

User manual for Pyvolve v1.0

Stephanie J. Spielman

Email: stephanie.spielman@gmail.com

Contents

1	Introduction	2
2	Basic Usage	2
3	Defining phylogenies	3
4	Defining Evolutionary Models	4
4.1	Nucleotide Models	5
4.2	Amino-acid models	6
4.3	Mechanistic (dN/dS) codon models	7
4.4	Mutation-selection models	8
4.5	Empirical codon model	9
4.6	Specifying equilibrium frequencies	10
4.6.1	EqualFrequencies class	11
4.6.2	RandomFrequencies class	11
4.6.3	CustomFrequencies class	11
4.6.4	ReadFrequencies class	12
4.6.5	EmpiricalModelFrequencies class	12
4.6.6	Restricting frequencies to certain states	13
4.6.7	Converting frequencies between alphabets	13
4.7	Specifying mutation rates	14
5	Defining Partitions	14
6	Evolving sequences	15
6.1	Evolver output files	15
6.2	Sequence post-processing	16
6.2.1	Interpreting the "site_rates.txt" output file	16
6.2.2	Interpreting the "site_rates_info.txt" output file	17
6.3	Simulating replicates	17
7	Implementing site-wise rate heterogeneity	17
7.1	Implementing site-wise heterogeneity for nucleotide and amino-acid models	18
7.1.1	Gamma-distributed rate categories	18
7.1.2	Custom-distributed rate categories	18
7.2	Implementing site-wise heterogeneity for mechanistic codon models	19
7.3	Implementing site-wise heterogeneity for mutation-selection models	19
7.4	Implementing site-wise heterogeneity for the Empirical Codon Model	20
8	Implementing branch (temporal) heterogeneity	20
9	Implementing branch-site heterogeneity	22

10 Unique Pyvolve Features	22
10.1 Specifying custom rate matrices	22
10.2 Matrix scaling options	23
10.3 Perturbing branch lengths	24

1 Introduction

Pyvolve (pronounced “pie-volve”) is an open-source Python module for simulating genetic data along a phylogeny using Markov models of sequence evolution, according to standard methods [20]. The module is available for download under a FreeBSD license from <http://sjspielman.org/pyvolve/>. Pyvolve has several dependencies, including `BioPython`, `NumPy`, and `SciPy`. These modules must be properly installed and in your Python path. Please file any and all bug reports on the github repository [Issues](#) section.

Pyvolve is written such that it can be seamlessly integrated into your Python pipelines without having to interface with external software platforms. However, please note that for extremely large (>1000 taxa) and/or extremely heterogenous simulations (e.g. where each site evolves according to a unique evolutionary model), Pyvolve may be quite slow and thus may take several minutes to run. Faster sequence simulators you may find useful are detailed in ref. [1], which gives an overview of various sequence simulation softwares (from 2012).

Pyvolve supports a variety of evolutionary models, including the following:

- Nucleotide Models
 - Generalized time-reversible model [16] and all nested variants
- Amino-acid exchangeability models
 - JTT [6], WAG [18], and LG [10]
- Codon models
 - Mechanistic (dN/dS) models (MG-style [11] and GY-style [3])
 - Empirical codon model [9]
- Mutation-selection models
 - Halpern-Bruno model [4], implemented for codons and nucleotides

Both site- and branch- (temporal) heterogeneity are supported. A detailed and highly-recommended overview of Markov process evolutionary models, for DNA, amino acids, and codons, is available in the book *Computational Molecular Evolution*, by Ziheng Yang [20].

Although Pyvolve does not simulate insertions and deletions (indels), Pyvolve does include several novel options not available (to my knowledge) in other sequence simulation softwares. These options, detailed in Section 10, include custom rate-matrix specification, novel matrix-scaling approaches, and branch length perturbations.

2 Basic Usage

Similar to other simulation platforms, Pyvolve evolves sequences in groups of **partitions** (see, for instance, the Indelible simulation platform [2]). Each partition has an associated size and model (or set of models, if branch heterogeneity is desired). Note that all partitions will evolve according to the same phylogeny¹.

The general framework for a simple simulation is given below. In order to simulate sequences, you must define the phylogeny along which sequences evolve as well as any evolutionary model(s) you’d like to use, and assign model(s) to partition(s). Each evolutionary model has associated parameters which you can customize, as detailed in Section 4.

¹If you wish to have different partitions evolve according to distinct phylogenies, I recommend performing several simulations and then merging the resulting alignments in the post-processing stage.

```

1 ##### General pyvolve framework #####
2 #####
3
4 # Import the Pyvolve module
5 import pyvolve
6
7 # Read in phylogeny along which Pyvolve should simulate
8 my_tree = pyvolve.read_tree(file = "file_with_tree_for_simulating.tre")
9
10 # Define evolutionary model(s) with the Model class
11 my_model = pyvolve.Model(<model_type>, <custom_model_parameters>)
12
13 # Define partition(s) with the Partition class
14 my_partition = pyvolve.Partition(models = my_model, size = 100)
15
16 # Evolve partitions with the callable Evolver class
17 my_evolver = pyvolve.Evolver(tree = my_tree, partitions = my_partition)
18 my_evolver() # evolve sequences

```

Each of these steps is explained below, in detail with several examples. For additional information, consult the API documentation, at <http://sjspielman.org/pyvolve>. Further, all functions and classes in Pyvolve have highly descriptive docstrings, which can be accessed with Python's `help()` function.

3 Defining phylogenies

Phylogenies must be specified as newick strings (see [this wikipedia page](#) for details) *with branch lengths*. Pyvolve reads phylogenies using the function `read_tree`, either from a provided file name or directly from a string:

```

1 # Read phylogeny from file with the keyword argument "file"
2 phylogeny = pyvolve.read_tree(file = "/path/to/tree/file.tre")
3
4 # Read phylogeny from string with the keyword argument "tree"
5 phylogeny = pyvolve.read_tree(tree = "(t4:0.785,(t3:0.380,(t2:0.806,(t5:0.612,t1:0.660
   ):0.762):0.921):0.207);")

```

To implement branch (temporal) heterogeneity, in which different branches on the phylogeny evolve according to different models, you will need to specify *model flags* at particular nodes in the newick tree, as detailed in Section 8.

Further, to assess that a phylogeny has been parsed properly (or to determine the automatically-assigned names of internal nodes), use the `print_tree` function:

```

1 # Read phylogeny from string
2 phylogeny = pyvolve.read_tree(tree = "(t4:0.785,(t3:0.380,(t2:0.806,(t5:0.612,t1:0.660
   ):0.762):0.921):0.207);")
3
4 # Print the parsed phylogeny
5 pyvolve.print_tree(phylogeny)

```

```

6  ## Output from the above statement: node_name  branch_length  model_flag
7  '''
8  >>>      t4 0.785 None
9  >>>      internal_node3 0.207 None
10 >>>      t3 0.38 None
11 >>>      internal_node2 0.921 None
12 >>>      t2 0.806 None
13 >>>      internal_node1 0.762 None
14 >>>      t5 0.612 None
15 >>>      t1 0.66 None
16 '''

```

In the above output, tabs represent nested hierarchies in the phylogeny. Each line shows the node name (either a tip name, "root", or an internal node), the branch length leading to that node, and the model flag associated with that node. This final value will be `None` if model flags are not provided in the phylogeny. Again, note that model flags are only required in cases of branch heterogeneity (see Section 8).

