

empiricIST: An Integrated Software and analysis Tool for analyzing time-sampled sequence data such as EMPIRIC

README

empiricIST: An Integrated Software and analysis Tool for analyzing time-sampled sequence data such as EMPIRIC

Sebastian Matuszewski, Jeffrey D. Jensen and Claudia Bank

Lausanne, January 11, 2016

Contents

1	Introduction	3
2	Input data	3
2.1	empiricIST_MCMC_Input.py	4
3	Usage	8
3.1	Compilation	8
3.2	Execution	11
4	MCMC output	14
4.1	Combining Files	19
5	DFE tail shape estimation	22

1 Introduction

This program serves as an extension of the Bayesian Monte Carlo Markov chain (MCMC) method described in Bank et al. (2014) for estimating selection coefficients (growth rates) from engineered-mutation-driven experimental evolution data. These data are based on methods – such as EMPIRIC – in which specific mutations are introduced and compared against each other and the wild-type. All mutants (and the wild-type) are assumed to have evolved together in bulk culture for a number of generations with samples taken throughout the course of the experiment. Growth rates – and thus the selection coefficient of each mutation – can then be estimated from the number of reads obtained from deep sequencing. The motivation for the *empiricIST* software package is to provide an integrative framework for the analysis of experimental evolution data, and includes separate programs for processing raw sequence data and correcting for sequencing errors, obtaining statistically meaningful estimates of selection coefficients in a fast and efficient manner, and for providing ready-to-use summary statistics of the MCMC analysis and its associated parameter estimates. It is highly advisable to read the two accompanying papers by Bank et al. (2014) and Matuszewski et al. (in prep) that describe the methods and underlying assumptions in more detail before using the software.

2 Input data

The `empiricIST_MCMC` program only takes files in csv-format (comma-separated values) as input, which, furthermore, need to match a specified format (e.g., line endings need to be UNIX specific; input data needs to be ordered in a specified way). However, as part of the *empiricIST* software package, we provide a python script – `empiricIST_MCMC_Input.py` – that adjusts the raw data to match the specific input format needed for the MCMC simulation program (including other options discussed in more detail below). A minimal example of the raw data that is required to generate the `empiricIST_MCMC` input file is depicted in Figure 1. Note that the raw data itself also needs to be csv-formatted with column entries separated by a comma (','), a semi-colon(';') or a tab ('t'). Furthermore, the raw data file needs to have (exactly) one column either called 'sequence', 'Sequence', 'seq' or 'Seq' with at least two rows (the wild-type reference and a mutant) and at least additional 3 columns (i.e., time points) corresponding to the number of sequencing reads, with respectively named header cells.

seq	4.8	7.2	9.6	12	16.8	26.4	36
CCGGTCAAAACGGTTGGTCTGCTAACATGGAAA	24901	28500	48710	58076	46121	52651	104330
CCGGTAACAACGGTTGGTCTGCTAACATGGAAA	626	738	1515	1497	1417	1928	2512
CCGGTAAGAACGGTTGGTCTGCTAACATGGAAA	579	499	1116	1510	1322	2080	3444
CCGGTAATAACGGTTGGTCTGCTAACATGGAAA	532	642	1198	1414	1151	1596	2210
CCGGTACAACGGTTGGTCTGCTAACATGGAAA	727	861	1721	1897	1752	2506	4040
CCGGTACCAACGGTTGGTCTGCTAACATGGAAA	1358	1536	2899	3046	3315	4906	7384
CCGGTACGAACGGTTGGTCTGCTAACATGGAAA	892	999	1979	2277	2194	3168	4896
CCGGTACTAACGGTTGGTCTGCTAACATGGAAA	880	1091	1957	2029	2112	3081	4727
CCGGTAGAAACGGTTGGTCTGCTAACATGGAAA	441	443	887	1235	1075	1645	2594
CCGGTAGCAACGGTTGGTCTGCTAACATGGAAA	505	633	1194	1631	1355	1948	2546
CCGGTAGGAACGGTTGGTCTGCTAACATGGAAA	418	431	907	1236	1082	1769	2924

Figure 1 – Schematic illustration of the minimal data needed to run the `empiricIST_MCMC` program.

2.1 `empiricIST_MCMC_Input.py`

The program `empiricIST_MCMC_Input.py` is written in Python (2.7) and serves as a link between the (raw) time-sampled sequence data (e.g., obtained from EM-PIRIC) and the `empiricIST_MCMC` simulation program for the estimation of mutant growth rates 'r'. While it primarily ensures that the input data matches the input format required by the MCMC simulation program, it comes with additional options that will be detailed here.

The general usage is as follows: After opening a command-line interface (e.g., Shell, Terminal) and navigating to the location of the `empiricIST_MCMC_Input.py` file, the program can be executed by typing

```
python empiricIST_MCMC_Input.py [options] .
```

Note that this requires that the 'PATHVARIABLE' for Python has been set correctly on your system. Please consult the online Python documentation for further details (<https://docs.python.org/2/>). Without specifying any options the program will exit with an error and provide a short documentation on its usage, as it requires the name of a (raw) data input file and the start of the sequencing read data table to be specified (by invoking the '-f' and '-s' option, respectively). All options and their usage are given in Table 1.

Table 1 – A summary of the options of the `empiricIST_MCMC_Input.py` program.

Short/Long option	Accepted values	Description
-h, -help	none	When the '-h'-option is invoked, a short documentation on the usage of the program is shown. Note that, if this option is invoked, the python program is not executed.

-f, --file=	string	The '-f'-option is a mandatory option, which passes the name of the (raw) data input file (csv formatted) to the python program. Files created by the python program will take the name of the input file and add option-dependent specific file identifiers.
-s, --skipcol=	integer	The '-s'-option is a mandatory option, which takes an integer value corresponding to the number of descriptive columns that precede the actual 'data matrix' of sequencing read counts. For example, for the raw data depicted in Figure 1, the user would have to pass '-s 1' (or equivalently '-skipcol=1'). Please note that the data matrix must always span all remaining columns.
-o, --outlier=	'detect' or 'impute'	When the '-o'-option is invoked, the python program will perform an outlier analysis. If '-o detect' (or equivalently '-outlier=detect'), the python program performs a log-linear regression analysis for all mutants. Data points are then classified as outliers on the basis of the DFBeta statistic with a cut-off value of 2 (with data points surpassing this cut-off regarded as outliers). For more details please consult Bank et al. (2014). If '-o impute' (or equivalently '-outlier=impute'), the python program performs a log-linear regression analysis for all mutants. Data points are then classified as outliers on the basis of the DFBeta statistic with a cut-off value of 2 and their studentized residuals with a cut-off value of 3 (with data points surpassing <i>both</i> cut-offs are regarded as outliers) and imputed as described in Matuszewski et al. (in prep). Furthermore, an additional output file – whose name consists of 'ImputedData' along with the input file name – is produced that lists all imputed data points. In particular, the output is a simple matrix where rows denote different mutants and columns correspond to the different time points. An entry of '1' indicates that the data has <i>not</i> been imputed; an entry of '0' indicates that the data point has been imputed.
-l, --leadseq=	integer	When the '-l'-option is invoked, the first 'integer' characters, that precede the original mutant sequence (e.g., sites that function as DNA barcode or sequence tag), are removed.
-t, --trailseq=	integer	When the '-l'-option is invoked, the first 'integer' characters, that trail the original mutant sequence (e.g., sites that function as DNA barcode or sequence tag), are removed.
-p, --pool	none	When the '-p'-option is invoked, the DNA-sequences (characterizing the different mutants) are translated into amino acids. The data is then pooled based on their amino acid sequence, assuming that identical amino acid sequences, though differing in their DNA sequence (synonymous mutants), have the same growth rate (but different initial population sizes). Note that even if the '-p'-option is not invoked, data is pooled based on the sequence name (which can be any string and not only letters from the DNA alphabet).

