

Find SNPs distinguish two populations

Jinliang Yang

Nov. 2nd, 2014

INTRODUCTION

Develop a piece of software to identify a set of SNPs that distinguish two populations.

INSTALL and USAGE

```
git clone git@github.com:RILAB/N2fixation.git
cp N2fixation/lib/fpsnp_v0.1.py ~/bin/fpsnp
export PATH=$PATH:~/bin/
fpsnp -h
```

input data

The input SNP data should be in BED-like format. The first three columns are **snpid**, **chr** and **pos**. The following columns in the first row should be any number of individuals with SNP genotype information. In the second row, the first three columns should be "N" and the followings are two levels of intergers (such as 1 and 2) to specify the populations they belonging to. See the example below.

```
dsf <- read.table("../data/simunsnp.txt", nrow = 5,
  header = TRUE)
head(dsf[, 1:7])
```

##	snpid	chr	pos	SAM105	SAM367	SAM70	SAM106
## 1	N	N	N	1	1	1	1
## 2	1_86	1	86	N	N	N	N
## 3	1_492	1	492	T	T	T	T
## 4	1_509	1	509	C	C	C	C
## 5	1_1825	1	1825	G	G	G	G

Run the following command to get the population frequencies:

```
fpsnp -p . -i data/simunsnp.txt -o tests/output.txt
```

test the correctness of the frequency calculations

The following R codes do the same population specific allele frequency calculation. But it is less efficient.

```
### simulate 1000 SNP data
source("../profiling/1.A.1.simulate_SNP_data.R")
### results computed from R were written to
### tests/routput.txt
source("../profiling/1.A.2.R_fpSNP_test.R")
```

After comparisons, only one SNP got different minor allele frequency values in population one and population two. While, this SNP has a overall MAF of 0.5, the minor allele was determined randomly by Python and R. So it is not a problem for our calculation.

```
### compare the results from python and R
setwd("../")
source("profiling/1.A.3.p_r_comp.R")
```

```
# idx <- which(pout$maf2 != rout$maf2)
# rout[idx,] pout[idx,]
```

FINDSNPs to distinguish two populations

With the computed allele frequencies in the two populations, we will be able to distinguish the two populations by select the most different SNPs in terms of their allele frequencies. To do this, a statistical test was conducted. The null hypothesis for this test is that for a given set of independent SNPs in two identical populations the allele frequencies are the same. If the possibility that the observed allele frequency differences for a given set of SNPs is smaller than a threshold (say $pvalue < 0.05$), we will have evidence to reject the null hypothesis and claim the two populations are different.

We used the paired t-test to conduct the analysis. To control the independence of the SNPs, we used a selected bin window (for example 1-Mb) to remove SNPs that might be in LD.

```
### simulate 1000 SNP data
source("../lib/getfpsnp.R")
test <- findSNP(frqfile = "../tests/output.txt",
  binsize = 100, pcutoff = 0.001, missingrate = 0.5)

## ###>>> With the allele frequencies of [ 20 ] selected SNPs, only [ 0.000732981206556325 ]
## possibility that the two populations are the same!

head(test)
```

