

# polyDFE

## User Manual, BETA Version

Paula Tataru

November 9, 2017

**polyDFE** infers the distribution of fitness effects (DFE) of new mutations from polymorphism, and, if available, can also incorporate divergence data obtained from an outgroup. The polymorphism data is provided as an unfolded site frequency spectrum (SFS). Note that, in order to obtain an unfolded SFS, at least one outgroup is needed. However, including the divergence data when inferring the DFE can lead to bias in the estimated parameters (see ? for further discussion).

Once the DFE is obtained, **polyDFE** also calculates  $\alpha$ , the rate of adaptive molecular evolution, commonly defined as the proportion of fixed adaptive mutations among all non-synonymous substitutions. This code implements the program described in ?.

## Contents

# 1 Authors

The polyDFE has been developed and implemented by Paula Tataru and Marco A.P. Franco.

# 2 License

polyDFE is open source, distributed under the GNU General Public License, version 3. See `LICENSE.txt` for details.

# 3 Requirements

polyDFE is implemented in C and uses the GSL library. This should be downloaded from <https://www.gnu.org/software/gsl/> and compiled accordingly. polyDFE has been tested using GSL-1.16.

# 4 Installation

The simplest way to compile polyDFE is

1. `cd` to the directory containing the `makefile` and type `make all` to compile the code.
2. You can remove the program binaries and object files from the source code directory by typing `make clean`.

If the GSL library is not installed in the default directory, than the installation path needs to be specified in the `makefile` by updating `USR_INC` and `USR_LIB` and uncommenting the lines at the top of the `makefile`.

# 5 Running polyDFE

Running polyDFE with the argument `-h` will print out the usage (required and optional arguments)

```
$ ./polyDFE -h
Usage: ./polyDFE -d data_file [-m model(A, B, C, D) [K]] [-t]
      {-s m_neut L_neut m_sel L_sel n -i init_file ID ||
      [-o optim(bfgs, conj_pr, conj_fr)] [-r range_file ID]
      [-i init_file ID [-j]] [-e] [-w] [-g grouping_file ID]
      [-p optim_file ID] [-b [basinhop_file ID]]
      [-l min] [-v verbose(0, 1, frequency)]}
```

polyDFE runs in two modes: `-s`, used for simulating data under the hierarchical probabilistic model, and `-o`, used for estimating the DFE and  $\alpha$ . If neither `-s` or `-o` are given, then `-o bfgs` is used by default (see argument `-o` on page ?? for details). If both `-s` and `-o` are given, polyDFE will terminate with an error. Optional arguments are given between `[]`.

## 5.1 Shared arguments

A few arguments are shared between the two modes of polyDFE.

Regardless of the mode polyDFE is ran in, the following argument is required

- **-d data\_file**: path to the data file.

**data\_file** is where **polyDFE** will write the simulated data, if **-s** is used, or will read the input data, if **-o** is used.

The file can contain any number of empty lines and comment lines that start with **#**. Excluding these, the structure of **data\_file** is

$$\begin{array}{c}
 \begin{array}{ccccccc}
 & m_{\text{neut}} & m_{\text{sel}} & n & & & \\
 m_{\text{neut}} \text{ lines} & \left\{ \begin{array}{ccccccc}
 p_{\text{neut}}(1) & p_{\text{neut}}(2) & \dots & p_{\text{neut}}(n-1) & l_{\text{neut}} & d_{\text{neut}} & l_{\text{neut}}^d \\
 \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots
 \end{array} \right. \\
 m_{\text{sel}} \text{ lines} & \left\{ \begin{array}{ccccccc}
 p_{\text{sel}}(1) & p_{\text{sel}}(2) & \dots & p_{\text{sel}}(n-1) & l_{\text{sel}} & d_{\text{sel}} & l_{\text{sel}}^d \\
 \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots
 \end{array} \right.
 \end{array}
 \end{array}$$