4 Defining Evolutionary Models

The evolutionary models built into Pyvolve are outlined in Table 1 of this manual. Pyvolve uses `Model` objects to store evolutionary models:

```

1  # Basic framework for defining a Model object (second argument optional)
2  my_model = Model(<model_type>, <custom_model_parameters_dictionary>)

```

A single argument, `<model_type>`, is required when defining a `Model` object. Available model types are shown in Table 1. Each model type has various associated parameters, which can be customized via the second optional argument to `Model`, written above as `<custom_model_parameters_dictionary>`. This argument should be a dictionary of parameters to customize, and each modeling framework has particular keys which can be included in this dictionary. Available model types and associated customizable parameters are shown in Table 1 and detailed in the subsections below.

Note that there are several additional optional keyword arguments which may be passed to `Model`, including arguments pertaining to site-rate heterogeneity (see Section 7) and the matrix scaling algorithm (see Section 10.2).

Modeling framework	Pyvolve model type(s)	Optional parameters ("key")
Nucleotide models	"nucleotide"	<ul style="list-style-type: none"> Equilibrium frequencies ("state_freqs") Mutation rates ("mu" or "kappa")
Empirical amino-acid models	"JTT", "WAG", or "LG"	<ul style="list-style-type: none"> Equilibrium frequencies ("state_freqs")
Mechanistic (dN/dS) codon models	"GY", "MG", or "codon"	<ul style="list-style-type: none"> Equilibrium frequencies ("state_freqs") Mutation rates ("mu" or "kappa") dN/dS ("alpha", "beta", and/or "omega")
Mutation-selection models	"MutSel"	<ul style="list-style-type: none"> Equilibrium frequencies ("state_freqs") OR fitness values ("fitness") Mutation rates ("mu" or "kappa")
Empirical codon model (ECM)	"ECMrest", "ECMunrest", or "ECM"	<ul style="list-style-type: none"> Equilibrium frequencies ("state_freqs") Transition-tranversion bias(es) ("k_ti" and/or "k_tv") dN/dS^\dagger ("alpha", "beta", and/or "omega")

Table 1. Accepted model types in Pyvolve with associated customizable parameters. Names given in the column "Pyvolve model type(s)" should be specified as the first argument to `Model` as strings (case-insensitive). Customizable parameters indicated in the column "Optional parameters" should be specified as keys in the custom model-parameters dictionary, the second argument using when defining a `Model` object.

[†]Note that the interpretation of this dN/dS value is different from the usual interpretation.

Subsections below explain each modeling framework in detail, with examples of parameter customizations.

4.1 Nucleotide Models

Nucleotide rate matrix elements, for the substitution from nucleotide i to j , are generally given by

$$q_{ij} = \mu_{ij}\pi_j \quad (1)$$

where μ_{ij} describes the rate of change from nucleotide i to j (i.e. mutation rate), and π_j represents the equilibrium frequency of the target nucleotide j . Note that mutation rates are symmetric, e.g. $\mu_{ij} = \mu_{ji}$.

By default, nucleotide models have equal equilibrium frequencies and equal mutation rates. A basic model can be constructed with,

```
1 # Simple nucleotide model
2 nuc_model = pyvolve.Model("nucleotide")
```

To customize a nucleotide model, provide a custom-parameters dictionary with optional keys "state_freqs" for custom equilibrium frequencies and "mu" for custom mutation rates (see Section 4.6 for details on frequency customization and Section 4.7 for details on mutation rate customization).

```
1 # Define mutation rates in a dictionary with keys giving the nucleotide pair
2 # Below, the rate from A to C is 0.5, and similarly C to A is 0.5
3 custom_mu = {"AC":0.5, "AG":0.25, "AT":1.23, "CG":0.55, "CT":1.22, "GT":0.47}
4
5 # Define custom frequencies, in order A C G T. This can be a list or numpy array.
6 freqs = [0.1, 0.45, 0.3, 0.15]
```

```

7
8 # Construct nucleotide model with custom mutation rates and frequencies.
9 nuc_model = pyvolve.Model( "nucleotide", {"mu":custom_mu, "state_freqs":freqs} )

```

As nucleotide model mutation rates are symmetric, if you provide a rate for $A \rightarrow T$ (key "AT"), it will automatically be applied as the rate for $T \rightarrow A$. Any unspecified mutation rate pairs will have a value of 1.

As an alternate to "mu", you can provide the key "kappa", which corresponds to the transition:transversion ratio (e.g. for an HKY85 model [5]), in the custom-parameters dictionary. When kappa is specified, transversion mutation rates are set to 1, and transition mutation rates are set to the provided "kappa" value.

```

1 # Construct nucleotide model with transition-to-transversion bias, and default
  frequencies
2 nuc_model = pyvolve.Model( "nucleotide", {"kappa":2.75, "state_freqs":freqs} )

```

4.2 Amino-acid models

Amino-acid exchangeability matrix elements, for the substitution from amino acid i to j , are generally given by

$$q_{ij} = r_{ij}\pi_j \quad (2)$$

where r_{ij} is a symmetric matrix that describes the probability of changing from amino acid i to j , and π_j is the equilibrium frequency of the target amino acid j . The r_{ij} matrix corresponds to an empirically determined model, such as WAG [18] or LG [10].

By default, Pyvolve assigns equal equilibrium frequencies. A basic amino-acid model can be constructed with,

```

1 # Simple amino-acid model
2 aa_model = pyvolve.Model("WAG") # Here, WAG can be one of JTT, WAG, LG (case-
  insensitive)

```

To customize an amino-acid model, provide a custom-parameters dictionary with the key "state_freqs" for custom equilibrium frequencies (see Section 4.6 for details on frequency customization). Note that amino-acid frequencies must be in the order A, C, D, E, ... Y. Further, to specify the model's default equilibrium frequencies, use Pyvolve's `EmpiricalModelFrequencies` class (described in-depth in Section 4.6):

```

1 # Define default WAG state frequencies
2 f = pyvolve.EmpiricalModelFrequencies("WAG") # model name is case-insensitive
3 freqs = f.compute_frequencies()
4
5 # Construct amino-acid model with WAG frequencies
6 aa_model = pyvolve.Model( "WAG", {"state_freqs":freqs} )

```

4.3 Mechanistic (dN/dS) codon models

GY-style [3] matrix elements, for the substitution from codon i to j , are generally given by

$$q_{ij} = \begin{cases} \mu_{o_i t_j} \pi_j \alpha & \text{synonymous change} \\ \mu_{o_i t_j} \pi_j \beta & \text{nonsynonymous change} \\ 0 & \text{multiple nucleotide changes} \end{cases}, \quad (3)$$

where $\mu_{o_i t_j}$ is the mutation rate (e.g. for a change AAA to AAC, the corresponding mutation rate would be $A \rightarrow C$), π_j is the frequency of the target *codon* j , α is the rate of synonymous change (dS), and β is the rate of nonsynonymous change (dN).