<code>-g, -group=</code>	integer	<p>When the <code>'-g'</code>-option is invoked, the data is grouped into subsets of mutants each of minimal size <code>'integer'</code>. This results in more data sets with less mutants, such that the per-data set computation time is reduced, without affecting parameter estimates or the shape of the log-likelihood surface (compared to analysis of the full data set). The program also ensures, that mutants with identical mutant or protein ID (i.e., mutants that have an identical DNA- or amino acid sequence) remain in the same data sub-set as they are assumed to evolve under an identical growth rate (<code>r</code>). In general, the last line of the resulting input file serves as a summary line (indicated by a <code>'-1'</code> as mutant/protein ID and by <code>'XXX'</code> in the sequence column) that gives the number of all sequencing reads not in the current sub-set (but which are still needed to be accounted for). Estimates for the growth rates and the selection coefficients for the summary line are not calculated (but simply equated to those of the reference strain), since MCMC simulation estimates for this 'summary mutant' will later be discarded. The name of the output file (i.e., the MCMC input file) is composed of the standard output file identifier <code>'MCMCInput'</code>, the grouping identifier (a consecutive number of the sub data sets created) and the name of the input file.</p>
<code>-i, -initialize</code>	none	<p>When the <code>'-i'</code>-option is invoked, and additional input file is created that specifies the initial growth rates (<code>r</code>) and initial population sizes (<code>c</code>) for all mutants based on the log-linear regression. Note that this could potentially bias the MCMC algorithm, since the initial starting point of the Markov chain could be trapped in a local log-likelihood optimum. Often, however, the median of the growth rates and the initial population sizes from the MCMC-DFE simulations are close enough to the corresponding estimates from the log-linear regression such that starting at these values could shorten the burn-in period and, thus, reduce the run time. For mutants with identical DNA or amino acid sequence, the mean initial population size is calculated (from these mutants) and taken as the starting value for the MCMC simulation program. Given the estimated mean initial population size, the growth rate is estimated from the log-linear regression. The name of the output file is composed of the standard output file identifier <code>'MCMCInput'</code>, an optional grouping identifier (see <code>'-g'</code>-option), the name of the input file, and an initialization file identifier <code>'_inputRC'</code>.</p>

An illustration of the output file produced by the python program (i.e., the input file for the MCMC simulation program) is depicted in Figure 2.

Depending on the invoked options, the name of the MCMC simulation input file that is produced by the python program is given by the standard output file identifier `'MCMCInput'`, an optional grouping identifier (see `'-g'`-option) and the name of the raw data file.

