

Population stratification

GBIO0009

Kridsakorn Chaichoompu

University of Liege

Population stratification

Population stratification is the presence of a **systematic difference in allele frequencies between subpopulations** in a population possibly due to **different ancestry**, especially in the context of association studies. Population stratification is also referred as population structure, in this context.

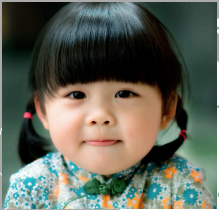




Diversity

- Human
- Plants
- Animals
- Bacteria
- etc

Human Diversity



How to group people?



Countries

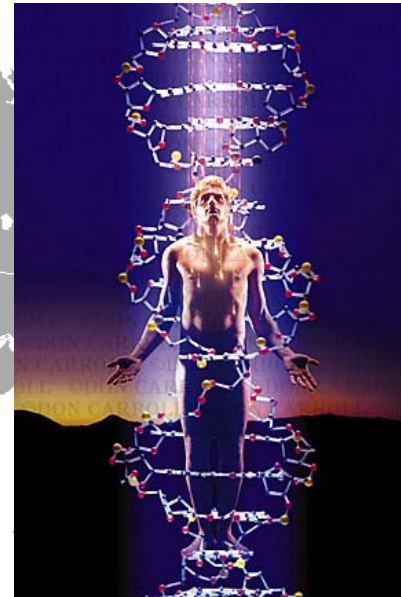


Languages

Physical appearances: Hair colors, Eye colors, Skin colors



DNA: the blueprint of our lives



PROPER DRUGS AND TREATMENT



Principal Component Analysis (PCA)

Principal component analysis (PCA) is a statistical procedure that uses an orthogonal transformation to convert a set of observations of possibly correlated variables into a set of values of linearly uncorrelated variables called **principal components (PCs)**.



PCA in R

- `prcomp(x, retx = TRUE, center = TRUE, scale. = FALSE, tol = NULL, ...)`
- `princomp(formula, data = NULL, subset, na.action, ...)`
- `eigen(x, symmetric, only.values = FALSE, EISPACK = FALSE)`
- `svd(x, nu = min(n, p), nv = min(n, p), LINPACK = FALSE)`

`library(rARPACK)`

- `svds(A, k, nu = k, nv = k, opts = list(), ...)`
- `eigs(A, k, which = "LM", sigma = NULL, opts = list(), ...)`

PCA for GWAS

Principal components analysis corrects for stratification in genome-wide association studies

Alkes L Price^{1,2}, Nick J Patterson², Robert M Plenge^{2,3}, Michael E Weinblatt³, Nancy A Shadick³ & David Reich^{1,2}

Population stratification—allele frequency differences between cases and controls due to systematic ancestry differences—can cause spurious associations in disease studies. We describe a method that enables explicit detection and correction of population stratification on a genome-wide scale. Our method uses principal components analysis to explicitly model ancestry differences between cases and controls. The resulting correction is specific to a candidate marker’s variation in frequency across ancestral populations, minimizing spurious associations while maximizing power to detect true associations. Our simple, efficient approach can easily be applied to disease studies with hundreds of thousands of markers.

PCA for GWAS (Price, 2006)

The above procedure is motivated by the decomposition $X = USV^T$, where U is an $M \times N$ matrix whose k th column contains coordinates of each SNP along the k th principal component, S is a diagonal matrix of singular values and V is an $N \times N$ matrix whose k th column contains ancestries a_{jk} of each individual j along the k th principal component. It follows that $X^T X = VS^2 V^T$; thus, the columns of V are the eigenvectors of the matrix $X^T X$. The matrix $X^T X$ is equivalent up to a constant to the covariance matrix Ψ , and the matrix S^2 of squared singular values is equivalent up to a constant to the diagonal matrix of eigenvalues of Ψ .

snpStats – Bioconductor Package

- <http://www.bioconductor.org/packages/release/bioc/html/snpStats.html>

Usually, principal components analysis is carried out by calculating the eigenvalues and eigenvectors of the correlation matrix. With N cases and P variables, if we write X for the $N \times P$ matrix which has been standardised so that columns have zero mean and unit standard deviation, we find the eigenvalues and eigenvectors of the $P \times P$ matrix $X^T.X$ (which is N or $(N - 1)$ times the correlation matrix depending on which denominator was used when calculating standard deviations). The first eigenvector gives the loadings of each variable in the first principal component, the second eigenvector gives the loadings in the second component, and so on. Writing the first C component loadings as columns of the $P \times C$ matrix B , the $N \times C$ matrix of subjects' principal component scores, S , is obtained by applying the factor loadings to the original data matrix, *i.e.* $S = X.B$. The sum of squares and products matrix, $S^T.S = D$, is diagonal with elements equal to the first C eigenvalues of the $X^T.X$ matrix, so that the variances of the principal components can be obtained by dividing the eigenvalues by N or $(N - 1)$.