##	snpid	chr	pos	major	minor	maf0
##	429	1_1040057	1 1040057	C	G	0.3529
##	428	1_1040037	1 1040037	A	G	0.3018

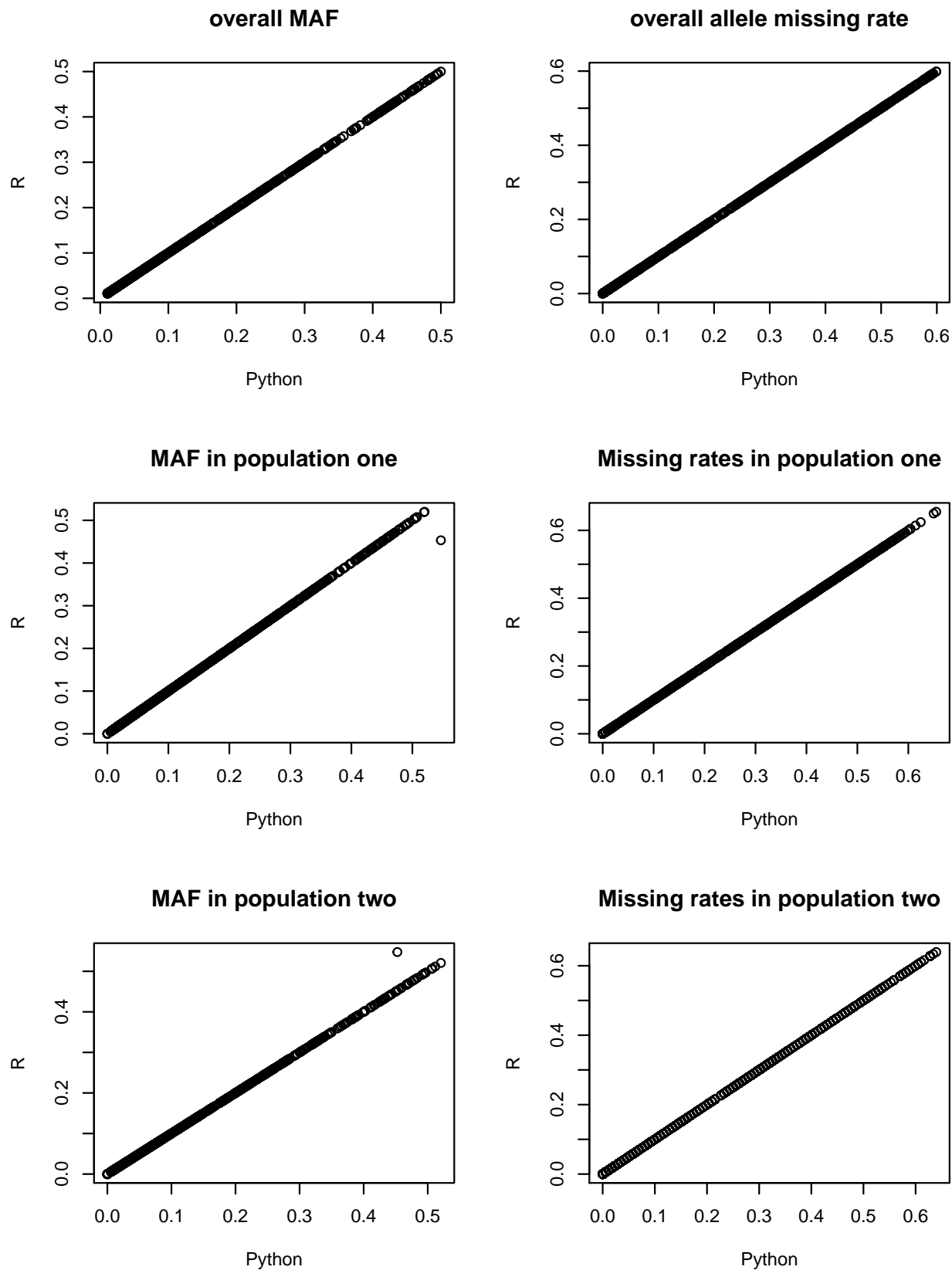


Figure 1: Allele frequency comparisons

```

## 348 1_1003542    1 1003542      G      A 0.3412
## 734 1_1860196    1 1860196      A      G 0.1980
## 360 1_1037879    1 1037879      C      G 0.2951
## 142 1_273881     1 273881       C      T 0.4670
##      missing0    maf1 missing1    maf2 missing2
## 429  0.3550 0.4426    0.3807 0.2586    0.3256
## 428  0.3984 0.3805    0.4264 0.2202    0.3663
## 348  0.4282 0.2793    0.4365 0.4100    0.4186
## 734  0.4661 0.2571    0.4670 0.1304    0.4651
## 360  0.3387 0.2419    0.3706 0.3500    0.3023
## 142  0.4661 0.5200    0.4924 0.4124    0.4360
##      diff      bin
## 429 0.1840 1_10401
## 428 0.1603 1_10400
## 348 0.1307 1_10035
## 734 0.1267 1_18602
## 360 0.1081 1_10379
## 142 0.1076 1_2739

```