protID	seq	aa	r	rCIL	rCIU	s	sCIL	sCIU	4.8	7.2	9.6	12	16.8	26.4	36x(4.8)	4x(7.2)	4x(9.6)	4x(12)	4x(16.8)	4x(26.4)	4x(36)
1	CAAAACGGTTGGTCTGCTAACATGAA	NGWSANME	1	0.9959568244	1.0040431756	0	0.0040431756	0.0040431756	26082	29923	51373	61182	49083	56396	111318	1	1	1	1	1	1
2	AACACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0226838373	1.0106965648	1.0346711097	0.0226838373	0.0106965648	0.0346711097	626	738	1515	1497	1417	1928	2512	1	1	1	1	1	0
3	TAATACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0001721595	1.0135174892	1.0206989296	0.0001721595	0.0135174892	0.0206989296	532	642	1196	1414	1151	1596	2210	1	1	1	1	1	0
4	AAGACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0332578812	1.0241684213	1.042347341	0.0332578812	0.0241684213	0.042347341	579	499	1116	1510	1322	2080	3444	0	1	1	1	1	1
5	AAAAACGGTTGGTCTGCTAACATGAA	NGWSANME	1.038444547	1.0340523822	1.0448665173	0.038444547	0.0340523822	0.0448665173	717	706	1403	1755	1599	2344	4431	0	1	1	1	1	1
6	CAACACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0311282584	1.0247358119	1.037494337	0.0311282584	0.0247358119	0.037494337	727	861	1721	1897	1752	2526	4040	1	1	1	1	1	0
7	ACCAACGGTTGGTCTGCTAACATGAA	NGWSANME	1.030305803	1.0188962326	1.0417553734	0.030305803	0.0188962326	0.0417553734	1358	1536	2899	3046	3315	4906	7384	1	1	1	1	1	1
8	ACGACGGTTGGTCTGCTAACATGAA	NGWSANME	1.038054621	1.0293154247	1.0427838172	0.038054621	0.0293154247	0.0427838172	862	969	1979	2277	2194	3188	4896	1	1	1	1	1	1
9	ACATACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0268965044	1.0174664707	1.0363352381	0.0268965044	0.0174664707	0.0363352381	880	1091	1957	2029	2112	3081	4727	1	1	1	1	1	1
10	SAGAACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0349584329	1.0214225183	1.0484963476	0.0349584329	0.0214225183	0.0484963476	441	443	887	1235	1075	1645	2594	1	1	1	1	1	1
11	SAGAACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0423237063	1.0293590295	1.0528638351	0.0423237063	0.0293590295	1.0528638351	418	451	907	1236	1082	1769	2924	1	1	1	1	1	1
12	SCGACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0314469336	1.0257708701	1.0371011172	0.0314469336	0.0257708701	0.0371011172	1285	1446	2717	3190	2893	4274	6744	1	1	1	1	1	1
13	SGGACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0387027387	1.034055574	1.0433499033	0.0387027387	0.034055574	0.0433499033	1538	2438	4479	5834	5114	7841	14455	0	1	1	1	1	1
14	SGGACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0452100346	1.0371389514	1.0533812179	0.0452100346	0.0371389514	1.0533812179	1210	1520	3143	3944	3525	5566	10297	1	1	1	1	1	1
15	SGTAAACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0417924631	1.032683319	1.0509016071	0.0417924631	0.032683319	1.0509016071	1334	1680	3334	4422	3916	5776	10622	1	1	1	1	1	1
16	SAGACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0382111179	1.0249112279	1.0515110079	0.0382111179	0.0249112279	0.0515110079	505	633	1194	1631	1355	1948	2546	1	1	1	1	1	0
17	6AGTACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0379829588	1.0279824283	1.0479824853	0.0379829588	0.0279824283	1.0479824853	367	504	907	1077	938	1493	1946	1	1	1	1	1	1
18	6TCAACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0327954085	1.0248414625	1.0407493544	0.0327954085	0.0248414625	1.0407493544	620	624	1059	1258	1289	1843	3050	0	1	1	1	1	1
19	6TCAACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0295706582	1.0114329184	1.0477084001	0.0295706582	0.0114329184	1.0477084001	923	979	1493	1770	1907	2960	4523	1	1	1	1	1	1
20	6TCAACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0267647167	1.0151518422	1.0383710811	0.0267647167	0.0151518422	1.0383710811	746	815	1421	1550	1492	2385	3707	1	1	1	1	1	1
21	6TCTAACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0275903062	1.0115202082	1.0436604043	0.0275903062	0.0115202082	1.0436604043	815	813	1418	1550	1741	2495	3859	1	1	1	1	1	1
22	7ATAACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0402679414	1.0286828378	1.0518530451	0.0402679414	0.0286828378	1.0518530451	287	314	619	872	808	1123	2036	1	1	1	1	1	1
23	7ATCAACGGTTGGTCTGCTAACATGAA	NGWSANME	1.04127329	1.0308020421	1.051780156	0.04127329	0.0308020421	1.051780156	545	586	1111	1260	1305	2166	3645	1	1	1	1	1	1
24	7ATTAACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0519667794	1.0390486163	1.0648849425	0.0519667794	0.0390486163	1.0648849425	360	412	898	1140	1034	1637	2620	1	1	1	1	1	0
25	8ATGAACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0381992172	1.0284265966	1.0479718378	0.0381992172	0.0284265966	1.0479718378	441	504	994	1286	1133	1630	2620	1	1	1	1	1	0
26	8ATGAACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0175184847	1.0127712281	1.022267414	0.0175184847	0.0127712281	1.022267414	1887	2896	5528	5844	5477	7000	11436	0	1	1	1	1	1
27	9CATAACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0204646084	1.0156041671	1.0253250497	0.0204646084	0.0156041671	1.0253250497	1423	1721	3295	3704	3254	4080	7373	1	1	1	1	1	1
28	10CCAACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0232659712	1.0174185675	1.029113375	0.0232659712	0.0174185675	1.029113375	2334	3381	6081	6693	6273	8407	13853	1	1	1	1	1	1
29	10CCAACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0257346803	1.0079868669	1.0434735036	0.0257346803	0.0079868669	1.0434735036	3548	5407	10407	9572	10159	14588	21039	1	1	1	1	1	1
30	10CCGACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0266821542	1.0221297293	1.0311945791	0.0266821542	0.0221297293	0.0311945791	2119	2453	4759	5527	5235	6761	11833	1	1	1	1	1	1
31	10CTAACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0253921974	1.0172026128	1.033581782	0.0253921974	0.0172026128	1.033581782	2435	3193	6144	6538	6404	8626	13741	1	1	1	1	1	1
32	11CTAACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0383530994	1.033589335	1.0391472652	0.0383530994	0.033589335	1.0391472652	1244	1508	2944	3465	3222	4516	8766	1	1	1	1	1	1
33	11CTAACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0369712551	1.0308964095	1.0430461007	0.0369712551	0.0308964095	1.0430461007	1837	2246	4632	5323	4912	7328	12988	1	1	1	1	1	1
34	11CTAACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0374775488	1.0323016836	1.0426534141	0.0374775488	0.0323016836	1.0426534141	1111	1248	2399	3041	2805	4136	7475	1	1	1	1	1	1
35	11CTAACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0375347989	1.0309096641	1.0436996336	0.0375347989	0.0309096641	1.0436996336	839	969	1905	2389	2224	2977	5973	1	1	1	1	1	1
36	11TTAACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0337755821	1.0242591192	1.0432920449	0.0337755821	0.0242591192	1.0432920449	551	578	1064	1460	1332	1819	3330	1	1	1	1	1	1
37	11TTAACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0426843064	1.038465358	1.0498832548	0.0426843064	0.038465358	1.0498832548	522	624	1212	1636	1423	2192	4065	1	1	1	1	1	1
38	12GAAACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0002616314	0.98403143	1.0064602486	0.0002616314	0.0064602486	1.0064602486	625	660	1181	1424	1159	1355	1889	1	1	1	1	1	1
39	12GAAACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0031157965	0.9952629826	1.0109686103	0.0031157965	0.0047370174	0.0109686103	501	522	965	1122	938	1085	1369	1	1	1	1	1	0
40	13GACACGGTTGGTCTGCTAACATGAA	NGWSANME	0.9511941181	0.9321370557	0.9714311804	0.0482158819	0.0678929443	0.0265889196	557	745	1115	985	796	664	509	1	1	1	1	1	0
41	13GATACGGTTGGTCTGCTAACATGAA	NGWSANME	0.9595611097	0.9264346852	0.9526953453	0.0604358903	0.075563148	0.0473094657	402	483	636	696	474	376	252	1	1	1	1	1	0
42	14GCAACGGTTGGTCTGCTAACATGAA	NGWSANME	1.033424934	1.0222683662	1.0445815019	0.033424934	0.0222683662	1.0445815019	582	651	1124	1367	1405	2149	3302	1	1	1	1	1	1
43	14GCCACGGTTGGTCTGCTAACATGAA	NGWSANME	1.0547949283	1.0487457244	1.0608441321	0.0547949283	0.0487457244	1.0608441321	1083	1474	2639	3364	3329	5378	7680	1	1	1	1	1	1

Figure 2 – Schematic illustration of the output produced by the python program and which serves as input for the MCMC simulation program. This input data has been created using the minimal raw data shown in Fig. 1, where the first five bases and the last base have been discarded (barcodes; '-l' and '-t' option). DNA sequences ('seq' column) have been translated to amino acids ('aa' column) and pooled ('-p' option), such that identical amino acid sequences have the same protein ID ('protID' column). Estimates of the growth rates 'r' ('r' column) and the selection coefficient 's' ('s' column) along with the 95%-confidence intervals are based on the log-linear regression (where 'rCIL' and 'rCIU' give the lower and upper boundary of the confidence interval for the growth rate r, respectively. Notation is analogous for the selection coefficient 's'). Please note, that while the MCMC simulation program assumes that mutants with identical sequence information (i.e., for sequences with identical 'protID') evolve at the same growth rate, log-linear estimates for 'r' and 's' are based on individual mutants. The columns '4.8', '7.2', '9.6', '12', '16.8', and '26.4' give the number of sequencing reads obtained from sampling at these time points for each mutant (row). If the '-o detect'-option is invoked, the matrix of sequencing reads is followed by an outlier matrix for the corresponding time points and mutants, where '0' indicate data points that were classified as outliers.

3 Usage

We provide executables for (Mac OS X), Windows and Linux. The C++ source code along with a system specific makefile are provided under a GNU General Public License as published by the Free Software Foundation. If you do not need to compile the program yourself you can skip the next subsection.

3.1 Compilation

Note that compilation requires that the (Gnu Scientific Library) (gsl-library) is installed on your system. Information on how to install the gsl-library can be found under <http://www.gnu.org/software/gsl/>. On Windows the easiest way to obtain the gsl-library is to install Cygwin (<http://www.cygwin.com/>) including the developers (all) packages. Alternatively, MinGW (<http://www.mingw.org/>) provides a "minimalist GNU for Windows" development environment. Under MinGW though, the gsl-library needs to be installed independently. A short instruction is given in the following paragraph.