snpStats - PCA

This standard method is rarely feasible for genome-wide data since P is very large indeed and calculating the eigenvectors of $X^T.X$ becomes impossibly onerous. However, the calculations can also be carried out by calculating the eigenvalues and eigenvectors of the $N \times N$ matrix $X.X^T$. The (non-zero) eigenvalues of this matrix are the same as those of $X^T.X$, and its eigenvectors are proportional to the principal component scores defined above; writing the first C eigenvectors of $X.X^T$ as the columns of the $N \times C$ matrix, U , then $U = S.D^{-1/2}$. Since for many purposes we are not too concerned about the scaling of the principal components, it will often be acceptable to use the eigenvectors, U , in place of the more conventionally scaled principal components. However some attention should be paid to the corresponding eigenvalues since, as noted above, these are proportional to the variances of the conventional principle components. The factor loadings may be calculated by $B = X^T.U.D^{-1/2}$.

The next step in the calculation is to obtain the SNP loadings in the components. This requires calculation of $B = X^T.S.D^{-1/2}$. Here we calculate the transpose of this matrix, $B^T = D^{-1/2}S^T.X$, using the special function `snp.pre.multiply` which pre-multiplies a `SnpMatrix` object by a matrix after first standardizing it to zero mean and unit standard deviation.

PCA for SNPs

- X is the $M \times N$ matrix, where M is a number of individuals and N is a number of SNPs.

$$XX^T = UDV^T$$

Note: In this case, U and V are equal because XX^T is a square matrix

U is the matrix of eigenvectors or PC scores.

$$B^T = D^{-1/2}U^TX$$

B is the factor loadings

$$\text{PCs} = X.B$$

Normalization

- Zero means

If X is a vector

$$M = X - \text{mean}(X)$$

- Unit variance

$$Y = M / \text{sd}(X)$$

- In R, it is more efficient to use `apply()` with `mean()` and `sd()`

Quality Control

- Missing data
- Linkage Disequilibrium (LD) pruning
- Hardy-Weinberg Equilibrium (HWE)