where

- $m_{\text{neut}}$  and  $m_{\text{sel}}$  are the number of fragments that are assumed to be evolving neutrally or under selection, respectively;
- $n$  is the ingroup sample size (number of sequences);
- $p_z(i)$  is the  $i^{\text{th}}$  entry in the SFS within the fragment, for sites that are assumed to be evolving neutrally,  $z = \text{neut}$ , or under selection,  $z = \text{sel}$ ;
- $l_z$  is the total number of sites within the fragment used for the SFS (number of successfully sequenced sites within the ingroup) that are assumed to be evolving neutrally,  $z = \text{neut}$ , or under selection,  $z = \text{sel}$ ;
- $d_z$  is the total number of fixed mutations in the ingroup relative to the outgroup (divergence counts) within the fragment, that are assumed to be evolving neutrally,  $z = \text{neut}$ , or under selection,  $z = \text{sel}$ ;
- $l_z^d$  is the total number of sites within the fragment used for the divergence counts (number of successfully sequenced sites within the ingroup and outgroup) that are assumed to be evolving neutrally,  $z = \text{neut}$ , or under selection,  $z = \text{sel}$ ;  $l_z$  and  $l_z^d$ , for a given fragment, could be different if certain sites are successfully sequenced in the ingroup, but not in the outgroup.

The divergence data (last two columns,  $d_z$  and  $l_z^d$ ) can be absent from the file. If it is absent, it should be absent for all fragments.

**polyDFE** uses data divided into fragments in order to model variability in mutation rates, as described in ?. If only one neutral and one selected fragments are provided, than variability in mutation rates is not modeled.

The provided files **input/example\_1**, **input/example\_2**, **input/example\_3**, **input/example\_4** and **input/example\_5** are examples of **data\_file**.

Two optional arguments are shared between the the modes of **polyDFE**

- **-m model(A, B, C, D) [K]**: assumed DFE model.

**polyDFE** assumes that the DFE  $\phi$  takes one of four specific functional forms, encoded **A**, **B**, **C**, and **D**, described below

- **A**: the DFE is given by a reflected displaced Gamma distribution, parameterized by  $\bar{S}$ ,  $b$  and  $S_{\text{max}}$ , with density

$$\phi_A(S; \bar{S}, b, S_{\text{max}}) = \begin{cases} f_{\Gamma}(S_{\text{max}} - S; S_{\text{max}} - \bar{S}, b) & \text{if } S \leq S_{\text{max}} \\ 0 & \text{otherwise} \end{cases}$$

where  $\bar{S}$  is the mean of the DFE,  $b$  is the shape of the Gamma distribution,  $S_{\text{max}}$  is the maximum value that  $S$  can take, and  $f_{\Gamma}(x; m, b)$  is the density of the Gamma distribution with mean  $m$  and shape  $b$ .

- **B**: the DFE is given by a mixture of a Gamma and discrete distributions, parameterized by  $S_d$ ,  $b$ ,  $p_b$  and  $S_b$ , where

$$\phi_B(S; S_d, b, p_b, S_b) = \begin{cases} (1 - p_b)f_\Gamma(-S; -S_d, b) & \text{if } S \leq 0 \\ p_b & \text{if } S = S_b \\ 0 & \text{otherwise} \end{cases}$$

where  $S_d$  is the mean of the DFE for  $S \leq 0$ ,  $b$  is the shape of the Gamma distribution,  $p_b$  is the probability that  $S > 0$ , and  $S_b$  is the shared selection coefficient of all positively selected mutations, and  $f_\Gamma(x; m, b)$  is the density of the Gamma distribution with mean  $m$  and shape  $b$ .

- **C**: the DFE is given by a mixture of a Gamma and Exponential distributions, parameterized by  $S_d$ ,  $b$ ,  $p_b$  and  $S_b$ , where

$$\phi_C(S; S_d, b, p_b, S_b) = \begin{cases} (1 - p_b)f_\Gamma(-S; -S_d, b) & \text{if } S \leq 0 \\ p_b f_e(S; S_b) & \text{if } S > 0 \end{cases}$$

where  $S_d$  is the mean of the DFE for  $S \leq 0$ ,  $b$  is the shape of the Gamma distribution,  $p_b$  is the probability that  $S > 0$ ,  $S_b$  is the mean of the DFE for  $S > 0$ , and  $f_\Gamma(x; m, b)$  is the density of the Gamma distribution with mean  $m$  and shape  $b$ , while  $f_e(x; m)$  is the density of the Exponential distribution with mean  $m$ .

- **D**: the DFE is given as a discrete DFE, where the selection coefficients can take one of  $S_i$  distinct values,  $1 \leq i \leq K$ , where each value  $S_i$  has probability  $p_i$ , with

$$\sum_{i=0}^K p_i = 1$$

The value of  $K$  needs to be provided in the command line for model **D**.

If **-m** is not given, model **C** is used by default.