Compilation of the program has successfully been tested on MacOSX (10.9.5) using 'clang' (version 6.0) and 'gcc' (version 4.9), on Ubuntu (14.04 LTS) using 'g++' (version 4.8.2), and on Windows (8.1) using 'g++' (version 4.8.1 under MinGW; version 4.9.2 under Cygwin).

Windows

Cygwin Please note that Cygwin is a large collection of GNU and open source tools, which provide functionality similar to a Linux distribution on Windows. Thus, when using Cygwin, commands and folder navigation are different to those from a native Windows Shell (e.g., Powershell). For example, under Cygwin you first need to type

```
cd /cygdrive/c
```

to navigate to the drive 'C:\'.

The easiest way to compile the program is by using the provided makefile. After navigating to the folder where the source code is stored, simply type

```
make -f Makefile_MCMC_Cygwin .
```


Please note that you might want to adjust the makefile, in particular, to change the name of the executable which is by default set to 'empiricIST_MCMC_Cygwin.exe'.

Alternatively, you can manually compile the program from source by navigating to the folder where the source code is stored and by typing

```
g++ empiricIST_MCMC.cpp -lgsl -lgslcblas -O3 -o 'NameOfProgram' ,
```

where 'NameOfProgram' should be replaced by the actual name of the compiled program (e.g., **empiricIST_MCMC**).

Note that Cygwin does not allow to link libraries statically.

MinGW Please note that MinGW provides a complete open source programming tool set for the development of native Windows applications (including both different compilers and a “minimal system” bourne shell command line interpreter system MSYS).

Compilation and installation of the gsl-library from source Before compilation and installation of the gsl-library please read and follow the installation instructions provided with the most recent version of gsl.

1. Download the latest version of gsl from <http://ftpmirror.gnu.org/gsl/>
2. Navigate to the place where the downloaded tar archive is stored and unpack it by typing

```
tar -zxvf gsl-x.xx.tar.gz ,
```

where 'x.xx' should be replaced by the version number.

3. Navigate to 'gsl-x.xx/' and carefully read the 'INSTALL' document and follow the instructions to configure, make and install the gsl-library:

```
./configure  
make  
make install
```

4. GSL binaries, headers and library files are installed automatically in the 'bin/', 'include/gsl/', and 'lib/' subdirectories (if not specified otherwise; in that case you would also need to adjust the linker and compiler flags in the makefile).

Compilation of the MCMC Open a MSYS shell and type

```
cd /c
```

to navigate to the drive 'C:\'.

The easiest way to compile the program is by using the provided makefile. After navigating to the folder where the source code is stored simply type

```
make -f Makefile_empiricIST_MCMC_WINDOWS .
```

Please note that you might want to adjust the makefile, in particular, to change the name of the executable which is by default set to '**empiricIST_MCMC.exe**'.

Alternatively, you can manually compile the program from source by navigating to the folder where the source code is stored and by typing

```
g++ empiricIST_MCMC.cpp -lgsl -lgslcblas -O3 -o 'NameOfProgram' ,
```

where 'NameOfProgram' should be replaced by the actual name of the compiled program (e.g., **empiricIST_MCMC**).

Linux

The easiest way to compile the program is by using the provided makefile. Use the Shell to navigate to the folder where the source code is stored and simply type

```
make -f empiricIST_MCMC_Linux .
```

Please note that you might want to adjust the makefile, in particular, to change the name of the executable which is by default set to '**empiricIST_MCMC**'.

Alternatively, you can manually compile the program from source by navigating to the folder where the source code is stored and by typing

```
g++ empiricIST_MCMC.cpp -lgsl -lgslcblas -O3 -o 'NameOfProgram' ,
```

where 'NameOfProgram' should be replaced by the actual name of the compiled program (e.g., **empiricIST_MCMC**).

Mac OS X

The easiest way to compile the program is by using the provided makefile. Use the Terminal to navigate to the folder where the source code is stored and simply type

```
make -f empiricIST_MCMC_MacOSX .
```

Please note that you might want to adjust the makefile, in particular, to change the name of the executable which is by default set to 'empiricIST_MCMC'.

Alternatively, you can manually compile the program from source by navigating to the folder where the source code is stored and by typing

```
g++ empiricIST_MCMC.cpp -lgsl -lgslcblas -O3 -o 'NameOfProgram' ,
```

where 'NameOfProgram' should be replaced by the actual name of the compiled program (e.g., `empiricIST_MCMC`).

3.2 Execution

Note that the *empiricIST_MCMC* program is a command line program which is run from a command-line interface (e.g., Shell, Terminal, Powershell), with arguments and parameters being passed over the command line. An overview and description of the parameters can be found in Table 2. Depending on your operating system there are different ways to call the program. These will be explained in the subsequent paragraphs.

Windows

Cygwin There are two ways to execute the program. First, the easiest way is to use the provided bash script 'empiricIST_MCMC_Cygwin.sh', where the parameters (e.g., 'PathToDataFile', skipCol, burnin) can be specified by the user. To begin, navigate to the folder where the executable and the bash script are stored and type

```
./empiricIST_MCMC_Cygwin.sh .
```

Second, you can manually execute the program by using Cygwin to navigate to the folder where the executable is stored and typing

```
./ 'NameOfProgram' 'PathToDatafile' outfileName skipCol  
outliersPresent burnin subsampling noSets set  
proposalDistCScale proposalDistRSD initialRC printLogLTS  
printESS printOutput seed ,
```

where the parameters (e.g., 'PathToDataFile', skipCol, burnin) need to be specified by the user.

Note that under Cygwin the 'PathToDataFile' needs to be preceded by '/cyg-drive/c/' (if the datafile is stored on the 'c' drive).

Powershell There are two ways to execute the program. First, the easiest way is to use the provided powershell script '*empiricIST_MCMC_WINDOWS.ps1*', where the parameters (e.g., 'PathToDataFile', skipCol, burnin) can be specified by the user. To begin, open Powershell and navigate to the folder where the executable and the powershell script are stored and type

```
& empiricIST_MCMC_WINDOWS.ps1 .
```

Second, you can manually execute the program by using Powershell to navigate to the folder where the executable is stored and typing

```
& 'NameOfProgram' 'PathToDatafile' outfileName skipCol  
outliersPresent burnin subsampling noSets set  
proposalDistCScale proposalDistRSD initialRC printLogLTS  
printESS printOutput seed ,
```

where the parameters (e.g., 'PathToDataFile', skipCol, burnin) need to be specified by the user.

Linux

There are two ways to execute the program. First, the easiest way is to use the provided bash script '*empiricIST_MCMC_Linux.sh*', where the parameters (e.g., 'PathToDataFile', skipCol, burnin) can be specified by the user. To begin, navigate to the folder where the executable and the bash script are store and type

```
./empiricIST_MCMC_Linux.sh .
```

Second, you can manually execute the program by using the shell to navigate to the folder where the executable is stored and typing

```
./ 'NameOfProgram' 'PathToDatafile' outfileName skipCol  
outliersPresent burnin subsampling noSets set  
proposalDistCScale proposalDistRSD initialRC printLogLTS  
printESS printOutput seed ,
```

where the parameters (e.g., 'PathToDataFile', skipCol, burnin) need to be specified by the user.