MG-style [11] matrix elements, for the substitution from codon i to j , are generally given by

$$q_{ij} = \begin{cases} \mu_{o_i t_j} \pi_{t_j} \alpha & \text{synonymous change} \\ \mu_{o_i t_j} \pi_{t_j} \beta & \text{nonsynonymous change} \\ 0 & \text{multiple nucleotide changes} \end{cases}, \quad (4)$$

where $\mu_{o_i t_j}$ is the mutation rate, π_{t_j} is the frequency of the target *nucleotide* t_j (e.g. for a change AAA to AAC, the target nucleotide would be C), α is the rate of synonymous change (dS), and β is the rate of nonsynonymous change (dN).

Both GY-style and MG-style codon models use symmetric mutation rates. Codon models *require* that you provide a dN/dS rate ratio as a parameter in the custom-parameters dictionary. There are several ways to specify this value:

- Specify a single parameter, "omega". This option sets the synonymous rate to 1.
- Specify a single parameter, "beta". This option sets the synonymous rate to 1.
- Specify a two parameters, "alpha" and "beta". This option sets the synonymous rate to α and the nonsynonymous rate to β .

By default, mechanistic codon models have equal mutation rates and equal equilibrium frequencies. Basic mechanistic codon models can be constructed with,

```
1 # Simple GY-style model (specify as GY)
2 gy_model = pyvolve.Model("GY", {"omega": 0.5})
3
4 # Simple MG-style model (specify as MG)
5 mg_model = pyvolve.Model("MG", {"alpha": 1.04, "beta": 0.67})
6
7 # Specifying "codon" results in a *GY-style* model
8 codon_model = pyvolve.Model("codon", {"beta": 1.25})
```

To customize a mechanistic codon model, provide a custom-parameters dictionary with optional keys "state_freqs" for custom equilibrium frequencies and "mu" for custom mutation rates (see Section 4.6 for details on frequency customization and Section 4.7 for details on mutation rate customization). Note that codon frequencies must be ordered alphabetically (AAA, AAC, AAG, ..., TTG, TTT) *without* stop codons.

```
1 # Define mutation rates in a dictionary with keys giving the nucleotide pair
2 # Below, the rate from A to C is 0.5, and similarly C to A is 0.5
```

```

3 custom_mu = {"AC":0.5, "AG":0.25, "AT":1.23, "CG":0.55, "CT":1.22, "GT":0.47}
4
5 # Construct codon model with custom mutation rates
6 codon_model = pyvolve.Model( "codon", {"mu":custom_mu, "omega":0.55} )

```

Mechanistic codon model mutation rates are symmetric; if you provide a rate for $A \rightarrow T$ (key "AT"), it will automatically be applied as the rate for $T \rightarrow A$. Any unspecified mutation rate pairs will have a value of 1.

As an alternate to "mu", you can provide the key "kappa", which corresponds to the transition:transversion ratio (e.g. for an HKY85 model [5]), in the custom-parameters dictionary. When kappa is specified, transversion mutation rates are set to 1, and transition mutation rates are set to the provided "kappa" value.

```

1 # Construct codon model with transition-to-transversion bias, and default frequencies
2 codon_model = pyvolve.Model( "codon", {"kappa":2.75, "alpha":0.89, "beta":0.95} )

```

4.4 Mutation-selection models

Mutation-selection (MutSel) model [4] matrix elements, for the substitution from codon (or nucleotide) i to j , are generally given by

$$q_{ij} = \begin{cases} \mu_{ij} \frac{S_{ij}}{1-1/S_{ij}} & \text{single nucleotide change} \\ 0 & \text{multiple nucleotide changes} \end{cases}, \quad (5)$$

where μ_{ij} is the mutation rate, and S_{ij} is the scaled selection coefficient. The scaled selection coefficient indicates the fitness difference between the target and source state, e.g. $fitness_j - fitness_i$. MutSel mutation rates are *not* constrained to be symmetric (e.g. μ_{ij} can be different from μ_{ji}).

MutSel models are implemented both for codons and nucleotides, and they may be specified *either* with equilibrium frequencies or with fitness values. Note that equilibrium frequencies must sum to 1, but fitness values are not constrained in any way. (The relationship between equilibrium frequencies and fitness values for MutSel models is detailed in refs. [4, 13]). Pyvolve automatically determines whether you are evolving nucleotides or codons based on the provided vector of equilibrium frequencies or fitness values; a length of 4 indicates nucleotides, and a length of 61 indicates codons. Note that, if you are constructing a codon MutSel model based on *fitness* values, you can alternatively specify a vector of 20 fitness values, indicating amino-acid fitnesses (in the order A, C, D, E, ... Y). These fitness values will be directly assigned to codons, such that all synonymous codons will have the same fitness.

Basic nucleotide MutSel models can be constructed with,

```

1 # Simple nucleotide MutSel model constructed from frequencies, with default (equal)
  mutation rates
2 nuc_freqs = [0.1, 0.4, 0.3, 0.2]
3 mutsel_nuc_model_freqs = pyvolve.Model("MutSel", {"state_freqs": nuc_freqs})
4
5 # Simple nucleotide MutSel model constructed from fitness values, with default (equal)
  mutation rates
6 nuc_fitness = [1.5, 0.88, -4.2, 1.3]
7 mutsel_nuc_model_fits = pyvolve.Model("MutSel", {"fitness": nuc_fitness})

```


Basic codon MutSel models can be constructed with,

```

1 import numpy as np # imported for convenient example frequency/fitness generation
2
3 # Simple codon MutSel model constructed from frequencies, with default (equal)
  mutation rates
4 codon_freqs = np.repeat(1./61, 61) # constructs a vector of equal frequencies, as an
  example
5 mutsel_codon_model_freqs = pyvolve.Model("MutSel", {"state_freqs": codon_freqs})
6
7 # Simple codon MutSel model constructed from codon fitness values, with default (equal)
  mutation rates
8 codon_fitness = np.random.normal(size = 61) # constructs a vector of normally
  distributed codon fitness values, as an example
9 mutsel_codon_model_fits = pyvolve.Model("MutSel", {"fitness": codon_fitness})
10
11 # Simple codon MutSel model constructed from *amino-acid* fitness values, with default
  (equal) mutation rates
12 aa_fitness = np.random.normal(size = 20) # constructs a vector of normally distributed
  amino-acid fitness values, as an example
13 mutsel_codon_model_fits2 = pyvolve.Model("MutSel", {"fitness": aa_fitness})

```

Mutation rates can be customized with either the "mu" or the "kappa" key in the custom-parameters dictionary. Note that mutation rates in MutSel models do not need to be symmetric. However, if you a rate for $A \rightarrow C$ (key "AC") and no rate for $C \rightarrow A$ (key "CA"), then Pyvolve will assume symmetry and assign $C \rightarrow A$ the same rate as $A \rightarrow C$. If *neither* pair is provided (e.g. both "AC" and "CA" are not defined), then both will be given a rate of 1.