- **-t**: the total running time of **polyDFE** will be measured and printed on the standard output at the end.

**-i** is also common to the two modes. However, **-i** is required when simulating data (mode **-s**), while when inferring the parameters (mode **-o**), it is an optional argument for DFE models **A**, **B** and **C**. Note though that if using model **D**, then **-i** has to be given, as it contains the information about the selection coefficients  $S_i$ .

- **-i init\_file ID**: path to a file containing the values of the DFE and other parameters, and **id**.

**init\_file** can contain any number of empty lines and comment lines that start with **#**. Excluding these, the structure of each line is

id [0|1]  $\epsilon_{an}$  [0|1]  $\epsilon_{cont}$  [0|1]  $\lambda$  [0|1]  $\bar{\theta}$  [0|1]  $a$  < DFE parameters with flags > [0|1]  $r_2 r_3 \dots r_n$

Each line starts with a positive numerical **id**. This is used such that multiple configurations of the parameters can be stored in the same file. When running **polyDFE**, the desired line for initializing the parameters is specified through the **id** by setting **ID** to the right value.

The **id** is followed by the values of the parameters, in the order  $\epsilon_{an}$ ,  $\epsilon_{cont}$ ,  $\lambda$ ,  $\bar{\theta}$ ,  $a$ , the DFE parameters (as detailed for argument **-m** above), and the  $r_i$ ,  $2 \leq i \leq n$ , parameters, where  $n$  is the ingroup sample size.

All parameters are preceded by a flag, which is either 0 or 1. In the **-s** mode, the flag is not used, but this becomes important for the **-o** mode, where the flag indicates if a parameter should be estimated (flag 0) or should be kept fixed to the value provided in **init\_file** (flag 1). All  $r_i$  parameters share one flag. The  $r_i$  parameters can also be provided for  $i > n$ , but those values will be ignored. This is useful if the same line in **init\_file** is used for multiple data sets that have different sample sizes.

For DFE model **D**, only the probabilities  $p_i$  are estimated, while the selection coefficients  $S_i$  are fixed to the given values.

For details of the meaning of each parameter, please see **?**. The provided files **input/init\_model\_A.txt**, **input/init\_model\_BandC.txt** and **input/init\_model\_D.txt** give examples of the structure of the **init\_file** for the different DFE models.

Note that the  $a$  parameter is directly linked to variability of mutation rates. If it is wished for mutation rates to be assumed constant, than  $a$  can be fixed to  $-1$ . If both  $m_{neut}$  and  $m_{sel}$  provided in the input

datafile (through argument `-d`) are equal to 1 (i.e. the data is not divided into fragments), then the mutation variability is not estimated by default. When mutation variability is allowed, then the observed data (counts) are assumed to follow a negative binomial distribution, otherwise they are assumed to follow a Poisson distribution. See ? for more details on mutation variability.

The parameter  $\epsilon_{\text{cont}}$  allows for modeling the SFS for the sites that are assumed to be evolving neutrally as a mixture of neutral and selected sites. Note that this parameter is not identifiable and it is therefore recommended that its flag is always 1. However, this parameter can be used to test the robustness of the method at the assumption of neutrality: one can check how much the estimated DFE and  $\alpha$  would change if  $\epsilon_{\text{cont}}$  was set to some value different than 0.

The provided files `input/init_model_A.txt`, `input/init_model_BandC.txt` and `input/init_model_D.txt` are examples of `init_file`.

## 5.2 Simulating data

The simulation mode (`-s`) also requires the following argument

- `-s m_neut L_neut m_sel L_sel n`: simulate data.

When using `-s`, simulated data will be printed in `data_file` with  $m_{\text{neut}} = \text{m\_neut}$ ,  $m_{\text{sel}} = \text{m\_sel}$ ,  $l_{\text{neut}} = l_{\text{neut}}^d = \text{L\_neut}$  and  $l_{\text{sel}} = l_{\text{sel}}^d = \text{L\_sel}$ , for all simulated fragments. For details on  $m_z$  and  $l_z$ , with  $z \in \{\text{neut}, \text{sel}\}$ , see argument `-d` on page ??.

## 5.3 Estimating DFE and $\alpha$

The estimation (optimization) mode (`-o`) allows for the estimation of DFE and  $\alpha$ . For this, next to the arguments described above, it allows for a series of optional arguments, which control the optimization, both which and how the parameters are estimated, but also when a set of parameters found is considered to be optimal.