Mac OS X

There are two ways to execute the program. First, the easiest way is to use the provided bash script '*empiricIST_MCMC_MacOSX.sh*', where the parameters (e.g., 'PathToDataFile', skipCol, burnin) can be specified by the user. To begin, navigate to the folder where the executable and the bash script are store and type

```
./empiricIST_MCMC_MacOSX.sh .
```

Second, you can manually execute the program by using the terminal to navigate to the folder where the executable is stored and typing

```
./ 'NameOfProgram' 'PathToDatafile' outfileName skipCol
outliersPresent burnin subsampling noSets set
proposalDistCScale proposalDistRSD initialRC printLogLTS
printESS printOutput seed ,
```

where the parameters (e.g., 'PathToDataFile', skipCol, burnin) need to be specified by the user.

Table 2 – A summary of the parameters of the MCMC program.

Category Parameter	Accepted values	Description
Data		
datafile	string	Give the full path to the datafile (e.g., /users/me/PathToData/reads.csv).
outfileName	string	Name of the output file. The program will take the name of the input file and add the 'outfileName', the time of execution and the identifier (e.g., '_C'). This produces for example /PathTo-Data/_ 'datafile' _ 'outfileName' _ 'date&time' _C.txt.
skipCol	integer, ≥ 0	Number of columns to skip in data file before read numbers start.
outliersPresent	bool	If '0' there is no outlier matrix in the data file. If '1' there is an outlier matrix in the data file.
MCMC		
burnin	integer, ≥ 0	Number of accepted values that are discarded (burn-in period). During the burn-in period the parameters of the proposal distribution are adjusted.
subSampling	integer, ≥ 0	After the burn-in period only every 'subSampling' accepted value is recorded (i.e., written to file).
noSets	integer, ≥ 0	Number of output data sets that are recorded each of size 'setSize'. The total chain length is given by 'burnin'+noSets' \times 'setSize' \times 'subSampling'.
setSize	integer, ≥ 0	Number of recorded samples per set. The total chain length is given by 'burnin'+noSets' \times 'setSize' \times 'subSampling'. Output to file is written every 'setSize' accepted and recorded samples.
proposalDistCScale	double, > 0	Scale parameter of the proposal distribution of initial population sizes 'c' (Cauchy distribution) .

proposalDistRSD	double, $\$ > 0$	Standard deviation of the proposal distribution of growth rates 'r' (Normal distribution) .
inititalRC	string	An alternative way to initialize the growth rates 'r', the initial population sizes 'c' and to (optionally) set the parameters of the proposal distributions. This option allows, for example, to continue an MCMC run that has not been run long enough from the previous accepted sample. Note that for continuing an MCMC run the burn-in has to be set to 0. If 'initialRC' = -1, the default initialization is used (i.e., r=1 and initial population size = first observed read number for all mutants). When 'initialRC' $\neq -1$ the program will search for the initialization file 'datafile'_'initialRC'. In particular if the 'datafile' has been specified as '/users/me/PathToData/reads.csv' and the name of the initialization file is '_reads_initialRC.txt', the initialRC parameter passed to the program should be 'initialRC.txt'. Note that this file needs to be located in the same directory as the datafile and that the name initialization file should be of the form XYZData_'initialRC'.
<hr/>		
Output		
printLogLTS	bool	If '0', no time series of log-likelihoods is written. If '1', a time series of log-likelihoods is written.
printESS	bool	If '0', no ESS statistics are written. If '1', ESS statistics are written. The effective sample size (ESS) is calculated every 1000 accepted samples.
printOutput	bool	If '0', no additional output will be written to screen. If '1', additional output will be written to screen. Mainly for inspection purposes.
<hr/>		
Random numbers		
seed	integer	Sets the random number seed. If 'seed' ≤ 0 a random number seed is created automatically based on computer run time.

4 MCMC output

The program creates different kinds of output files that contain time-series data for different parameters, summary and diagnostic statistics for those parameters and for the entire MCMC run. Note that, even though all output files are plain 'txt'-files, they are all formatted as '.tsv'-files (tab separated) and can be displayed nicely with any spreadsheet application (e.g., Excel). Furthermore, all output files are structured such that they start with a list of the input/parameters for the MCMC run which are followed by the simulation results.

Table 3 – A summary of the output of the MCMC program.

File Parameter	Description
<hr/>	
.*_R	
sample	Consecutive number of samples. Sample '0' gives the initial values.
r.*	Sampled value for the growth rate 'r' for all mutants.
.*_C	
sample	Consecutive number of samples. Sample '0' gives the initial values.
c.*	Sampled value for the initial population size 'c' for all mutants.
.*_logLTS	
sample	Consecutive number of samples. Sample '0' gives the initial values.
logL	Log-likelihood for the current sampled values for the initial population size 'c' and the growth rate 'r'.
.*_R_quantiles	
protID	Protein ID as specified by the input file.
mutant	Consecutive number of mutant identifier 'r.*'.
i%	Values for the $i\%$ -quantile of all samples of the growth rate 'r' for each mutant, where $i = 0, 1, 2.5, 5, 25, 50, 75, 95, 97.5, 99, 100$.
.*_C_quantiles	
protID	Protein ID as specified by the input file.
mutant	Consecutive number of mutant identifier 'c.*'.
i%	Values for the $i\%$ -quantile of all samples of the initial population size 'c' for each mutant, where $i = 0, 1, 2.5, 5, 25, 50, 75, 95, 97.5, 99, 100$.
.*_logL_quantiles	
logL	Log-Likelihood identifier.
i%	Values for the $i\%$ -quantile of all samples of the log-likelihood, where $i = 0, 1, 2.5, 5, 25, 50, 75, 95, 97.5, 99, 100$.
.*_ess	
sample	Number of samples after which effective sample size (ESS) is calculated. Note that the ESS is calculated every 'setSize' samples.
minESS	Minimum effective sample size computed for any parameter of interest (i.e., growth rate 'r', initial population size 'c' and log-likelihood). Note that $ESS \leq sample$.
r.*	ESS for growth rate 'r.*' for all mutant'.
c.*	ESS for initial population size 'c.*' for all mutants.
logL	ESS for log-likelihood.
acceptRatio	Overall acceptance ratio. To ensure high efficiency of the MCMC, the width of the proposal distributions should be chosen such that the acceptance ratio is between 0.15 – 0.45. Performance is maximal with an acceptance ratio around 0.25. During the burn-in period the width of the proposal distributions is automatically tuned – based on the acceptance ratio – such that the acceptance ratio, when recording samples, has close to maximal efficiency. Thus, a sufficiently long burn-in period not only increases the chance that recorded samples are actually taken from the posterior distribution, but also that the width of the proposal distribution is set appropriately.