4.5 Empirical codon model

Matrix elements of the empirical codon model (ECM) [9] are given by,

$$q_{ij} = \begin{cases} s_{ij}\pi_j\kappa(i,j)\alpha & \text{synonymous change} \\ s_{ij}\pi_j\kappa(i,j)\beta & \text{nonsynonymous change} \end{cases}, \quad (6)$$

where s_{ij} is the symmetric, empirical matrix indicating the probability of changing from codon i to j , π_j is the equilibrium frequency of the target codon j , $\kappa(i,j)$ is a mutational parameter indicating transition and/or tranversion bias, and α and β represent dS and dN , respectively. Importantly, because this model is empirically-derived, the parameters $\kappa(i,j)$, α , and β as used in ECM each represent the transition-tranversion bias, synonymous rate, and nonsynonymous rate, respectively, *relative* to the average level present in the PANDIT database [17], from which this model was constructed². The parameter $\kappa(i,j)$ is described in depth in ref. [9], specifically in the second half of the Results section *Application of the ECM*.

Importantly, there are two versions of this model: **restricted** and **unrestricted**. The restricted model restricts instantaneous change to single-nucleotide only, whereas the unrestricted model also allows for double- and triple-nucleotide changes. Pyvolve refers to these models, respectively, as ECMrest and ECMunrest.

²Personally, I would not recommend using any of these parameters when simulating (although they have been fully implemented in Pyvolve), as their interpretation is neither straight-forward nor particularly biological.

By default, Pyvolve assumes that $\kappa(i, j)$, α , and β all equal 1, as well as equal equilibrium codon frequencies. Basic ECM can be constructed by specifying either "ECMrest" or "ECMunrest" (case-insensitive) when defining a Model object,

```
1 # Simple restricted ECM
2 ecm_model = pyvolve.Model("ECMrest")
3
4 # Simple unrestricted ECM
5 ecm_model = pyvolve.Model("ECMunrest")
6
7 # Specifying "ECM" results in a *restricted ECM* model
8 ecm_model = pyvolve.Model("ECM")
```

As with mechanistic codon models, the dS and dN parameters can be specified with custom model parameter dictionary keys α , β , and/or ω (but again, these parameters do not correspond to dN/dS in the traditional sense!):

```
1 # Restricted ECM with dN/dS parameter of 0.75
2 ecm_model = pyvolve.Model("ECMrest", {"omega":0.75})
```

The $\kappa(i, j)$ parameter is specified using the keys "k_ti" for transition bias and "k_tv", for transversion bias. Specifically, "k_ti" corresponds to t_s , and "k_tv" corresponds to t_v in equations 9-11 in ref. [9]. Thus, each of these parameters can be specified as either 0, 1, 2, or 3 (the Pyvolve default is 1).

Finally, equilibrium frequencies can be customized with the "state_freqs" key in the custom model parameters dictionary (see Section 4.6 for details on frequency customization). To specify either the ECMrest or ECMunrest default equilibrium frequencies, use Pyvolve's EmpiricalModelFrequencies class:

```
1 # Define default restricted ECM state frequencies
2 f = pyvolve.EmpiricalModelFrequencies("ECMrest") # model name is case-insensitive. Use
   "ECMunrest" for unrestricted
3 freqs = f.compute_frequencies()
4
5 # Construct restricted ECM model with its default frequencies
6 ecm_model = pyvolve.Model("ECMrest", {"state_freqs":freqs})
```

4.6 Specifying equilibrium frequencies

Equilibrium frequencies can be specified for a given Model object with the key "state_freqs" in the custom-parameters dictionary. This key's associated value should be a list (or numpy array) of frequencies, summing to 1. The values in this list should be ordered alphabetically. For nucleotides, the list should be ordered ACGT. For amino-acids, the list should be ordered alphabetically, with regards to single-letter amino-acids abbreviations: ACDEFGHIKLMNPQRSTVWY. Finally, for codons, the list should be ordered AAA, AAC, AAG, AAT, ACA, ... TTT, *excluding* stop codons.

By default, Pyvolve assumes equal equilibrium frequencies (e.g. 0.25 for nucleotides, 0.05, for amino-acids, 1/61 for codons). These conditions are not, however, very realistic, so I strongly recommend that you specify custom equilibrium frequencies for your simulations. Pyvolve provides a convenient class, called StateFrequencies, to help you with this step, with several child classes:

- **EqualFrequencies** (default)

- Computes equal frequencies
- **RandomFrequencies**
 - Computes (semi-)random frequencies
- **CustomFrequencies**
 - Computes frequencies from a user-provided dictionary of frequencies
- **ReadFrequencies**
 - Computes frequencies from a sequence or alignment file
- **EmpiricalModelFrequencies**³
 - Sets frequencies to default values for a given *empirical* model

All of these classes should be used with the following setup (the below code uses EqualFrequencies as a representative example):

```
1 # Define frequency object
2 f = pyvolve.EqualFrequencies("nucleotide") # or "amino_acid" or "codon", depending on
   your simulation
3 frequencies = f.compute_frequencies() # returns a vector of equilibrium frequencies
```

The constructed vector of frequencies (named "frequencies" in the example above) can then be provided to the custom model parameters dictionary with the key "state_freqs". In addition, to conveniently save this vector of frequencies to a file, use the argument `savefile = <name_of_file>` when calling `.construct_frequencies()`:

```
1 # Define frequency object
2 f = pyvolve.EqualFrequencies("nucleotide")
3 frequencies = f.compute_frequencies(savefile = "my_frequency_file.txt") # returns a
   vector of equilibrium frequencies and saves them to file
```

4.6.1 EqualFrequencies class

Pyvolve uses this class to construct the default equilibrium frequencies. Usage should be relatively straightforward, according to the example above.

4.6.2 RandomFrequencies class

This class is used to compute "semi-random" equilibrium frequencies. The resulting frequency distributions are not entirely random, but rather are virtually flat distributions with some amount of noise.

4.6.3 CustomFrequencies class

With this class, you can provide a dictionary of frequencies, using the argument `freq_dict`, from which a vector of frequencies is constructed. The keys for this dictionary are the nucleotides, amino-acids (single letter abbreviations!), or codons, and the values should be the frequencies. Any states not included in this dictionary will be assigned a 0 frequency, so be sure the values in this dictionary sum to 1.

³Note that this is not actually a child class of `StateFrequencies`, but its behavior is virtually identical.

In the example below, `CustomFrequencies` is used to create a vector of amino-acid frequencies in which aspartate and glutamate each have a frequency of 0.25, and tryptophan has a frequency of 0.5. All other amino acids will have a frequency of 0.

```
1 # Define CustomFrequencies object
2 f = pyvolve.CustomFrequencies("amino_acid", freq_dict = {"D":0.25, "E":0.25, "W":0.5})
3 frequencies = f.compute_frequencies()
```

4.6.4 ReadFrequencies class

The `ReadFrequencies` class can be used to compute equilibrium frequencies from a file of sequences and/or multiple sequence alignment. Frequencies can be computed either using all data in the file, or, if the file contains an alignment, using specified alignment column(s). Note that Pyvolve will ignore all ambiguous characters present in this sequence file.