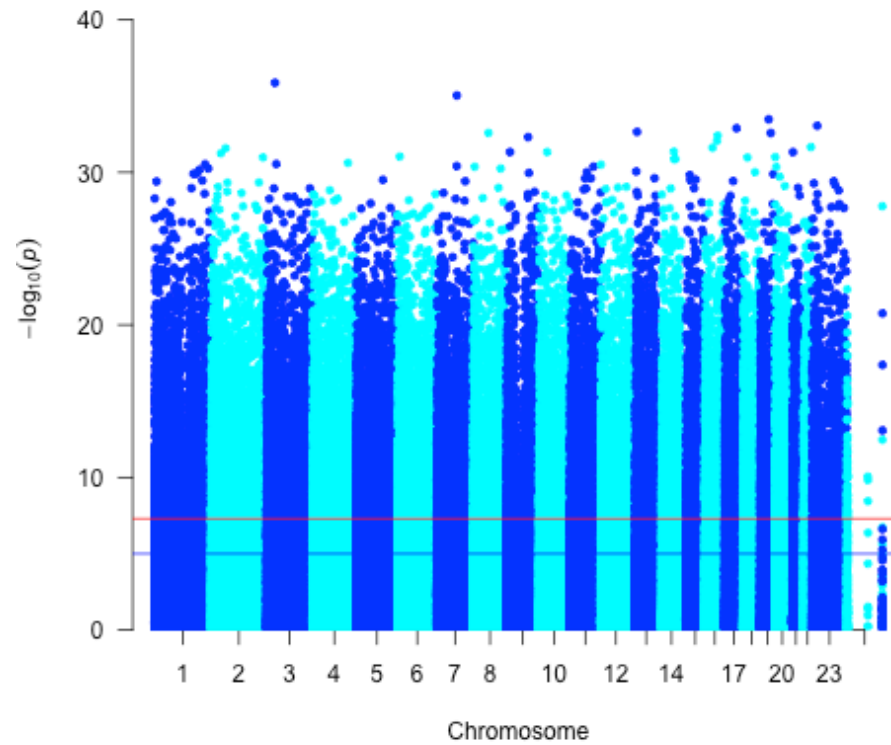
Suggestion: use PLINK

<http://pngu.mgh.harvard.edu/~purcell/plink/>

Exercise - PCA

- Calculate PCs for the example data - simSNP_rep1, more information:
 - Non-redundant SNPs, no LD
 - No missing data
 - Follow HWE
- Plot the first two eigenvectors
- Plot the first two PCs

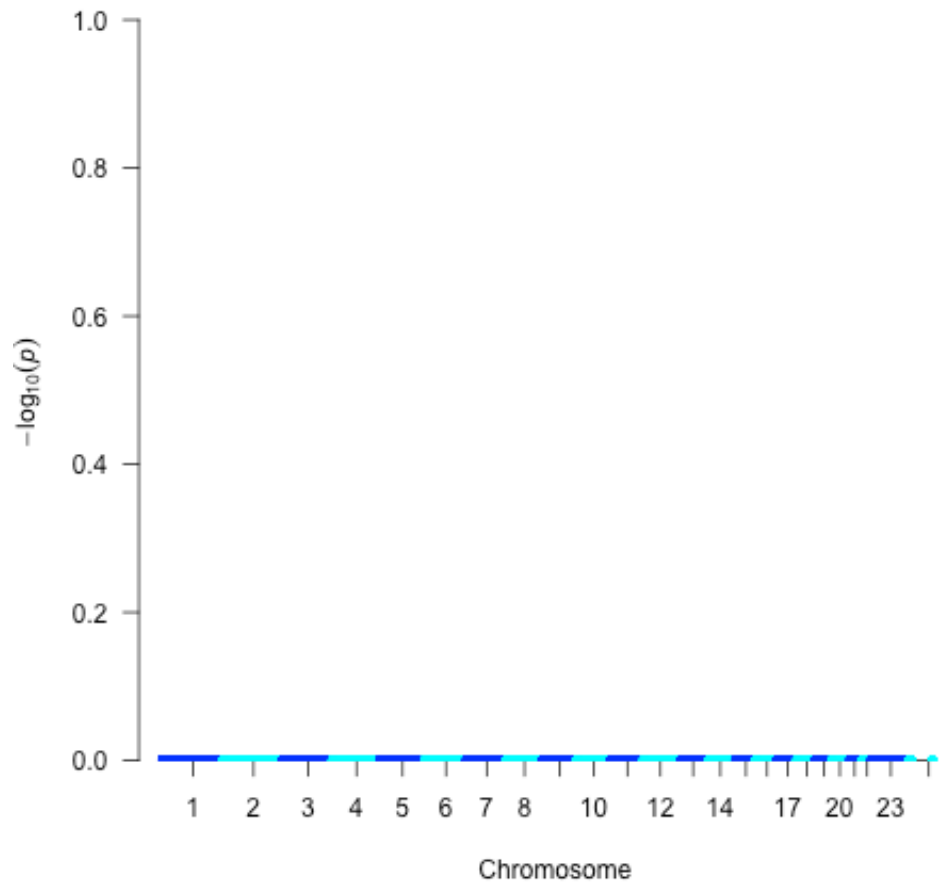
HW1



PC Adjustment in PLINK

- For quantitative traits, use
`plink --bfile mydata --linear`
- For disease traits, specify logistic regression with
`plink --bfile mydaya --logistic`
- Adjust with covariates, then the command
`plink --bfile mydata --linear --genotypic --covar mycov.txt`

Adjusted Manhattan plot of HW1



Linear Regression in R

Linear models

`lm(formula, data, subset, ...)`

Example in help page:

```
ctl <- c(4.17, 5.58, 5.18, 6.11, 4.50, 4.61, 5.17, 4.53, 5.33, 5.14)
trt <- c(4.81, 4.17, 4.41, 3.59, 5.87, 3.83, 6.03, 4.89, 4.32, 4.69)
group <- gl(2, 10, 20, labels = c("Ctl", "Trt"))
weight <- c(ctl, trt)
lm.D9 <- lm(weight ~ group)
plot(lm.D9)
```