- `-o (bfgs, conj_pr, conj_fr)`: optimization method to use for estimating the parameters.

`polyDFE` uses `GSL` to run a local optimization algorithm and find estimates of the parameters. `GSL` implements multiple optimization algorithms. Generally, optimization is more efficient if it also uses the derivative of the function to be optimized. Therefore, in `polyDE` we have chosen to use three of the algorithms offered by `GSL`, which all rely on the function (in this case, the data likelihood) derivatives. These are:

- `bfgs`: the Broyden-Fletcher-Goldfarb-Shanno (BFGS) algorithm.
- `conj_pr`: the Polak-Ribiere conjugate gradient algorithm;
- `conj_fr`: the Fletcher-Reeves conjugate gradient algorithm.

For details on the different optimization algorithms, please see the `GSL` manual. When `-o` is not given, `-o BFGS` is used by default. `polyDFE` has been tested using the BFGS algorithm.

- `-r range_file ID`: path to a file containing the allowed range of all parameters, and id.

`range_file` can contain any number of empty lines and comment lines that start with `#`. Excluding these, the structure of each line is

id     $k$      $\epsilon_{\text{an}}^{\min}$   $\epsilon_{\text{an}}^{\max}$      $\epsilon_{\text{cont}}^{\min}$   $\epsilon_{\text{cont}}^{\max}$      $\lambda^{\min}$   $\lambda^{\max}$      $\bar{\theta}^{\min}$   $\bar{\theta}^{\max}$      $a^{\min}$   $a^{\max}$     < ranges for DFE parameters >     $r^{\min}$   $r^{\max}$

Each line starts with a positive numerical id. This is used such that multiple configurations of the parameters can be stored in the same file. When running `polyDFE`, the desired line for the range of the parameters is specified through the id by setting `ID` to the right value.

The id is followed by  $k$ , which controls the transformation of the parameters (see below) and the range within each parameter should be estimated. Some parameters are intrinsically constrained (for example,  $p_b$  for DFE models B and C is, by construction, constrained to be between 0 and 1), while many parameters have to be positive (or negative, for example,  $S_d$  for DFE models B and C). In order to treat all parameters in the same way and to also allow the user to narrow the range in which the optimum parameters are searched, all parameters are then constrained to the range provided by the user through

the `range_file`. However, the optimization algorithms `polyDFE` uses (see argument `-o` on page ?? for details) are unconstrained, in that the parameters can take any value in the interval  $(-\infty, +\infty)$ . In order to change the original constrained optimization problem to an unconstrained one, `polyDFE` transforms the parameters from the range provided by the user to the interval  $(-\infty, +\infty)$ . For this, a generalized logistic function is used. This function is parameterized by a growth rate (or steepness of curve)  $k$ , which controls the spread of the transformation: the smaller the  $k$ , the more quicker the transformed values move from  $-\infty$  to  $+\infty$ . `polyDFE` has been tested with a value of  $k = 0.01$ . If the optimization algorithm cannot find a good solution (i.e. the gradient is small, see argument `-p` on page ?? for details), the  $k$  parameter can be changed in the hope of better performance.

Note that, in order for the optimization algorithm to successfully find the optimum parameters (the parameters that have the height likelihood), the ranges should be large enough to ensure that they cover the optimum values.

If initial values of the parameters provided through the `-i` argument or as calculated when the argument `-e` is provided are not within the ranges, the program will automatically update the ranges.

The provided files `input/range_model_A.txt`, `input/range_model_BandC.txt` and `input/range_model_D.txt` are examples of `range_file`.

- `-e`: automatically estimate the initial value of the parameters.

The optimization algorithms require an initial value for the parameters to be optimized. This can either be provided using `-i` if the user has a good idea for such values, or can, by using the argument `-e`. If `-e` is used without `-i`, then all the parameters (except for  $\epsilon_{\text{cont}}$ , which, by default, is not estimated and set to 0) will be automatically initialized and then optimized. However, if `-i` and `-e` are used together, the initial value of the parameters is estimated automatically only for those parameters that are flagged with 0 and are consequently optimized.

The parameters  $\bar{\theta}$ ,  $a$ ,  $\lambda$  and  $r_i$  can be estimated deterministically. The rest of the parameters are estimated using a grid search approach. However, for the DFE model D, the probabilities  $p_i$  are set to be uniform when `-e` is used.

If `-i` is not used, `-e` is used by default.