When specifying a file, use the argument `file`, and to specify the file format (e.g. "fasta" or "phylip"), use the argument `format`. Pyvolve uses BioPython to read the sequence file, so consult the BioPython AlignIO module documentation (or this nice [wiki](#)) for available formats. Pyvolve assumes a default file format of FASTA, so the `format` argument is not needed when the file is FASTA.

```
1 # Build frequencies using *all* data in the provided file
2 f = pyvolve.CustomFrequencies("amino_acid", file = "a_file_of_sequences.fasta")
3 frequencies = f.compute_frequencies()
```

To read frequencies from a specific column in a multiple sequence alignment, use the argument `columns`, which should be a list (indexed from 1) of integers giving the column(s) which should be considered in frequency calculations.

```
1 # Build frequencies using alignment columns 1 through 5 (inclusive)
2 f = pyvolve.CustomFrequencies("amino_acid", file = "alignment_file.fasta", columns =
    range(1,6))
3 frequencies = f.compute_frequencies()
4
5 # Build frequencies using only phylip-formatted alignment column 15
6 f = pyvolve.CustomFrequencies("amino_acid", file = "alignment_file.phy", format = "
    phylip", columns = 15)
7 frequencies = f.compute_frequencies()
```

4.6.5 EmpiricalModelFrequencies class

The `EmpiricalModelFrequencies` class will return the default vector of equilibrium frequencies for a given empirical model [amino-acid models JTT, WAG, and LG and the codon model ECM, restricted and unrestricted versions (see ref. [9] for details)]. These default frequencies correspond to the frequencies originally published with each respective empirical model. Provide `EmpiricalModelFrequencies` with the name of the desired empirical model to obtain these frequencies:

```
1 # Obtain frequencies for the WAG model
2 f = pyvolve.EmpiricalModelFrequencies("WAG")
3 frequencies = f.compute_frequencies()
```

```

4
5 # For the ECM models, use the argument "ECMrest" for restricted, and "ECMunrest" for
   unrestricted
6 f = pyvolve.EmpiricalModelFrequencies("ECMrest") # restricted ECM frequencies
7 frequencies = f.compute_frequencies()

```

4.6.6 Restricting frequencies to certain states

When using the classes `EqualFrequencies` and `RandomFrequencies`, it is possible to specify that only certain states be considered during calculations using the `restrict` argument, when defining the object. This argument takes a list of states (nucleotides, amino-acids, or codons) which should have non-zero frequencies. All states not included in this list will have a frequency of zero. Thus, by specifying this argument, frequencies will be distributed *only* among the indicated states.

The following example will return a vector of amino-acid frequencies evenly divided among the five specified amino-acids; therefore, each amino acid in the `restrict` list will have a frequency of 0.2.

```

1 # Compute equal frequencies among 5 specified amino acids
2 f = pyvolve.EqualFrequencies("amino_acid", restrict = ["A", "G", "V", "E", "F"])
3 frequencies = f.compute_frequencies()

```

Note that specifying this argument will have no effect on the `CustomFrequencies`, `ReadFrequencies`, or `EmpiricalModelFrequencies` classes.

4.6.7 Converting frequencies between alphabets

When defining a `StateFrequencies` object, you always have to indicate the alphabet (nucleotide, amino acid, or codon) in which frequency calculations should be performed. However, it is possible to have the `.construct_frequencies()` method return frequencies in a different alphabet, using the argument `type`. This argument takes a string specifying the desired type of frequencies returned (either "nucleotide", "amino_acid", or "codon").

This functionality is probably most useful when used with the `ReadFrequencies` class; for example, you might want to obtain amino-acid frequencies from multiple sequence alignment of codons:

```

1 # Define frequency object
2 f = pyvolve.ReadFrequencies("codon", file = "my_codon_alignment.fasta")
3 frequencies = f.compute_frequencies(type = "amino_acid")

```

As another example, you might want to obtain amino-acid frequencies which correspond to equal codon frequencies of 1/61 each:

```

1 f = pyvolve.EqualFrequencies("codon")
2 frequencies = f.compute_frequencies(type = "amino_acid") # returns a vector of amino-
   acid frequencies that correspond to equal codon frequencies

```

Alternatively, you can also go the other way (amino acids to codons):

```

1 f = pyvolve.EqualFrequencies("amino_acid")
2 frequencies = f.compute_frequencies(type = "codon")

```

When converting amino acid to codon frequencies, Pyvolve assumes that there is *no codon bias* and assigns each synonymous codon the same frequency.

4.7 Specifying mutation rates

Nucleotide, mechanistic codon (dN/dS), and mutation-selection (MutSel) models all use nucleotide mutation rates as parameters. By default, mutation rates are equal for all nucleotide changes (e.g. the Jukes Cantor model [7]). These default settings can be customized, in the custom model parameters dictionary, in one of two ways:

1. Using the key "mu" to define custom rates for any/all nucleotide changes
2. Using the key "kappa" to specify a transition-to-transversion bias ratio (e.g. the HKY85 mutation model. [5])

The value associated with the "mu" key should itself be a dictionary of mutation rates, with keys "AC", "AG", "AT", etc, such that, for example, the key "AC" represents the mutation rate from A to C. Importantly, nucleotide and codon models use symmetric mutation rates; therefore, if a rate for "AC" is defined, the same value will automatically be applied to the change C to A. Thus, there are a total of 6 nucleotide mutation rates you can provide for a custom nucleotide and/or mechanistic codon model. Note that any rates not specified will be set to 1.

Alternatively, MutSel models do not constrain mutation rates to be symmetric, and thus, for instance, the "AC" rate may be different from the "CA" rate. Thus, there are a total of 12 nucleotide mutation rates you can provide for a custom MutSel model. Again, if a rate for "AC" but not "CA" is defined, then the "AC" rate will be automatically applied to "CA". Any unspecified nucleotide rate pairs will be set to 1.

```
1 # Example using customized mutation rates to construct a nucleotide model
2 custom_mutation_rates = {"AC":1.5, "AG":0.5, "AT":1.75, "CG":0.6, "CT":1.25, "GT":1.88
3 }
4 my_model = pyvolve.Model("nucleotide", {"mu": custom_mutation_rates})
```

If, instead, the key "kappa" is specified, then the mutation rate for all transitions (e.g. purine to purine or pyrimidine to pyrimidine) will be set to the specified value, and the mutation rate for all transversions (e.g. purine to pyrimidine or vice versa) will be set to 1. This scheme corresponds to the HKY85 [5] mutation model.

```
1 # Example using customized kappa to construct a nucleotide model
2 my_model = pyvolve.Model("nucleotide", {"kappa": 3.5})
```

5 Defining Partitions

Partitions are defined using the `Partition()` class, with two required keyword arguments: `models`, the evolutionary model(s) associated with this partition, and `size`, the number of positions (sites) to evolve within this partition.

```
1 # Define a default nucleotide model
2 my_model = pyvolve.Model("nucleotide")
3
4 # Define a Partition object which evolves 100 positions according to my_model
```

```
5 my_partition = pyvolve.Partition(models = my_model, size = 100)
```

In cases of branch homogeneity (all branches evolve according to the same model), each partition is associated with a single model, as shown above. When branch heterogeneity is desired, a list of models used should be provided to the `models` argument (as detailed, with examples, in Section 8).

