

# Instructions of how to use PGEToolbox

## *Installation and Getting Started*

PGEToolbox can be downloaded from <http://bioinformatics.org/pgettoolbox/>

Step by step instruction:

1. Download 'PGEToolbox.zip' or 'PGEToolbox.tar.gz' file;
2. Unzip the files in a local directory keeping the structure of sub-directories. A
3. directory called `pgettoolbox` will be created;
4. Start up Matlab;
5. Add the local directory `./pgettoolbox` to the Matlab path;
6. From Matlab run command `PGEGUI`. This should bring up the main menu GUI of the toolbox.

## *Tutorial and Help System*

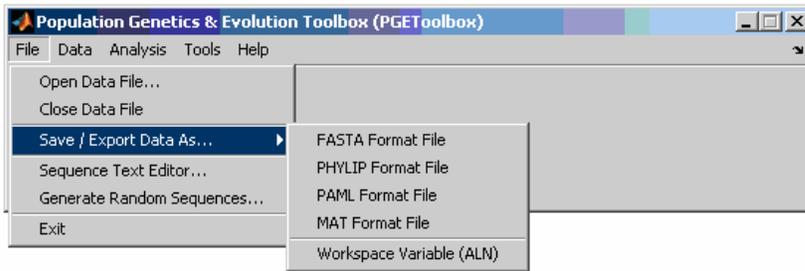
PGEToolbox provides a comprehensive tutorial and help system. Command `PGEDEMO` brings up slideshow-style demos for sequence-based and SNP-related analysis. Command `help pgettoolbox` lists the name of functions and brief introduction of those functions in the command line window. Command `help` or `edit` followed by a function name gives the usage description or the source code of the function. PGEToolbox website contains a step-by-step tutorial and documentation of functions. `PGEGUI` contains a help menu, in which user can bring up Matlab build-in help browser. The version checker under the same menu allows user to check for the latest updates.

## *Supplementary Figures*

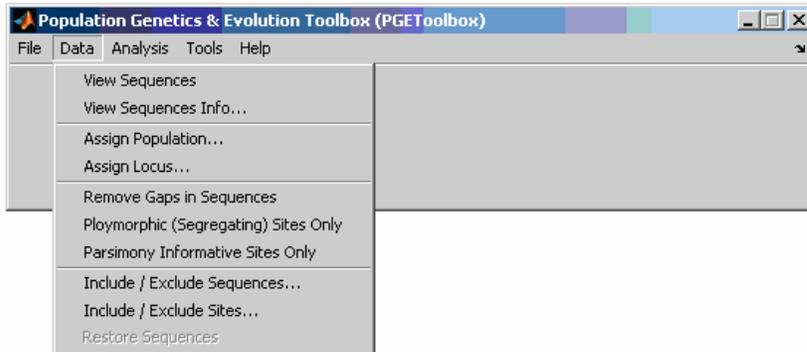
*Figure S1 - PGEToolbox GUI*

(A) File submenu; (B) Data submenu; (C) Analysis submenu; and (D) Tools submenu

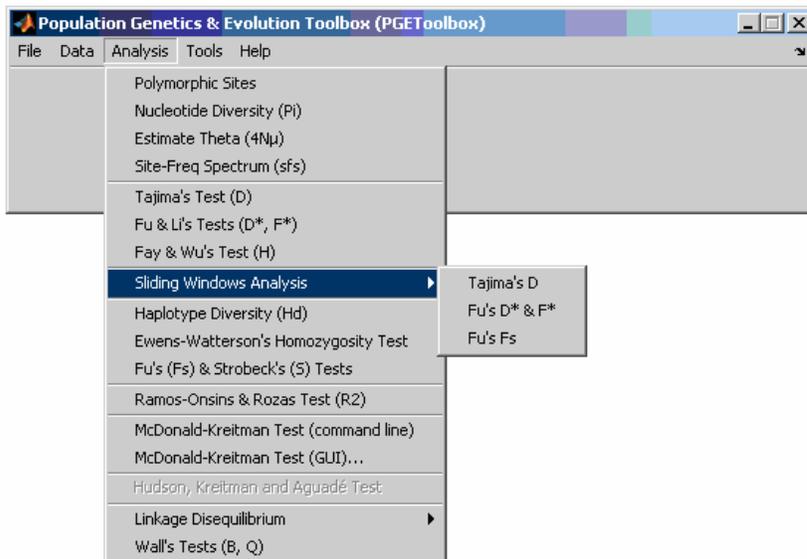
(A)