- `-j`: use jointly both the provided initial value and automatically estimated value of parameters.

When the argument `-j` is used together with `-i`, then the optimization algorithm is ran twice, once where the initial value of the parameters is the one given in the `init_file`, and a second time where the automatic estimation from `-e` is used. At the end, the parameters that have the best likelihood, found from one of the two runs, will be reported.

- `-w`: do not use divergence data.

When the argument `-w` is used, then the DFE and other parameters are estimated without using divergence data, as described in (?). If the divergence data is absent from the `data_file`, `-w` will be used by default.

- `-g grouping_file ID`: path to a file containing grouping information for the  $r_i$  parameters, and id.

`grouping_file` can contain any number of empty lines and comment lines that start with `#`. Excluding these, the structure of each line is

$$\text{id} \quad G \quad i_1 \quad i_2 \quad \dots \quad i_G$$

Each line starts with a positive numerical id. This is used such that multiple configurations of the grouping can be stored in the same file. When running `polyDFE`, the desired line for the range of the parameters is specified through the id by setting `ID` to the right value.

The `grouping_file` is used to provide information on grouping the  $r_i$  parameters. By default, each entry  $1 \leq i < n$  in the SFS, and the divergence count, have one value of  $r_i$  associated to it. If  $n$  is very large, this leads to a lot of parameters that need to be estimated. One could expect that the distortion that is incurred on the counts that are modeled using the  $r_i$  parameters could be similar for neighboring values (for example,  $i - 1$ ,  $i$  and  $i + 1$ ). Then only one  $r$  value can be estimated for those counts that share a similar distortion.

The grouping given in `grouping_file` will result in  $r$  values corresponding to the SFS entries as follows (where one interval represents one range for the SFS entries):

$$\begin{array}{ccccccc} r_1 & r_2 & r_3 & \dots & r_{G+1} \\ [1] & [2, i_1] & [i_1 + 1, i_2] & \dots & [i_G, n'] \end{array}$$

where  $n'$  is either  $n$ , if divergence data is used in the estimation (see argument `-w` on page ?? for details), or  $n - 1$  otherwise. For identifiability reasons,  $r_1$  is always set to 1 and not estimated.

The provided file `input/params.grouping` is an example of a `grouping_file`.

- `-p optim_file ID`: path to a file containing parameters that control the optimization algorithm, and id.

`optim_file` can contain any number of empty lines and comment lines that start with `#`. Excluding these, the structure of each line is

```
id       $\epsilon_{\text{abs}}$   step size  tol  max iter
```

Each line starts with a positive numerical id. This is used such that multiple configurations for the optimization algorithms can be stored in the same file. When running `polyDFE`, the desired line for the range of the parameters is specified through the id by setting `ID` to the right value.

The `optim_file` is used to control the behavior of the optimization algorithms from `GSL`.

- $\epsilon_{\text{abs}}$ : this is the absolute tolerance value that the gradient of the function to be optimized (here, the likelihood) is tested against. When the gradient is smaller than this value, than `GSL` considers that it has successfully found an optimum. If this value is too large, it will lead to an early termination of the algorithm. If the value is too small, it might lead to a longer running time of the algorithm without, potentially, much improvement on the function. However, it is always best to have an absolute tolerance value that is too small than too large. If a set of values found is a true optimum, then the gradient should be 0. Default value: 0.00001.
- step size: the `GSL` optimization algorithms use a line search and the step size determines the size of the first step performed in the search. Default value: 2.
- tol: specifies the accuracy of the line search, with the precise meaning depending on the algorithm used. Default value: 0.1.
- max iter: the `GSL` optimization algorithms are iterative. If the number of performed iterations reaches max iter, then the algorithm is stopped, regardless of the gradient value. Default value: 1500.

For more details on the above, please consult the `GSL` manual.

The provided file `input/params.optim` is an example of an `optim_file`.

- `-b [basinhop_file ID]`: path to a file containing parameters that control the basin-hopping algorithm, and id.

`basinhop_file` can contain any number of empty lines and comment lines that start with `#`. Excluding these, the structure of each line is

```
id  max same  max iter  temp  step  accept rate  interval  factor
```

Each line starts with a positive numerical id. This is used such that multiple configurations for the basin-hopping algorithm can be stored in the same file. When running `polyDFE`, the desired line for the range of the parameters is specified through the id by setting `ID` to the right value.