<https://stat.ethz.ch/R-manual/R-devel/library/stats/html/lm.html>

Generalized Linear Models - GLM

`glm(formula, family = gaussian, data, weights, ...)`

Example from help page:

```
counts <- c(18,17,15,20,10,20,25,13,12)
outcome <- gl(3,1,9)
treatment <- gl(3,3)
print(d.AD <- data.frame(treatment, outcome, counts))
glm.D93 <- glm(counts ~ outcome + treatment, family =
poisson())
```

<http://stat.ethz.ch/R-manual/R-patched/library/stats/html/glm.html>

Models for GLM

```
glm(formula, family=familytype(link=linkfunction), data=)
```

Family	Default Link Function
binomial	(link = "logit")
gaussian	(link = "identity")
Gamma	(link = "inverse")
inverse.gaussian	(link = "1/mu^2")
poisson	(link = "log")
quasi	(link = "identity", variance = "constant")
quasibinomial	(link = "logit")
quasipoisson	(link = "log")

<http://www.statmethods.net/advstats/glm.html>

Exercise – Linear regression

- Do linear regression with the example data using

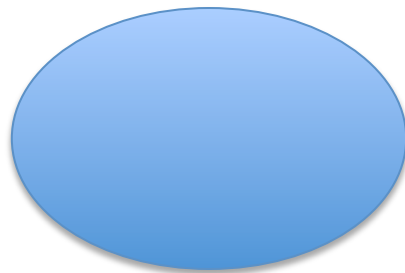
$$\text{SNPs} \sim \text{PC1} + \text{PC2} + \text{PC3} + \text{PC4} + \text{PC5}$$

- Calculate PCs from the residuals
- Check PC plot
- Try with `glm()` with logistic model

Fixation index (F_{ST})

- F_{ST} can be used to describe a distance among population.
- F_{ST} can be biased due to the allele frequencies and the number of independent SNPs.