6 Evolving sequences

The callable class `Evolver` is Pyvolve's engine for all sequence simulation. Defining an `Evolver` object requires two keyword arguments: `partitions`, either the name of a single partition or a list of partitions to evolve, and `tree`, the phylogeny along which sequences are simulated.

Examples below show how to define an `Evolver` object and then evolve sequences. The code below assumes that the variables `my_partition` and `my_tree` were previously defined using `Partition` and `read_tree`, respectively.

```
1 # Define an Evolver instance to evolve a single partition
2 my_evolver = pyvolve.Evolver(partitions = my_partition, tree = my_tree)
3 my_evolver() # evolve sequences
4
5 # Define an Evolver instance to evolve several partitions
6 my_multipart_evolver = pyvolve.Evolver(partitions = [partition1, partition2, partition3
7     ], tree = my_tree)
8 my_multipart_evolver() # evolve sequences
```

6.1 Evolver output files

Calling an `Evolver` object will produce three output files to the working directory:

1. **simulated_alignment.fasta**, a FASTA-formatted file containing simulated data
2. **site_rates.txt**, a tab-delimited file indicating to which partition and rate category each simulated site belongs (described in Section 6.2.1)
3. **site_rates_info.txt**, a tab-delimited file indicating the rate factors and probabilities associated with each rate category (described in Section 6.2.2)

In the context of complete homogeneity, in which all sites and branches evolve according to a single model, the files "site_rates.txt" and "site_rates_info.txt" will not contain much useful information. However, when sites evolve under site-wise and/or branch heterogeneity, these files will provide useful information for any necessary post-processing.

To change the output file names for any of those files, provide the arguments `seqfile` ("simulated_alignment.fasta"), `ratefile` ("site_rates.txt"), and/or `infofile` ("site_rates_info.txt") when *calling* an `Evolver` object:

```
1 # Define an Evolver object
2 my_evolver = pyvolve.Evolver(tree = my_tree, partitions = my_partition)
3 # Evolve sequences with custom file names
4 my_evolver(ratefile = "custom_ratefile.txt", infofile = "custom_infofile.txt", seqfile
5     = "custom_seqfile.fasta" )
```

To suppress the creation of any of these files, define the argument(s) as either `None` or `False`:

```
1 # Only output a sequence file (suppress the ratefile and infofile)
2 my_evolver = pyvolve.Evolver(tree = my_tree, partitions = my_partition)
3 my_evolver(ratefile = None, infofile = None)
```

The output sequence file's format can be changed with the argument `seqfmt`. Pyvolve uses BioPython to write sequence files, so consult the BioPython AlignIO module documentation (or this nice [wiki](#)) for available formats.

```
1 # Save the sequence file as seqs.phy, in phyliip format
2 my_evolver = pyvolve.Evolver(tree = my_tree, partitions = my_partition)
3 my_evolver(seqfile = "seqs.phy", seqfmt = "phyliip")
```

By default, the output sequence file will contain only the tip sequences. To additionally output all ancestral (including root) sequences, provide the argument `write_anc = True` when calling an `Evolver` object. Ancestral sequences will be included with tip sequences in the output sequence file (not in a separate file!). When ancestral sequences are written, the root sequence is denoted with the name "root", and internal nodes are named "internal_node1", "internal_node2", etc. To see precisely to which node each internal node name corresponds, it is useful to print the parsed newick tree with the function `print_tree`, as explained in Section 3.

```
1 # Output ancestral sequences along with the tip sequences
2 my_evolver = pyvolve.Evolver(tree = my_tree, partitions = my_partition)
3 my_evolver(write_anc = True)
```

6.2 Sequence post-processing

In addition to saving sequences to a file, `Evolver` can also return sequences back to you for post-processing in Python. Sequences can be easily obtained using the method `.get_sequences()`. This method will return a dictionary of sequences, where the keys are IDs and the values are sequences (as strings). Note that you must evolve sequences by calling your `Evolver` object before sequences can be returned!

```
1 # Return simulated sequences as dictionary
2 my_evolver = pyvolve.Evolver(tree = my_tree, partitions = my_partition)
3 my_evolver()
4 simulated_sequences = my_evolver.get_sequences()
```

By default, `.get_sequences()` will contain only the tip (leaf) sequences. To include ancestral sequences (root and internal node sequences) in this dictionary, specify the argument `anc = True`:

```
1 simulated_sequences = my_evolver.get_sequences(anc = True)
```

6.2.1 Interpreting the "site_rates.txt" output file

The output file "site_rates.txt" has three columns of data:

- **Site_Index**
 - Indicates a given position in the simulated data (indexed from 1)

- **Partition_Index**
 - Indicates the partition associated with this site
- **Rate_Category**
 - Indicates the rate category index associated with this site

The values in "Partition_Index" are ordered, starting from 1, based on the `partitions` argument list specified when setting up the `Evolver()` instance. Similarly, the values in "Rate_Category" are ordered, starting from 1, based on the rate heterogeneity lists (see Section 7 for details) specified when setting up the `Model()` / `CodonModel()` objects used in the respective partition.

6.2.2 Interpreting the "site_rates_info.txt" output file

The output file "site_rates_info.txt" provides more detailed rate information for each partition. This file has give columns of data:

- **Partition_Index**
 - Indicates the partition index (can be mapped back to the Partition_Index column in "site_rates.txt")
- **Model_Name**
 - Indicates the model name (note that, if no name provided, this is None. Also, only relevant for branch het)
- **Rate_Category**
 - Indicates the rate category index (can be mapped back to the Rate_Category column in "site_rates.txt")
- **Rate_Probability**
 - Indicates the probability of a site being in the respective rate category
- **Rate_Factor**
 - Indicates either the rate scaling factor (for nucleotide and amino-acid models), or dN/dS value for this rate category for codon models

6.3 Simulating replicates

The callable `Evolver` class makes simulating replicates of given modeling scheme straight-forward: simply define an `Evolver` object, and then call this object in a for-loop as many times as needed.

```
1 # Simulate 50 replicates
2 my_evolver = pyvolve.Evolver(tree = my_tree, partitions = my_partition)
3 for i in range(50):
4     my_evolver(seqfile = "simulated_replicate" + str(i) + ".fasta") # Change seqfile
                               name to avoid overwriting!
```

7 Implementing site-wise rate heterogeneity

This section details how to implement heterogeneity in site-wise rates within a partition.

7.1 Implementing site-wise heterogeneity for nucleotide and amino-acid models

In the context of nucleotide and amino-acid models, rate heterogeneity is applied by multiplying the rate matrix by scalar factors. Thus, sites evolving at different rates exhibit the same evolutionary patterns but differ in how quickly evolution occurs. Two primary parameters govern this sort of rate heterogeneity: the rate factors used to scale the matrix, and the probability associated with each rate factor (in other words, the probability that a given site is in each rate category).

