

## **slidingWindows\_v4.py**

*David Field 31/05/2017*  
*david.field@univie.ac.at*

*What does this script do?*

This script calculates various diversity ( $\pi_w$ ) and differentiation statistics ( $F_{ST}$ ,  $\pi_b$  /  $D_{xy}$ ) for individual or pooled sequence data for two or more populations. Two types of output are generated: (i) individual site statistics, (ii) sliding window averages of the same statistics (with adjustable overlaps). Different data manipulations with regards to data quality can be set by the user including, minimum and maximum depth and the minimum number of populations with required depth range for site inclusion. The script can also work with individual based data, as long as this is in the required format.

*How to install*

Simply place the script and all required input files in a common directory. Python modules required include: *re*, *sys*, *getopt*, *itertools*, *os*, *math*.

Ideally place this in a directory included in your PATH. See online how to edit your PATH

Some Python modules that may not come with base installation are required. If some modules are missing, download and install the required modules with, for example:

```
> sudo apt-get install python-numpy python-scipy
```

*Input files*

- **Genomic data:** in this version 4, the format is fixed on \*.sync files as used for pooled data in the program Popoolation2 (Kofler et al., 2010).
- **Scaffold data:** a file listing the set of scaffolds or chromosomes to analyse. This will assume \*.sync files exist with the same names.
- **Population details:** a file listing population names and sample sizes.

*How to run the program*

Assuming you have the script and the required input files in your PATH directory, you can run the program as follows:

```
> python SlidingWindows_v4.py scaffold_RosSample.txt  
popDetails_RosSample.txt LG6_RosSample 15 200 2 2 10000 0 1 0 1
```

where the system arguments for the script include (corresponding values in example):

1. list of scaffolds to analyse (scaffold\_RosSample.txt).
2. file with population names and sample size (popDetails\_RosSample.txt).
3. output file name for sliding window statistics (LG6\_RosSample.txt).
4. the minimum depth to include a site (15).
5. the maximum depth to include a site (200).
6. minimum copies for an allele call, i.e.  $> 1$ , no singletons (2).
7. minimum number of populations (or pools) with the minimum depth to include a site (2).
8. window size in base pairs (10000).
9. overlap size between windows (5000).
10. keep site specific file, whether to keep files for each scaffold/chromosome (1 = Y, 0 = N).
11. window analysis only (1 = Y, 0 = N). If repeated window analysis on same sites file.
12. site analysis only (1=Y, 0 = N)
13. keep all positions (1 = Y, 0 = N). For more compact sites output file (polymorphic sites only) then choose 0, for all sites choose 1.

*Some important details:*

- min depth: if running on individual based data, make minimum depth the size of the smallest number of haploid genomes (i.e.  $4 \times 2 = 8$ ),
- min copies: if running on individual based data, make minimum copies = 1 (i.e. allow for singletons),
- input: scaffold input list must have line breaks saved as Unix (LF). Use text wrangler or similar text editor,
- window analysis only: only works if a site specific file(s) already created and present in the same folder,
- min scaffold size: skips over any scaffold  $< 15\text{kB}$ . Ensures at least two windows minimum.

*Description of output for site output file*

- **scaffold**: genome scaffold/chromosome
- **position**: position on scaffold/chromosome
- **LG**: linkage group
- **cM**: physical distance along linkage group (centri Morgans)
- **ref**: reference genome allele
- **bases\_t**: variants detected at the site
- **alleles\_num**: number of allelic variants detected
- **1or2bases**: whether 1-2 bases detected (1) or more than two (0)
- **minDepth\_pools**: how many populations/pools passed filters

- **minDepth\_pass**: whether the minimum depth threshold was achieved at this site for the minimum number of populations (1=yes, 0=no).