Pop1 = 2,000 individuals



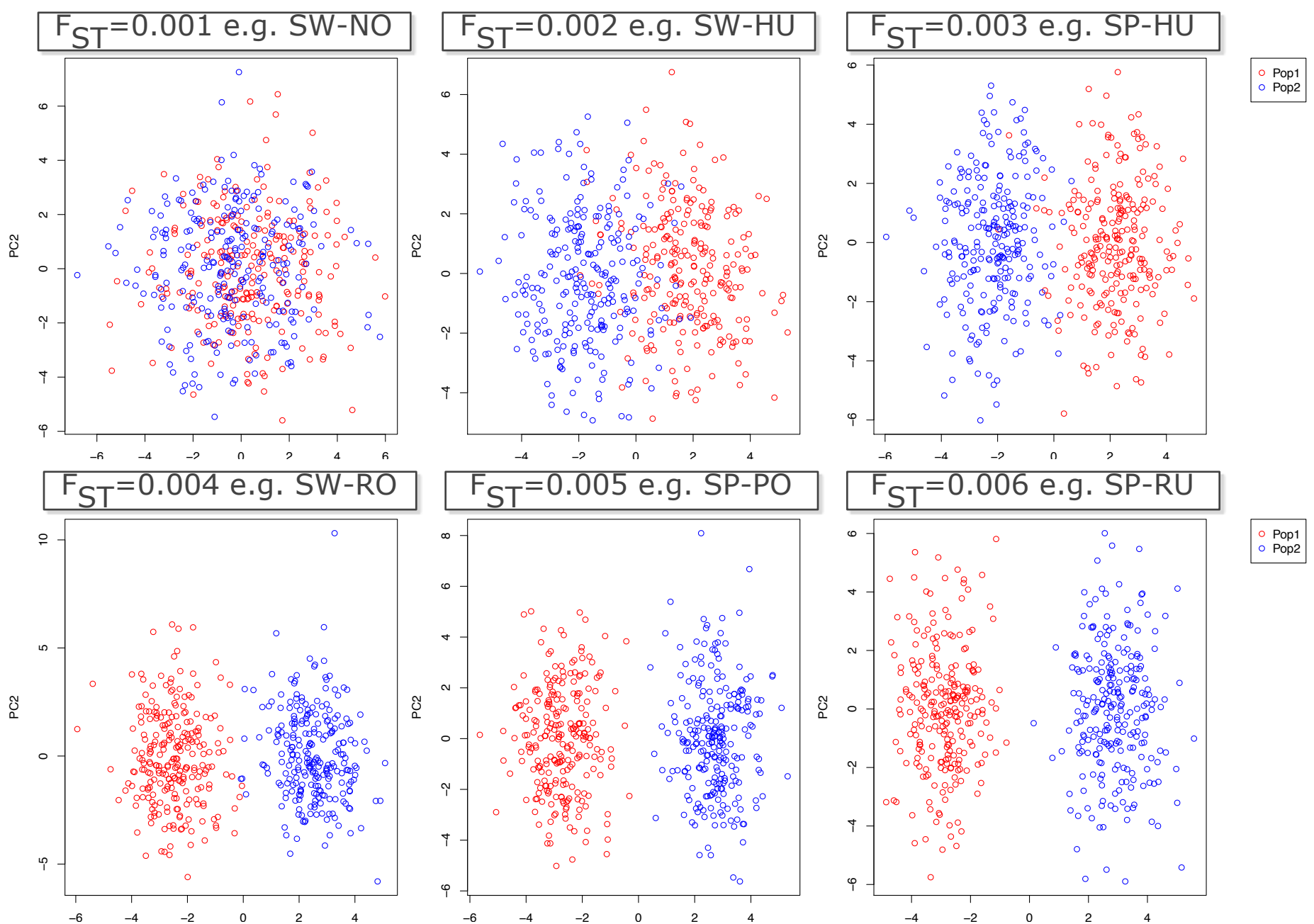
Pop2 = 500 individuals



F_{ST} among European populations

	<i>Sp</i>	<i>Fr</i>	<i>Be</i>	<i>UK</i>	<i>Sw</i>	<i>No</i>	<i>Ge</i>	<i>Ro</i>	<i>Cz</i>	<i>SI</i>	<i>Hu</i>	<i>Po</i>	<i>Ru</i>	<i>CEU</i>	<i>CHB</i>	<i>JPT</i>
<i>Fr</i>	0.0008															
<i>Be</i>	0.0015	0.0002														
<i>UK</i>	0.0024	0.0006	0.0005													
<i>Sw</i>	0.0047	0.0023	0.0018	0.0013												
<i>No</i>	0.0047	0.0024	0.0019	0.0014	0.0010											
<i>Ge</i>	0.0025	0.0008	0.0005	0.0006	0.0011	0.0016										
<i>Ro</i>	0.0023	0.0017	0.0018	0.0028	0.0041	0.0044	0.0016									
<i>Cz</i>	0.0033	0.0016	0.0013	0.0014	0.0016	0.0024	0.0003	0.0016								
<i>SI</i>	0.0034	0.0017	0.0015	0.0017	0.0019	0.0026	0.0005	0.0014	0.0001							
<i>Hu</i>	0.0030	0.0015	0.0013	0.0016	0.0020	0.0026	0.0004	0.0011	0.0001	0.0001						
<i>Po</i>	0.0053	0.0032	0.0028	0.0027	0.0023	0.0034	0.0012	0.0028	0.0004	0.0004	0.0006					
<i>Ru</i>	0.0059	0.0037	0.0034	0.0032	0.0025	0.0036	0.0016	0.0030	0.0008	0.0007	0.0009	0.0003				
<i>CEU</i>	0.0026	0.0008	0.0005	0.0002	0.0011	0.0012	0.0006	0.0028	0.0014	0.0016	0.0016	0.0026	0.0031			
<i>CHB</i>	0.1096	0.1094	0.1093	0.1096	0.1073	0.1081	0.1085	0.1047	0.1080	0.1069	0.1058	0.1086	0.1036	0.1095		
<i>JPT</i>	0.1118	0.1116	0.1114	0.1117	0.1095	0.1103	0.1107	0.1068	0.1102	0.1091	0.1079	0.1108	0.1057	0.1117	0.0069	
<i>YRI</i>	0.1460	0.1493	0.1496	0.1513	0.1524	0.1531	0.1502	0.1463	0.1503	0.1498	0.1490	0.1520	0.1504	0.1510	0.1901	0.1918

Simon et al. 2008



To understand F_{ST} , here are simulated data using Balding method and the examples of EU populations as reported in (Simon et al. 2008)

F_{ST} – R Packages

Package ‘PopGenome’