(B)



(C)



(D)

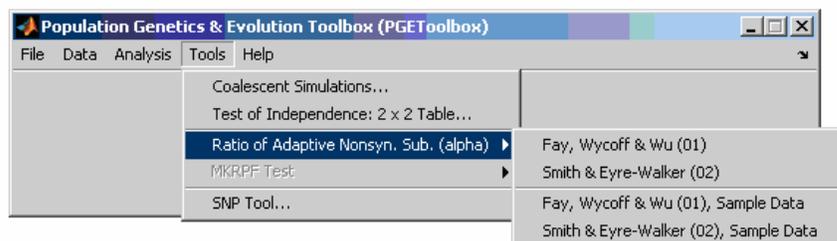


Figure 2 - Example output from function `estimatetheta`

The polymorphism sequences were randomly generated. The results are consistent with those from DnaSP.

```

=====
Estimation of theta (i.e., 4Mu)
=====
Theta-W (from Eta*) = 3.065693, var = 3.981063 (per sequence)
                    = 0.005135 (per site)
Theta-W (from S**)  = 3.065693, var = 3.981063 (per sequence)
                    = 0.005135 (per site)
Theta-Pi***        = 2.333333, var = 2.903704 (per sequence)
                    = 0.003908 (per site)

* Eta is the total number of mutations
** S is number of segregating sites
*** Theta-Pi equals nucleotide diversity
=====

```

Figure 3 - Relationship among calculation, simulation and testing functions

Using functions `tajima89d`, `tajima89d_simu` and `tajima89d_test` as example.