.*_Diag_R	
protID	Protein ID as specified by the input file.
mutant.*	Consecutive number of mutant identifier 'r.*'.
HD(*)	<p>Hellinger distance (HD) between sets of samples from two probability distributions. Note that HD is bounded by $0 \leq HD \leq 1$ and can be used to inspect the similarity between two distributions, where $HD = 0$ corresponds to no divergence and $HD = 1$ corresponds to no common support between the distributions. The HD can be used to diagnose the MCMC in terms of its burn-in and whether samples obtained at different points of time came (most likely) from the same (posterior) distribution. Note that one cannot determine if the MCMC chain has truly converged, but only if a chain is internally similar. Here, the HD is calculated for up to 10 equally sized sets of consecutive samples from the MCMC simulation. To obtain sufficient statistical power, the HD between two sets of samples is calculated only if each set consisted of at least 1000 samples. If the total number of samples is less than $10 \times 1000 = 10000$, the number of batches is chosen such that the total number of samples is divided into sets of samples of size 1000 each. If the total number of samples exceeds 10000, the number of samples per set is given by the total number of samples divided by 10 (i.e., the maximal number of batches). If the HD between sets of samples is less than 0.1 the distribution of posterior samples shows a high degree of similarity; if $0.1 \leq HD \leq 0.3$ the distribution of posterior samples are still quite similar, but may require closer inspection; if $0.3 \leq HD \leq 0.5$ sets of samples are vaguely similar and should be inspected more closely; a $HD > 0.5$ indicates strong dis-similarity between sets of samples and could be an indicator that all samples that were taken before might not be from the posterior distribution and should be discarded as burn-in. Note that the HD depends on the degree of autocorrelation between samples. Thus, a high HD might not necessarily indicate that samples were obtained from different sampling distributions, but poor mixing (i.e., a low ESS) for the parameter of interest. For details see Boone et al. (2014).</p>
mean	The mean of the posterior distribution for the parameter of interest.
SD	The standard deviation (SD) of the posterior distribution for the parameter of interest calculated with respect to the total number of samples.
median	The median of the posterior distribution for the parameter of interest.
2.5%	The 2.5% quantile of the posterior distribution for the parameter of interest.
97.5%	The 97.5% quantile of the posterior distribution for the parameter of interest.
ESS	The effective sample size for the parameter of interest.
minHD	The minimum HD calculated between consecutive batches. If there are not enough samples (more than 2000) to calculate the HD this field will read -1.
maxHD	The minimum HD calculated between consecutive batches. If there are not enough samples (more than 2000) to calculate the HD this field will read -1.
.*_Diag_C	
protID	Protein ID as specified by the input file.
mutant	Consecutive number of mutant identifier 'c.*'.

	<p>Hellinger distance (HD) between sets of samples from two probability distributions. Note that HD is bounded by $0 \leq HD \leq 1$ and can be used to inspect the similarity between two distributions, where $HD = 0$ corresponds to no divergence and $HD = 1$ corresponds to no common support between the distributions. The HD can be used to diagnose the MCMC in terms of its burn-in and whether samples obtained at different points of time came (most likely) from the same (posterior) distribution. Note that one cannot determine if the MCMC chain has truly converged, but only if a chain is internally similar. Here, the HD is calculated for up to 10 equally sized sets of consecutive samples from the MCMC simulation. To obtain sufficient statistical power, the HD between two sets of samples is calculated only if each set consisted of at least 1000 samples. If the total number of samples is less than $10 \times 1000 = 10000$, the number of batches is chosen such that the total number of samples is divided into sets of samples of size 1000 each. If the total number of samples exceeds 10000, the number of samples per set is given by the total number of samples divided by 10 (i.e., the maximal number of batches). If the HD between sets of samples is less than 0.1 the distribution of posterior samples shows a high degree of similarity; if $0.1 \leq HD \leq 0.3$ the distribution of posterior samples are still quite similar, but may require closer inspection; if $0.3 \leq HD \leq 0.5$ sets of samples are vaguely similar and should be inspected more closely; a $HD > 0.5$ indicates strong dis-similarity between sets of samples and could be an indicator that all samples that were taken before might not be from the posterior distribution and should be discarded as burn-in. Note that the HD depends on the degree of autocorrelation between samples. Thus, a high HD might not necessarily indicate that samples were obtained from different sampling distributions, but poor mixing (i.e., a low ESS) for the parameter of interest. For details see Boone et al. (2014).</p>
HD(*)	
mean	The mean of the posterior distribution for the parameter of interest.
SD	The standard deviation (SD) of the posterior distribution for the parameter of interest calculated with respect to the total number of samples.
median	The median of the posterior distribution for the parameter of interest.
2.5%	The 2.5% quantile of the posterior distribution for the parameter of interest.
97.5%	The 97.5% quantile of the posterior distribution for the parameter of interest.
ESS	The effective sample size for the parameter of interest.
minHD	The minimum HD calculated between consecutive batches. If there are not enough samples (more than 2000) to calculate the HD this field will read -1 .
maxHD	The minimum HD calculated between consecutive batches. If there are not enough samples (more than 2000) to calculate the HD this field will read -1 .
<hr/>	
.*_Diag_logL	
logL	Log-likelihood tag.
mutant	Consecutive number of mutant identifier 'c.*'.

Hellinger distance (HD) between sets of samples from two probability distributions. Note that HD is bounded by $0 \leq HD \leq 1$ and can be used to inspect the similarity between two distributions, where $HD = 0$ corresponds to no divergence and $HD = 1$ corresponds to no common support between the distributions. The HD can be used to diagnose the MCMC in terms of its burn-in and whether samples obtained at different points of time came (most likely) from the same (posterior) distribution. Note that one cannot determine if the MCMC chain has truly converged, but only if a chain is internally similar. Here, the HD is calculated for up to 10 equally sized sets of consecutive samples from the MCMC simulation. To obtain sufficient statistical power, the HD between two sets of samples is calculated only if each set consisted of at least 1000 samples. If the total number of samples is less than $10 \times 1000 = 10000$, the number of batches is chosen such that the total number of samples is divided into sets of samples of size 1000 each. If the total number of samples exceeds 10000, the number of samples per set is given by the total number of samples divided by 10 (i.e., the maximal number of batches). If the HD between sets of samples is less than 0.1 the distribution of posterior samples shows a high degree of similarity; if $0.1 \leq HD \leq 0.3$ the distribution of posterior samples are still quite similar, but may require closer inspection; if $0.3 \leq HD \leq 0.5$ sets of samples are vaguely similar and should be inspected more closely; a $HD > 0.5$ indicates strong dis-similarity between sets of samples and could be an indicator that all samples that were taken before might not be from the posterior distribution and should be discarded as burn-in. Note that the HD depends on the degree of autocorrelation between samples. Thus, a high HD might not necessarily indicate that samples were obtained from different sampling distributions, but poor mixing (i.e., a low ESS) for the parameter of interest. For details see Boone et al. (2014).

HD(*)

mean The mean of the posterior distribution for the parameter of interest.

SD The standard deviation (SD) of the posterior distribution for the parameter of interest calculated with respect to the total number of samples.

median The median of the posterior distribution for the parameter of interest.

2.5% The 2.5% quantile of the posterior distribution for the parameter of interest.

97.5% The 97.5% quantile of the posterior distribution for the parameter of interest.

ESS The effective sample size for the parameter of interest.

minHD The minimum HD calculated between consecutive batches. If there are not enough samples (more than 2000) to calculate the HD this field will read -1.

maxHD The minimum HD calculated between consecutive batches. If there are not enough samples (more than 2000) to calculate the HD this field will read -1.