Pyvolve models site-rate heterogeneity discretely, using either a discrete gamma distribution or a user-specified discrete rate distribution. Rate heterogeneity is incorporated into a `Model` object with several additional keyword arguments, detailed below.

7.1.1 Gamma-distributed rate categories

Gamma-distributed heterogeneity is specified with two (or three) keyword arguments when initializing a `Model` object:

- `alpha`, the shape parameter of the discrete gamma distribution from which rates are drawn (Note: following convention, $\alpha = \beta$ in these distributions [20]).
- `num_categories`, the number of rate categories to draw
- `rate_probs`, an optional list of probabilities for each rate category. If unspecified, all rate categories are equally probable. This list should sum to 1!

Examples for specifying gamma-distributed rate heterogeneity are shown below.

```
1 # Gamma-distributed heterogeneity for a nucleotide model. Gamma shape parameter is 0.5
  , and 6 categories are specified. All categories have an equal probability
2 nuc_model_het = pyvolve.Model("nucleotide", alpha = 0.5, num_categories = 6)
3
4 # Gamma-distributed heterogeneity for a nucleotide model. Gamma shape parameter is 0.5
  , and 6 categories are specified. Categories are assigned specified probabilities
5 nuc_model_het = pyvolve.Model("nucleotide", alpha = 0.5, num_categories = 6,
  rate_probs = [0.2, 0.3, 0.3, 0.1, 0.05, 0.05])
6
7 # Gamma-distributed heterogeneity for an amino-acid model. Gamma shape parameter is 0.
  5, and 6 categories are specified. All categories have an equal probability
8 aa_model_het = pyvolve.Model("WAG", alpha = 0.5, num_categories = 6)
```

7.1.2 Custom-distributed rate categories

A user-determined heterogeneity distribution is specified with one (or two) arguments when initializing a `Model` object:

- `rate_factors`, a list of scaling factors for each category
- `rate_probs`, an optional list of probabilities for each rate category. If unspecified, all rate categories are equally probable. This list should sum to 1!

Examples for specifying custom rate heterogeneity distributions are shown below.

```

1 # Custom heterogeneity for a nucleotide model, with four equiprobable categories
2 nuc_model_het = pyvolve.Model("nucleotide", rate_factors = [0.4, 1.87, 3.4, 0.001])
3
4 # Custom heterogeneity for a nucleotide model, with four categories, each with a
   specified probability (i.e. rate 0.4 occurs with a probability of 0.15, etc.)
5 nuc_model_het = pyvolve.Model("nucleotide", rate_factors = [0.4, 1.87, 3.4, 0.001],
   rate_probs = [0.15, 0.25, 0.2, 0.5])
6
7 # Gamma-distributed heterogeneity for an amino-acid model, with four equiprobable
   categories
8 aa_model_het = pyvolve.Model("WAG", rate_factors = [0.4, 1.87, 3.4, 0.001])

```

7.2 Implementing site-wise heterogeneity for mechanistic codon models

Due to the nature of mechanistic codon models, rate heterogeneity is not modeled with scalar factors, but with a distinct model for each rate (i.e. dN/dS value) category. To setup heterogeneous codon models, you must define models using a `CodonModel`, rather than a `Model`, object. Defining such objects is virtually the same as defining `Model` objects, except a list of dN/dS values should be provided to the custom-parameters dictionary to account for rate heterogeneity. As with standard codon models, you can provide dN/dS values with keys "omega", "beta", or "alpha" and "beta" together (to incorporate both synonymous and nonsynonymous rate variation) in the custom model parameters dictionary.

By default, each discrete dN/dS category will have the same probability. To specify custom probabilities, provide the argument `rate_probs`, a list of probabilities, when initializing a `Model` object.

Examples for specifying heterogeneous mechanistic codon models are shown below (note that a GY-style model is shown in the examples, but as usual, both GY-style and MG-style are allowed.)

```

1 # Define a heterogeneous codon model with dN/dS values of 0.1, 0.5, 1.0, and 2.5 .
   Categories are, by default, equally likely.
2 codon_model_het = pyvolve.CodonModel("GY", {"omega": [0.1, 0.5, 1.0, 2.5]})
3
4 # Define a heterogeneous codon model with two dN/dS categories: 0.102 (from 0.1/0.98)
   and 0.49 (from 0.5/1.02). Categories are, by default, equally likely.
5 codon_model_het = pyvolve.CodonModel("GY", {"beta": [0.1, 0.5], "alpha": [0.98, 1.02]}
   )
6
7 # Define a heterogeneous codon model with dN/dS values of 0.102 (with a probability of
   0.4) and 0.49 (with a probability of 0.6).
8 codon_model_het = pyvolve.CodonModel("GY", {"beta": [0.1, 0.5], "alpha": [0.98, 1.02]}
   , rate_probs = [0.4, 0.6])

```

7.3 Implementing site-wise heterogeneity for mutation-selection models

Due to the nature of MutSel models, site-wise heterogeneity should be accomplished using a series of partitions, in which each partition evolves according to a unique MutSel model. These partitions can then be provided as a list when defining an `Evolver` object.

7.4 Implementing site-wise heterogeneity for the Empirical Codon Model

Due to the peculiar features of this model (both empirically-derived transition probabilities and "mechanistic" parameters such as dN/dS), site-wise heterogeneity is not supported for these models, at this time. Pyvolve will simply ignore any provided arguments for site-rate heterogeneity with this model. Feel free to email the author to discuss and/or request this feature.

8 Implementing branch (temporal) heterogeneity

This section details how to implement branch (also known as temporal) heterogeneity within a partition, thus allowing different branches to evolve according to different models. To implement branch heterogeneity, your provided newick phylogeny should contain *model flags* at particular nodes of interest. Model flags must be in the format `_flagname_` (i.e. with both a leading and a trailing underscore), and they should be placed after *branch lengths* (not after taxon names or nodes!).

For example, a tree specified as `(t4:0.785, (t3:0.380, (t2:0.806, (t5:0.612, t1:0.660):0.762_m1_):0.921_m2_):0.207);` will be interpreted as in Figure 1. Trees with model flags, just like any other tree, are defined with the function `read_tree`:

```
1 # Define a tree with model flags m1 and m2, as read from a file
2 het_tree = pyvolve.read_tree(file = "/path/to/file/containing/tree/with/flags.tre")
3
4 # Define a tree with model flags m1 and m2, with a string
5 het_tree = pyvolve.read_tree(tree = "(t4:0.785, (t3:0.380, (t2:0.806, (t5:0.612, t1:0.660)
6     :0.762_m1_):0.921_m2_):0.207);")
7
8 # Print het_tree to see how model flags are applied:
9 pyvolve.print_tree(het_tree)
10
11 >>> root None None
12 >>>     t4 0.785 None
13 >>>     internal_node3 0.207 None
14 >>>         t3 0.38 None
15 >>>         internal_node2 0.921 m2
16 >>>             t2 0.806 m2
17 >>>             internal_node1 0.762 m1
18 >>>                 t5 0.612 m1
19 >>>                 t1 0.66 m1
```

Note that model flags may be repeated throughout the tree, but the model associated with each model flag will always be the same. Once a model flag has been placed at a given node, all of that node's children will inherit that model. If a new model is specified in a child node, however, then this model will be applied downstream.