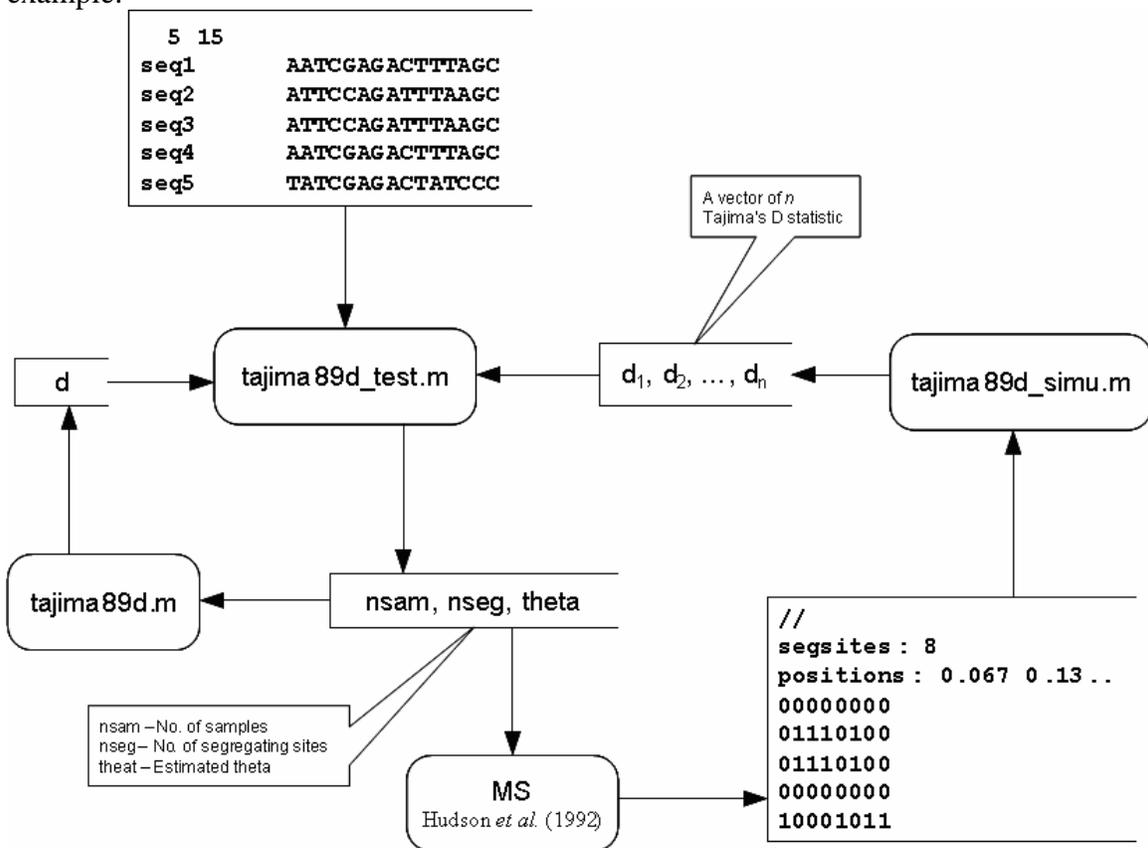


Figure 4 - Coalescent simulation dialog and histogram of result

(A) Dialog of parameter input for coalescent simulation. Simulated data will be generated by using 30, the observed number of segregating sites in the sample and under the conservative assumption of no recombination. (B) Histogram of estimated statistics (here, Tajima's  $D$ ) from simulated replicates.

(A)

Coalescent Simulation

Sample size (nsam): 10      No. replicates (nreps): 1000

Theta (-t): 12.5

Segregating Sites (-s): 30

No Recombination

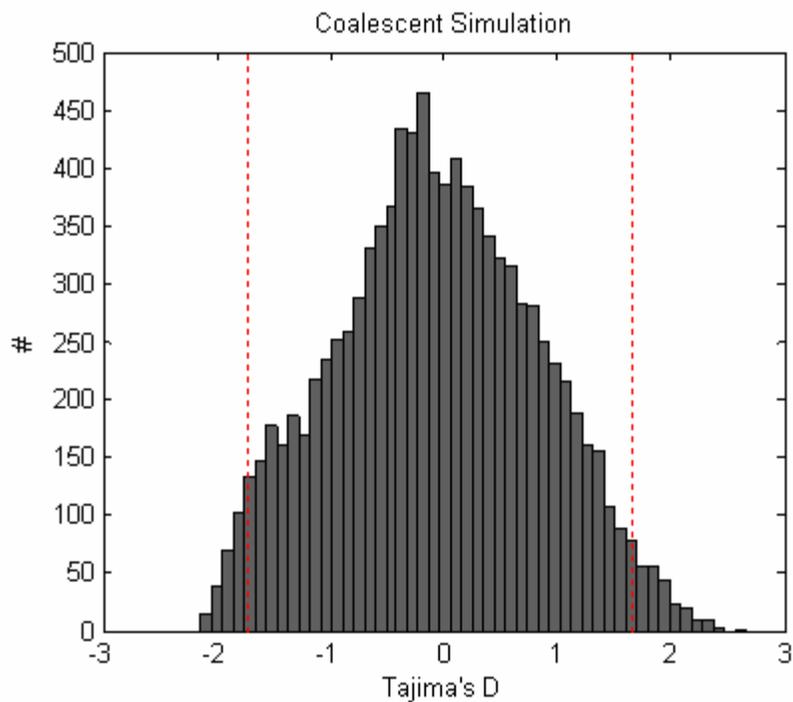
rho, 4Nr (-r): 10.0

No. Sites (nsites): 2000

Statistic: Tajima's D

Run      Cancel

(B)



*Figure 5 - Example results from SNP-related functions*

(A) Visual genotype (VG) view. In each panel, a graphical representation of genotypes is shown for the CHB (Chinese individuals from Beijing) and CEU (CEPH trios from Utah) samples. Rows correspond to individuals and columns denote SNPs. For each SNP, blue, yellow, and red boxes indicate whether the individual is homozygous for the common allele, heterozygous, or homozygous for the rare allele, respectively. Cyan boxes indicate missing data. The SNPs are from human locus EDAR, in which strong signature of positive selection has been identified in the CHB sample [32]. (B) Plot of EHH for two core haplotypes of the single SNP, rs9819197, with haplotype data for HapMap CEU population. Red dash line indicates the position of the core SNP.