.*_Diag_summary

samples Absolute number of accepted samples taken during the MCMC run.

minESS(c) The minimum effective sample size (ESS) that was observed for any initial population size 'c.*'.

maxACT(c) The maximal auto-correlation time (ACT) that was observed for any initial population size 'c.*'.

minESS(r) The minimum effective sample size (ESS) that was observed for any growth rate 'r.*'.

maxACT(r) The maximal auto-correlation time (ACT) that was observed for any growth rate 'r.*'.

minESS(logL) The minimum effective sample size (ESS) that was observed for the log-likelihood.

maxACT(r) The maximal auto-correlation time (ACT) that was observed for the log-likelihood.

minESS(all) The minimum effective sample size (ESS) that was observed for all parameters.

maxACT(all) The maximal auto-correlation time (ACT) that was observed for all parameters.

acceptRatio	The overall acceptance ratio of accepted (and recorded) samples. To ensure high efficiency of the MCMC, the width of the proposal distributions should be chosen such that the acceptance ratio is between 0.15 – 0.45. Performance is maximal with an acceptance ratio around 0.25. During the burn-in period the width of the proposal distributions is automatically tuned – based on the acceptance ratio – such that the acceptance ratio, when recording samples, has close to maximal efficiency. Thus, a sufficiently long burn-in period not only increases the chance that recorded samples are actually taken from the posterior distribution, but also that the width of the proposal distribution is set appropriately.
jumpSDR	Standard deviation of the proposal distribution of growth rates 'r' (Normal distribution) after auto-tuning.
jumpSDC	Scale parameter of the proposal distribution of initial population sizes 'c' (Cauchy distribution) after auto-tuning.
.*_initialRC	This file prints the last accepted values of the MCMC run so that these could be used as initial values, e.g., to continue an MCMC run that has not yielded enough independent samples. The first line gives the last sampled growth rates 'r.*', the second line gives the last sampled initial population sizes 'c.*', and the third line gives the standard deviation of the proposal distribution of growth rates 'r' (Normal distribution) and the scale parameter of the proposal distribution of initial population sizes 'c' (Cauchy distribution) after auto-tuning.

4.1 Combining Files

In case data has been split into multiple subsets to enhance computational performance (see '-g' option in Tab. 1), we provide scripts to assemble the individual MCMC output files for each sub-data set to a single file which contains all information necessary for further analysis. Note that we only provide scripts to assemble 'diagnostic', 'quantiles' and 'posterior sample' files (i.e., files which contain the MCMC samples for the growth rate (r) and the initial population size (c)). The `Combine_All.sh` script is a wrapper which executes all 'Combine' scripts (i.e., `Combine_Diagnostic_C.sh`, `Combine_Diagnostic_R.sh`, `Combine_Quantiles_C.sh`, `Combine_Quantiles_R.sh`, `Combine_PopSizes_C.sh` and `Combine_GrowthRates_R.sh`). It furthermore renames the `empiricIST_MCMC` output files (by removing the time stamp) and deletes the files that will not be combined, if specified by the user. More information on all the script files is given in Table 4.

Note that the bash scripts need to have the appropriate rights to perform the operations. In order to set the correct permissions, open a Shell and navigate to the folder where the scripts are stored and type

```
chmod 755 Combine_*.sh .
```

Table 4 – A summary of the 'combine files scripts'.

Script	No. of passed variables	Description
--------	-------------------------	-------------

Combine_All.sh	6	<p>This script renames all MCMC output files, executes all individual combine scripts, and deletes the files that will not be combined. To execute the script open a Shell and navigate to the folder where the scripts are stored and type</p> <pre>./Combine_All.sh pathToPerlRename pathToData prefixFileName suffixFileName maxIndex deleteFiles .</pre> <p>pathToPerlRename: Provide full path to 'rename.pl'.</p> <p>pathToData: Provide full path to data folder.</p> <p>prefixFileName: Provide file name prefix. For an example see Figure 3.</p> <p>suffixFileName: Provide file name suffix. For an example see Figure 3.</p> <p>maxIndex: The number of sub-data sets the original data has been split-up to. For an example see Figure 3.</p> <p>deleteFiles: When set to '1' files that will not be combined will be deleted.</p> <p>Information on the individual combine scripts can be found below.</p>
Combine_Diagnostic_C.sh	4	<p>This script combines all MCMC output files of type 'Diagnostic_C' of all sub-data sets to a single file. Note that the reference sequence and the 'summary mutants' (i.e., protID=-1) will be deleted, since they do not contain any relevant information. To execute the script open a Shell and navigate to the folder where the scripts are stored and type</p> <pre>./Combine_Diagnostic_C.sh pathToData prefixFileName suffixFileName maxIndex .</pre> <p>pathToData: Provide full path to data folder.</p> <p>prefixFileName: Provide file name prefix. For an example see Figure 3.</p> <p>suffixFileName: Provide file name suffix. For an example see Figure 3.</p> <p>maxIndex: The number of sub-data sets the original data has been split-up to. For an example see Figure 3.</p>
Combine_Diagnostic_R.sh	4	<p>As above, but for growth rates 'R'. Note that if estimates with identical estimates are only represented once. That is the case, for instance, when data has been analyzed on the level of amino acids (see '-p' option in Tab. 1).</p>

Combine_Quantiles_C.sh	4	<p>This script combines all MCMC output files of type 'Quantiles_C' of all sub-data sets to a single file. Note that the reference sequence and the 'summary mutants' (i.e., protID=-1) will be deleted, since they do not contain any relevant information. To execute the script open a Shell and navigate to the folder where the scripts are stored and type</p> <pre>./Combine_Quantiles_C.sh pathToData prefixFileName suffixFileName maxIndex .</pre> <p>pathToData: Provide full path to data folder. prefixFileName: Provide file name prefix. For an example see Figure 3. suffixFileName: Provide file name suffix. For an example see Figure 3. maxIndex: The number of sub-data sets the original data has been split-up to. For an example see Figure 3.</p>
Combine_Quantiles_R.sh	4	<p>As above, but for growth rates 'R'. Note that if estimates with identical estimates are only represented once. That is the case, for instance, when data has been analyzed on the level of amino acids (see '-p' option in Tab. 1).</p>
Combine_PopSizes_C.sh	4	<p>This script combines all MCMC output files of type '_C' of all sub-data sets to a single file. Note that the reference sequence and the 'summary mutants' (i.e., protID=-1) will be deleted, since they do not contain any relevant information. To execute the script open a Shell and navigate to the folder where the scripts are stored and type</p> <pre>./Combine_PopSizes_C.sh pathToData prefixFileName suffixFileName maxIndex .</pre> <p>pathToData: Provide full path to data folder. prefixFileName: Provide file name prefix. For an example see Figure 3. suffixFileName: Provide file name suffix. For an example see Figure 3. maxIndex: The number of sub-data sets the original data has been split-up to. For an example see Figure 3.</p>
Combine_GrowthRates_R.sh	4	<p>As above, but for growth rates 'R'. Note that if estimates with identical estimates are only represented once. That is the case, for instance, when data has been analyzed on the level of amino acids (see '-p' option in Tab. 1).</p>
RenameMCMCOutput.sh	3	<p>This script renames the <code>empiricIST_MCMC</code> output files (by removing the time stamp) To execute the script open a Shell and navigate to the folder where the scripts are stored and type</p> <pre>./RenameMCMCOutput.sh pathToPerlRename pathToData prefixFileName .</pre> <p>pathToPerlRename: Provide full path to 'rename.pl'. pathToData: Provide full path to data folder. prefixFileName: Provide file name prefix. For an example see Figure 3.</p>

$$\underbrace{\text{MCMCInput}}_{\text{prefixName}} - \overbrace{99}^{\text{maxIndex}} - \underbrace{\text{HGW_VIE}}_{\text{suffixName}} \overbrace{\text{08_06_2015_142000}}^{\text{date stamp}} \underbrace{\text{R}}_{\text{type}} . \text{txt}$$

Figure 3 – Illustration of how to set the command line arguments for the combine scripts.