#### *Single population statistics*

- **p\_1,..., p\_i**: allele frequencies in the *i*th population/pool.
- **poly\_pool\_1,..., poly\_pool\_i**: the average proportion of polymorphic sites in the *i*th population/pool.
- **pi\_1, ..., pi\_i**:  $\pi_w$  for individual based population data in the *i*th population.  
Here,  $\pi_i = \pi_w = 1 - \sum_{i=1}^n p_i^2 = 2(p_j)(1 - p_j)$ .
- **piAdj\_1,..., piAdj\_i**:  $\pi_w$  for pooled data in the *i*th population.  
Here, adjustment for binomial sampling where  $\pi_{Adj\_i} = \pi_w[(M-1)/(M-2b+1)]$ , where  $M$  = depth,  $b$  = minimum read count for allele call.
- **dpthAdj\_1,..., dpthAdj\_i**: mean depth within window for each population *I*, adjusted for binomial sampling
- **dpthPass\_1,..., dpthPass\_i**: what proportion of the window is covered in population *i*

#### *Population pair-wise statistics*

- **piBar\_1\_2,..., piBar\_i\_j**: average  $\bar{\pi}_w$  across each pairwise population/pool comparison,  $(\pi_1 + \pi_2)/2$ .
- **piBarAdj\_1\_2,..., piBarAdj\_i\_j**: average  $\bar{\pi}_w$  (as above) but adjusted for binomial sampling in pooled data. Where,  $\pi_{BarAdj\_i\_j} = \pi_{Bar\_i\_j}[(M-1)/(M-2b+1)]$
- **pi\_T**: total for each pairwise pool comparisons, e.g.  $\pi_{T\_1\_2} \dots$  **pi total = 2 x pbar x qbar** (i.e. 2 x mean allele p frequency x mean q allele frequency)  
 $\pi_{T\_adj\_1\_2}$  (i.e. 2 x mean allele p frequency x mean q allele frequency) adjusted (binomial errors) for each pairwise pool comparisons. e.g.  $\pi_{T\_adj\_1\_2}$ .
- **pi\_TAdj = (2 x pbar x qbar) x [(M-1)/(M-2b+1)]**
- **Dxy\_raw**:  $d_{XY\_raw\_1\_2} \dots$  **Dxy = (p1\*q2 + q1\*p2)**. Not recommended
- **Dxy\_fromFstadj**: but instead calculated backwards from  $F_{st}$  and  $\pi_{Bar}$ . e.g.  $d_{XY\_fromFst\_1\_2} \dots$  **Dxy = piBar\*(1-Fst/1+Fst)**
- pairwise relative divergence  $F_{st}$ , calculated from each site from  $\pi_i$  values then averaged for the window (note this is not what Popoolation does). e.g.  $F_{st\pi\_1\_2} \dots$  **Fst = (piTotal - pibar)/piT**
- **FstPiAdj**: pairwise relative divergence  $F_{st}$  but with adjusted values of  $\pi_i$ , calculated from each site from  $\pi_i$  values then averaged for the window (note this is NOT what Popoolation does). e.g.  $F_{st\pi_{Adj}\_1\_2} \dots$  **Fst = (piTotal\_adj - pibarAdj)/piTotal\_adj**
- **Fst** pairwise relative divergence  $F_{st}$  back calculated from raw  $D_{xy}$  values. e.g.  $F_{st\_fromDxy\_1\_2} \dots$  **Fst = (Dxy - piBar)/(Dxy + piBar)**. Not recommended as this comes from site averages.

- **FstfromMeanPi**: pairwise relative divergence Fst from average Pi values from the window,  $F_{st} = \frac{(\pi_{Total\_window} - \pi_{bar\_window})}{\pi_{Total\_window}}$ . (this is what we use for individual sequence data)
- **FstfromMeanPiAdj\_1\_2**: pairwise relative divergence from average adjusted PiAdj values from the window,  $F_{st} = \frac{(\pi_{Total\_adj\_window} - \pi_{barAdj\_window})}{\pi_{Total\_adj\_window}}$ . (this IS what we use for the pooled data)
- **DxyfromFst**: pairwise divergence Dxy, but instead calculated from backwards from Fst values above. e.g.  $d_{XY} = \pi_{bar} \times [(1 + F_{st}) / (1 - F_{st})]$ . (this is what we could use for individual data)
- **DxyfromFstAdj\_1\_2,...**: pairwise divergence Dxy, but instead calculated backwards from Fst adjusted values above,  $d_{XY} = \pi_{barAdj} \times [(1 + F_{stAdj}) / (1 - F_{stAdj})]$ . (this is what we could use for pooled data)
- last two are Da values, see Cruickshank and Hahn 2014 for details