By default, `polyDFE` runs the optimization algorithm once (or twice, see argument `-j` on page ?? for details), and reports the results. However, the optimization algorithms are local, in that they only find an optimum that is in the neighborhood of the initial values of the parameters. To look more thoroughly for a global optimum, `polyDFE` allows the user to run the basin-hopping algorithm (???). `polyDFE` contains a reimplementation of the basin-hopping algorithm found in the Python library, `scipy` (??).

Basin-hopping attempts to find the global optimum by iterating the following steps

1. perturbing randomly the current value of the parameters;
2. running the local optimization algorithm (as chosen through `-o`, see page ?? for details) using the new values as initial values;
3. accepting or rejecting the found optimum based on the likelihood.

Basin-hopping has been shown to be very efficient for a wide variety of problems in physics and chemistry. Unfortunately, there is no way to determine if the true global optimum has actually been found, and therefore it is left to the user to ensure that. The algorithm is controlled through

- max same: If after max same iterations, the basin-hopping algorithm does not find an improved set of parameters, than the algorithm stops. Default value: 50.
- max iter: The basin-hopping algorithm is ran for a maximum of max iter iterations. Default value: 500.
- temp: The temp is the temperature used in the metropolis criterion when accepting or rejecting the new values. For best results, the temperature should be comparable to the typical difference in likelihood between local optima. Default value: 1.
- step: The step controls how far away from the current value the perturbed value is. This is crucial for the algorithm's performance. Ideally, it should be comparable to the typical separation between local optima of the likelihood. The algorithm implemented in `polyDFE` will automatically adjust the step, but it make take many iterations to find an optimal value. Default value: 50.
- accept rate: The target accept rate (percentage of new values that are accepted in step 3) for when adjusting the step. Default value: 0.5.
- interval: Every interval iterations, the basin-hopping algorithm adjusts the step. Default value: 10.
- factor: When the step is adjusted, it is done with this factor. Default value: 0.9.

The argument `-b` can be ran without `basinhop_file` ID. If only `-b` is given, then `polyDFE` runs basin-hopping with the default parameter values, given above.

The provided file `input/params.basinhop` is an example of a `basinhop_file`.

- `-l min`: limit running time to min.

To control the running time of `polyDFE`, one run of the optimization algorithm is stopped if it has used more than `min`. By default, `min` is set to 60. Note that if `min` is set too low, there is a high chance that the optimization algorithm will terminate without finding good optimal parameters.

- `-v verbose(0, 1, frequency)`: verbose frequency.

The argument `-v` controls how often information about the status of the optimization algorithm is printed on the screen. By default, `verbose` is set to 0, and only the initial and optimum values of the parameters are printed. Otherwise, every `verbose` iterations of the optimization algorithm, information is printed about the current value of the parameters, the corresponding likelihood, gradient and status of the optimization. Apart from the possible status values that `GSL` uses, `polyDFE` has the additional values:

- -3: local optimization is stuck in the same parameter values (it cannot further improve on the current values).
- -4: restarting the local optimization lead to the same parameter values.
- -5: the local optimization is being restarted; sometimes, when the local optimization cannot further improve on the current values, restarting the optimization can allow it to proceed further. If the local optimization is stuck in the same values, restarting is done at most 3 consecutive times.
- -6: the local optimization used at least `min` minutes and will be stopped.
- -7: the likelihood evaluation returned NAN and the local optimization is therefore stopped. This can happen for certain values of the parameters, where the numerical integration over the DFE fails.
- -8: a maximum number of likelihood evaluations has been reached, and the local optimization has been forcefully stopped. If this status appears, is indicates an issue with the local optimization procedure and the obtained results are not reliable.

## 6 Post-processing using R

`polyDFE` is accompanied by a series for R functions that allow

- parsing of the output of `polyDFE` for easy manipulation in R
- creating `init_file` containing the best found parameters
- creating grouping for the  $r$  parameters
- calculating  $\alpha$  using alternative definitions
- performing model testing



## 6.1 Parsing the output of polyDFE

The function

```
parseOutput(filename, init = FALSE)
```

allows the user to parse the output of **polyDFE** from multiple runs stored in **filename**. The function returns a list with entries for each separate run of **polyDFE** found, containing

- the DFE model used;
- the best likelihood found and the corresponding gradient;
- the values of all parameters (including those that have been fixed using the flag 1 in **init\_file**);
- which parameters have been estimated;
- the value of  $n$  (the sample size of the data);
- expected SFS and, if relevant, divergence counts and missattributed polymorphism;
- estimates of  $\alpha$ .

