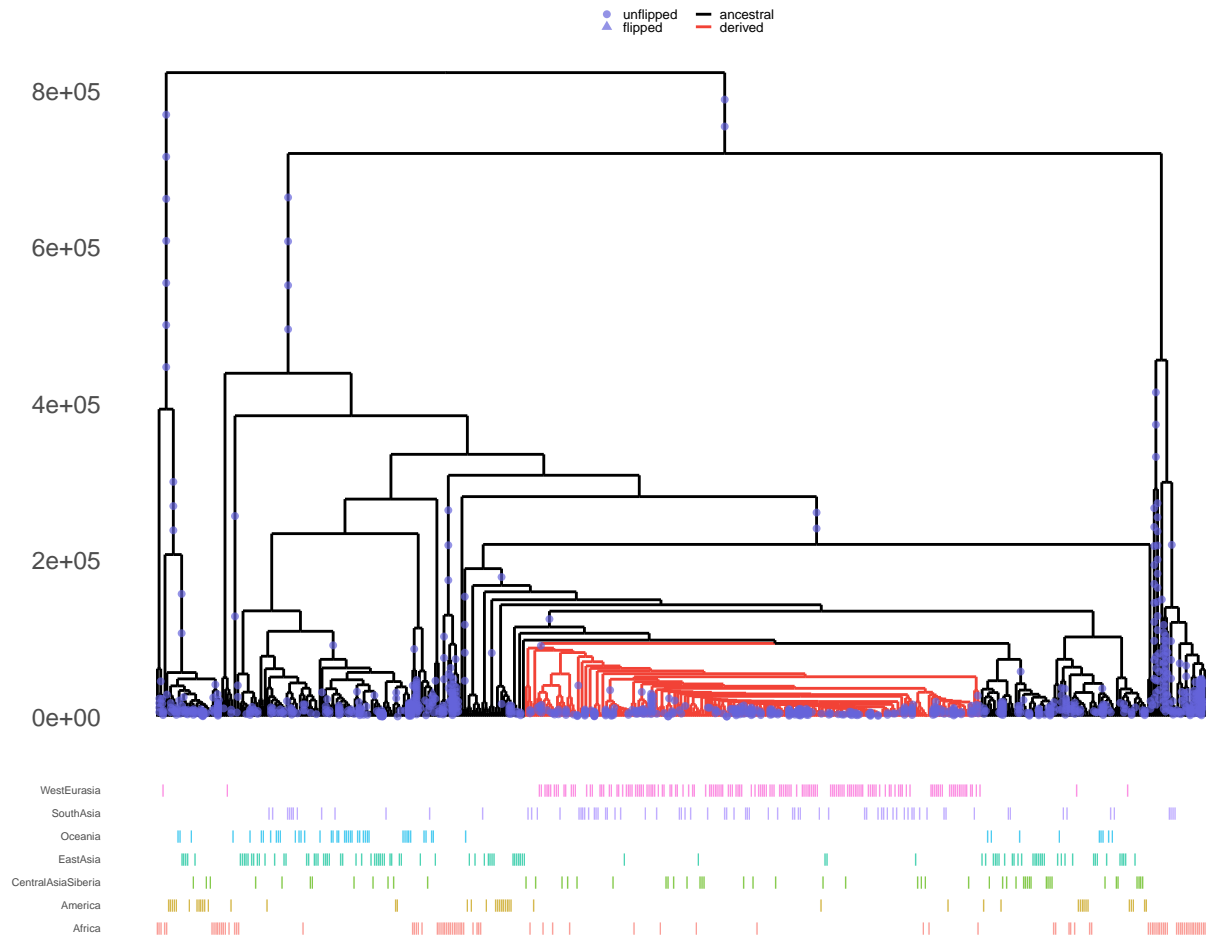


Figure 3: Marginal genealogical tree for rs1426654, a non-synonymous substitution associated with light skin pigmentation. The SNP is located in the SLC24A5 gene on chromosome 15 (BP=48,426,484). Derived allele carriers are shown in red.

4 Detecting evidence for positive selection

Positive natural selection on a derived allele is expected to lead to this allele spreading rapidly in a population, reflected in a burst of coalescence events. A well-known example is positive selection of a variant associated with lighter skin pigmentation in Europe and South Asia (Fig. 3). We have implemented a simple statistic that captures such events and measures the extent to which a mutation has out-competed other lineages.

To analyse selection, we will first apply functions provided with Relate to extract the relevant statistics, and then use the R package `relater` to analyse the output.

Note 8

We subset the sample into subgroups (e.g., WestEurasia) because the selection p-values assume a panmictic population and may get confounded by genetic structure.

In principle, the selection p-values are robust to misspecification of demographic histories, so we can use genealogies assuming a constant population size here, although we will use genealogies with jointly fitted population sizes and branch lengths below.

```
#assuming you are in worksheet_data
mkdir selection
cd selection

chr=15
lab="WestEurasia"
output="SGDP_${lab}_ne_chr${chr}"

# Copy anc/mut files for WestEurasia
cp ../precomp_results/WestEurasia/${output}* .

# Using these genealogies, calculate frequencies through time
# This will output a *.freq and *.lin file
${PATH_TO_RELATE}/bin/RelateSelection \
    --mode Frequency \
    -i ${output} \
    -o ${output}

# Next, calculate selection p-values.
# This will take the output of the previous step as input and output a *.sele file.
${PATH_TO_RELATE}/bin/RelateSelection \
    --mode Selection \
    -i ${output} \
    -o ${output}

# Also calculate a *.qual file storing statistics about the quality of trees
# (e.g., number of mutations mapping to the tree)
${PATH_TO_RELATE}/bin/RelateSelection \
    --mode Quality \
    -i ${output} \
    -o ${output}
```

Once we have generated these files, we use relater to analyse them in R.

```
library(relater)
output <- "./SGDP_WestEurasia_ne_chr15"

# parse files
mut <- read.mut(paste0(output, ".mut.gz"))
freq <- read.freq(paste0(output, ".freq"))
sele <- read.sele(paste0(output, ".sele"))
qual <- read.qual(paste0(output, ".qual"))

# combine these into a single data frame called allele_ages
allele_ages <- get.allele_ages(mut = mut, freq = freq, sele = sele)
# use the *qual file to filter out SNPs mapping to "bad" trees
allele_ages <- filter.allele_ages(allele_ages, qual)
allele_ages <- subset(allele_ages, !is.na(pvalue))

head(allele_ages[order(allele_ages$pvalue),])
```

Exercise 4

Let's analyse the allele_ages data frame. What is the top SNP with the lowest selection p-value?