Additionally, there is a shell script 'FormatTracer.sh' that formats the posterior sample output file (e.g., containing the initial population size 'c' or the growth rate 'r') such that it can be read and analyzed by *Tracer* (?). *Tracer* is a graphical tool for visualization and diagnostics of MCMC output that, for instance, displays the posterior distribution and its credibility interval, calculates the effective sample size (ESS; note that values might be slightly different from those calculated by the `empiricIST_MCMC` program since we use a more accurate but computational more intensive algorithm), and shows the trace of the posterior samples. Note that the input file is expected to be formatted as the output file created from the `Combine_PopSizes_C.sh/Combine_GrowthRates_R.sh` script. Otherwise the provided `Create_TailShapeFileR.sh` script can be used to obtain a correctly formatted input file (please see below for details).

`./FormatTracer.sh pathToData fileName .`

pathToData: Provide full path to data folder.

fileName: Provide file name without file extension (e.g., '.txt').

5 DFE tail shape estimation

The growth rate posterior samples obtained from the `empiricIST_MCMC` program can be used to estimate the shape of the beneficial tail of the DFE. As part of the *empiricIST* software package, we provide a python script – `empiricIST_MCMC_TailShape.py` that fits a generalized pareto distribution to the observed beneficial mutations by maximizing the log-likelihood with respect to the shape and scale parameter κ and ψ , respectively. Based on the shape parameter κ one can discriminate between three different domains of attraction – the Weibull, Gumbel and Fréchet domain – each corresponding to a different extreme value distribution. In biological terms, these different domains quantify the level of adaptedness of the organism in its (experimental) environment. In particular, the Gumbel domain ($\kappa = 0$; null model corresponding to 'normal' level of adaptedness) is characterized by an exponential tail, whereas the Weibull domain ($\kappa < 0$; better adapted) has lighter than exponential tails, and the Fréchet domain ($\kappa > 0$; less well adapted) has heavier than exponential tails. For more details please consult ?.

While the python script primarily estimates the DFE tail shape parameter κ , there are some additional options that will be detailed here.

The general usage is as follows: After opening a command-line interface (e.g., Shell, Terminal) and navigating to the location of the `empiricIST_MCMC_TailShape.py` file, the program can be executed by typing

```
python empiricIST_MCMC_TailShape.py [options] .
```

Without specifying any options the program will exit with an error and provide a short documentation on its usage, as it requires the name of the data input file (by invoking the '-f' option). Note that the input file needs to be formatted in a specific way. If the data set had been split into multiple subsets (see '-g' option in Tab. 1), and has been re-assembled by the using the provided shell scripts (see 'Combining Files'), the input file is already correctly formatted and there is nothing that needs to be done. However, if the data set has been analyzed as a whole (i.e., without being split into multiple subsets) the file containing the posterior growth rate samples needs to be reformatted to match the required input format. This can be done by using the provided shell script `Create_TailShapeFileR.sh`, which creates the input file for the `empiricIST_MCMC_TailShape.py` program and can be executed by opening a Shell and navigating to the folder where the scripts are stored and typing

```
./Create_TailShapeFileR.sh pathToData fileName .
```

pathToData: Provide full path to data folder.

fileName: Provide file name without file extension (e.g., '.txt').

Note that to execute the bash script, it needs to have the appropriate rights to perform the operations. In order to set the correct permissions, open a Shell and navigate to the folder where the script is stored and type

```
chmod 755 Create_TailShapeFile_R.sh .
```

All options and their usage of the `empiricIST_MCMC_TailShape.py` program are given in Table 5.

Table 5 – A summary of the options of the `empiricIST_MCMC_TailShape.py` program.

Short/Long option	Accepted values	Description
-h, -help	none	When the '-h'-option is invoked, a short documentation on the usage of the program is shown. Note that, if this option is invoked, the python program is not executed.

-f, -file=	string	The '-f'-option is a mandatory option, which passes the name of the data input file (tsv formatted) to the python program. Files created by the python program will take the name of the input file and add option-dependent specific file identifiers. Note that even if a random data set is created (see 'r'-option), a file name must be provided since it serves as the prefix for the output file name.
-m, -missing	none	When the '-m'-option is invoked, the distribution of measured fitnesses is shifted relative to the smallest observed selection coefficient to account for missing data (i.e., selection coefficients too small to have been observed, though this might not be a problem with EMPIRIC data).
-s, -samples=	integer	When the '-s'-option is invoked, the python program will only consider samples with more than 'samples' beneficial mutations for maximum likelihood estimation. Note that when this option is not specified, the default is set to 10.
-l, -lrt	none	When the '-l'-option is invoked, a likelihood-ratio test with null hypotheses $\kappa = 0$ against median($\hat{\kappa}$) is performed. Note that the distribution of the test statistic is generated by using a parametric bootstrap approach with 10.000 samples (? , see).
-r, -random=	float	When the '-r'-option is invoked, 100 random data sets are created with 100 samples each drawn from a generalized pareto distribution with scale parameter $\psi = 1$ and κ passed as command line argument.

By default the program will always create two output files that contain the κ and ψ estimates called `<InputFileName>_TailShape_KappaShape.txt` and `<InputFileName>_TailShape_PsiShape.txt`, respectively. Note that even if a random data set is created (see 'r'-option), a file name must be provided since it serves as the prefix for the output file name. When performing a likelihood-ratio test (by invoking the 'l'-option) an additional file called `<InputFileName>_LRT.txt` is created containing the $\hat{\kappa}$ against which $H_0 : \kappa = 0$ is evaluated. Furthermore, the output file gives the maximum-likelihood estimate $\hat{\psi}_{\kappa_0}$, i.e., the estimated scale parameter ψ restricting $\kappa = 0$, the value of the test statistic

$$-2 \log(\Lambda) = 2(\mathcal{L}(\mathbf{X} \mid \hat{\kappa}, \hat{\psi}) - \mathcal{L}(\mathbf{X} \mid o, \hat{\psi}_{\kappa_0})), \quad (1)$$

where \mathbf{X} denotes a single vector of posterior samples of growth rates and \mathcal{L} is the log-likelihood function, and the associated p-value along with the sample size. Note that power critically depends on sample size as has been discussed in ? and ? (the latter rather in the context of estimating κ accurately).

References

- Bank, C., R. T. Hietpas, A. Wong, D. N. Bolon, and J. D. Jensen, 2014. A bayesian mcmc approach to assess the complete distribution of fitness effects of new mutations: Uncovering the potential for adaptive walks in challenging environments. *Genetics* 196:841–852.
- Boone, E. L., J. R. Merrick, and M. J. Krachey, 2014. A hellinger distance approach to mcmc diagnostics. *Journal of Statistical Computation and Simulation* 84:833–849.
- Matuszewski, S., J. D. Jensen, and C. Bank, in prep. empiricist: An integrated software and analysis tool for analyzing time-sampled sequence data such as empiric .