May 4, 2015

Type Package

Title An Efficient Swiss Army Knife for Population Genomic Analyses

Version 2.1.6

Date 2015-05-1

Package ‘hierfstat’

December 4, 2015

Version 0.04-22

Date 2015-11-24

Title Estimation and Tests of Hierarchical F-Statistics

Package ‘StAMPP’

July 6, 2015

Type Package

Title Statistical Analysis of Mixed Ploidy Populations

Depends R (>= 2.14.0), pegas

Imports parallel, doParallel, foreach, adegenet, methods, utils

Version 1.4

Date 2015-06-30

Estimating F_{ST}

Method

Estimating and interpreting F_{ST} : The impact of rare variants

Gaurav Bhatia,^{1,2,6,7} Nick Patterson,^{2,6,7} Sriram Sankararaman,^{2,3} and Alkes L. Price^{2,4,5,7}

¹Harvard–Massachusetts Institute of Technology (MIT), Division of Health, Science, and Technology, Cambridge, Massachusetts 02139, USA; ²Broad Institute of Harvard and MIT, Cambridge, Massachusetts 02142, USA; ³Department of Genetics, Harvard Medical School, Boston, Massachusetts 02115, USA; ⁴Department of Epidemiology, Harvard School of Public Health, Boston, Massachusetts 02115, USA; ⁵Department of Biostatistics, Harvard School of Public Health, Boston, Massachusetts 02115, USA

In a pair of seminal papers, Sewall Wright and Gustave Malécot introduced F_{ST} as a measure of structure in natural populations. In the decades that followed, a number of papers provided differing definitions, estimation methods, and interpretations beyond Wright's. While this diversity in methods has enabled many studies in genetics, it has also introduced confusion regarding how to estimate F_{ST} from available data. Considering this confusion, wide variation in published estimates of F_{ST} for pairs of HapMap populations is a cause for concern. These estimates changed—in some cases more than twofold—when comparing estimates from genotyping arrays to those from sequence data. Indeed, changes in F_{ST} from sequencing data might be expected due to population genetic factors affecting rare variants. While rare variants do influence the result, we show that this is largely through differences in estimation methods. Correcting for this yields estimates of F_{ST} that are much more concordant between sequence and genotype data. These differences relate to three specific issues: (1) estimating F_{ST} for a single SNP, (2) combining estimates of F_{ST} across multiple SNPs, and (3) selecting the set of SNPs used in the computation. Changes in each of these aspects of estimation may result in F_{ST} estimates that are highly divergent from one another. Here, we clarify these issues and propose solutions.

Hudson's F_{ST}

Definition

Hudson et al. (1992) defined F_{ST} in terms of heterozygosity. The fundamental difference between these estimators is that for Hudson, the total variance is based upon the ancestral population and not the current sample.

Estimator

Hudson's estimator for F_{ST} is given by

$$\hat{F}_{ST}^{Hudson} = 1 - \frac{H_w}{H_b}, \quad (9)$$

where H_w is the mean number of differences within populations, and H_b is the mean number of differences between populations. While Hudson did not give explicit equations for H_w and H_b , we cast his description into an explicit estimator (see Supplemental Material for a derivation). The estimator that we analyze is

$$\hat{F}_{ST}^{Hudson} = \frac{(\tilde{p}_1 - \tilde{p}_2)^2 - \frac{\tilde{p}_1(1 - \tilde{p}_1)}{n_1 - 1} - \frac{\tilde{p}_2(1 - \tilde{p}_2)}{n_2 - 1}}{\tilde{p}_1(1 - \tilde{p}_2) + \tilde{p}_2(1 - \tilde{p}_1)}, \quad (10)$$

where n_i is the sample size and \tilde{p}_i is the sample allele frequency in population i for $i \in \{1, 2\}$. Analyzing this estimator using the definition of Weir and Hill (2002), we show (see Supplemental Material) that F_{ST} estimated using Hudson's estimator will tend toward Equation 3 (see Results), which is exactly the average of population-specific F_{ST} values that we seek to estimate. This emerges naturally, as the proposed estimator is the simple average of the population-specific estimators given in Weir and Hill (2002). This estimator has the desirable properties that it is (1) independent of sample composition, and (2) does not overestimate F_{ST} (it has a maximum value of 1). We recommend its use to produce estimates of F_{ST} for two populations.

Exercise – F_{ST} estimation

- Implement Hudson's method
- Estimate the average pairwise F_{ST} values for Pop1-6.