If **init = TRUE**, then the function returns the values for the local optimization where the parameter values found in **init\_file** where used as initial values for the parameters during the optimization (see argument **-e** on page ?? for details), regardless if this returned the best likelihood or not.

## 6.2 Creating init\_file

The function

```
createInitLines(estimates, outputfile, startingID, fix, groupingDiff = 0)
```

creates an **init\_file** with the given values of the parameters.

**estimates** can either be a list returned by **parseOutput**, or a file name, as for **parseOutput**. For each run of **polyDFE** found, a new line is appended to **<outputfile>\_init** containing the values of the best parameters found in the format described previously (see argument **-i** on page ?? for details). The IDs for each line are consecutive and start at **startingID**. The parameters given in **fix** are flagged with 1, while the rest of the parameters are flagged with 0.

Grouping of  $r$  parameters (see argument **-g** on page ?? for details) can be calculated automatically by setting **groupingDiff**. If the difference between  $r_i$  and  $r_{i+1}$  is smaller than **groupingDiff**, than they are placed in the same group. The grouping information is appended to **<outputfile>\_grouping**.

## 6.3 Calculating $\alpha$

The function

```
estimateAlpha(estimates, supLimit = 0, div = NULL, poly = FALSE)
```

Using the parameters found in **estimates**, which is one entry from the list returned by **parseOutput**, this function calculates both  $\alpha_{\text{div}}$  (when **div = NULL**) and  $\alpha_{\text{dfe}}$  (when **div** is given) (?).

**?** argues that mutations with positive selection coefficients that are not higher than a certain threshold should, effectively, be treated as neutral mutations when calculating  $\alpha$  (see also ?). This threshold can be controlled here by setting **supLimit**.

When used, **div** should contain divergence data, which can be read from **filename** using the function

```
parseDivergenceData(filename)
```

As described in ?, divergence data can contain missattributed polymorphism. This can be corrected for by setting `poly = TRUE`.

## 6.4 Model testing

The function

```
compareModels(est1, est2 = NULL, nested = FALSE)
```

compares the models found in `est1` and `est2` by calculating the AIC and, where relevant, the LRT.

Both `est1` and `est2` can either be lists returned by `parseOutput`, or file names, as for `parseOutput`. If these contain more than one run of `polyDFE`, then run *i* from `est1` is compared to run *i* from `est2`.

The function can automatically detect if models are nested for calculating the LRT, however nestedness can be enforced by setting `nested = TRUE`. Caution should be used when the *r* parameters have been grouped (see argument `-g` on page ?? for details), as in this case the function cannot correctly detect nestedness.

The function returns a list containing two entries

- AIC is a matrix containing the number of estimated parameters, the log likelihood and AIC for `est1` in columns 1 – 3 and for `est2` in columns 4 – 6, for all runs of `polyDFE` found;
- LRT is a matrix containing the degrees of freedom in column 1, the likelihood of `est1` in column 2, the likelihood of `est2` in column 3 and the p-value from the LRT test in column 4, for all runs of `polyDFE` found.

If `est2 = NULL`, only the AIC is calculated for `est1`.

## 6.5 Examples

## References

- N. Galtier. Adaptive protein evolution in animals and the effective population size hypothesis. *PLoS genetics*, 12(1), 2016.
- E. Jones, T. Oliphant, P. Peterson, et al. SciPy: Open source scientific tools for Python, 2001–. URL <http://www.scipy.org/>. [Online; accessed 2016-09-21].
- E. O. Purisima and H. A. Scheraga. An approach to the multiple-minima problem in protein folding by relaxing dimensionality: Tests on enkephalin. *Journal of Molecular Biology*, 196(3):697–709, 1987.
- P. Tataru, M. Mollion, S. Glemin, and T. Bataillon. Inference of distribution of fitness effects and proportion of adaptive substitutions from polymorphism data. *Genetics*, accepted, 2017.
- D. J. Wales and J. P. Doye. Global optimization by basin-hopping and the lowest energy structures of Lennard-Jones clusters containing up to 110 atoms. *The Journal of Physical Chemistry A*, 101(28):5111–5116, 1997.
- D. J. Wales and H. A. Scheraga. Global optimization of clusters, crystals, and biomolecules. *Science*, 285(5432):1368–1372, 1999.