All model flags specified in the newick phylogeny must have corresponding models. To link a model to a model flag, specify a given model's name using the keyword argument `name` when initializing a `Model` object. This name must be identical to a given model flag, *without* the leading and trailing underscores (e.g. the name "m1" corresponds to the flag `_m1_`).

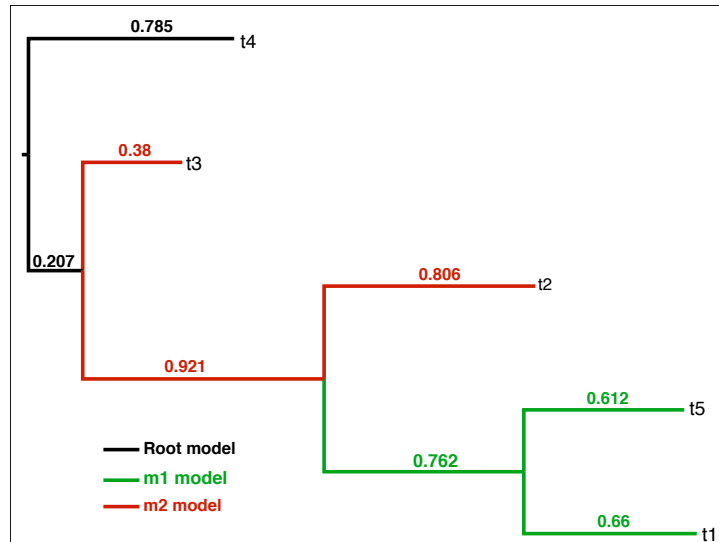


Figure 1: The newick tree with model flags given by

"(t4:0.785,(t3:0.380,(t2:0.806,(t5:0.612,t1:0.660):0.762_m1_):0.921_m2_):0.207);" indicates the model assignments shown.

The model at the root of the tree will not have a specific model flag, but nonetheless a model must be used at the root (obviously), and indeed at all other nodes which are not assigned a model flag (not that all branches on the tree which are not assigned a model flag will evolve according to the model used at the root). To specify a model at the root of the tree, simply create a model, with a name, and indicate this name when defining your partition.

Examples for defining models with names are shown below (for demonstrative purposes, nucleotide models with extreme state frequency differences are used here):

```
1 # Define the m1 model, with frequencies skewed for AT-bias
2 m1_model = pyvolve.Model("nucleotide", {"state_freqs":[0.4, 0.1, 0.1, 0.4]}, name = "
   m1")
3
4 # Define the m2 model, with frequencies skewed for GC-bias
5 m2_model = pyvolve.Model("nucleotide", {"state_freqs":[0.1, 0.4, 0.4, 0.1]}, name = "
   m2")
6
7 # Define the root model, with default equal nucleotide frequencies
8 root_model = pyvolve.Model("nucleotide", name = "root")
```

Alternatively, you can assign/re-assign a model's name with the `.assign_name()` method:

```
1 # (Re-)assign the name of the root model
2 root_model.assign_name("new_root_model_name")
```

Finally, when defining the partition that uses all of these models, provide all `Model` objects in list to the `models` argument. In addition, you *must* specify the name of the model you wish to use at the root of the tree with the keyword argument `root_model_name`.

```
1 # Define partition with branch heterogeneity, with 50 nucleotide positions
```

```

2 temp_het_partition = pyvolve.Partition(models = [m1_model, m2_model, root_model], size
    = 50, root_model_name = root_model.name)

```

9 Implementing branch-site heterogeneity

Simulating according to so-called "branch-site" models, in which there are both site-wise and branch heterogeneity, is accomplished using the same strategies shown for each individual aspect (branch, Section 8 and site, Section 7). However, there is a critical caveat to these models: all models within a given partition *must* have the same number of rate categories. Furthermore, the rate probabilities must be the same across models within a partition; if different values for `rate_probs` are indicated, then the probabilities provided for the *root model* will be applied to all subsequent branch models. (Note that this behavior is identical for other simulation platforms, like Indelible [2].)

The example below shows how to specify a branch-site heterogeneous nucleotide model with two models, root and model1 (note that this code assumes that the provided phylogeny contained the flag `_model1_`), when the rate categories are *not* equiprobable.

```

1 # Shared rate probabilities. Must be explicitly specified for all models (not just the
  #   root model)!
2 shared_rate_probs = [0.25, 0.3, 0.45]
3
4 # Construct a nucleotide model with 3 rate categories
5 root = Model("nucleotide", name = "root", rate_probs = shared_rate_probs, rate_factors
  = [1.5, 1.0, 0.05])
6
7 # Construct a second nucleotide model with 3 rate categories
8 model1 = Model("nucleotide", name = "model1", rate_probs = shared_rate_probs,
  rate_factors = [0.06, 2.5, 0.11])
9
10 # Construct a partition with these models, defining the root model name as "root"
11 part = Partition(models = [root, model1], root_model_name = "root", size = 50)

```

10 Unique Pyvolve Features

This section details features unique to Pyvolve, relative to other sequence simulation platforms.

10.1 Specifying custom rate matrices

Rather than using a built-in modeling framework, you can specify a custom rate matrix. Any provided matrix must be a square matrix with dimensions of 4×4 , 20×20 , or 61×61 (for nucleotide, amino-acid, or codon evolution, respectively). Further, all rows in this matrix must sum to 0. Pyvolve will perform limited sanity checks on your matrix to ensure that these conditions are met, but beyond this, Pyvolve takes your matrix at face-value. In particular, Pyvolve will not scale the matrix in any manner.

Note that you can still specify custom equilibrium frequencies, but these frequencies will not be applied to any matrix scaling - they will be used only for sampling the root sequence in simulation.

To specify a custom rate matrix, provide the argument "custom" as the first argument when defining a `Model` object, and provide your matrix in the custom-parameters dictionary using the key `matrix`. Any custom matrix specified should be either a 2D numpy array or a python list of lists. Pyvolve orders nucleotides in the order ACGT, amino acids in the order ACDEFGHIKLMNPQRSTVWY, and codons in the order AAA, AAC, AAG, ... TTT (without stop codons!). Below is an example of specifying a custom nucleotide rate matrix:

```
1 import numpy as np # import to construct matrix
2
3 # Define the custom rate matrix (4x4 for nucleotide evolution)
4 custom_matrix= np.array([[-1.0, 0.33, 0.33, 0.34],
5                          [0.25, -1.0, 0.25, 0.50],
6                          [0.10, 0.80, -1.0, 0.10],
7                          [0.34, 0.33, 0.33, -1.0]] )
8
9 # Construct a model using the custom rate-matrix
10 custom_model = pyvolve.Model("custom", {"matrix":custom_matrix})
```

Importantly, Pyvolve will extract equilibrium frequencies from this matrix and check that they are acceptable. This frequency vector will be automatically saved to a file called "custom_matrix_frequencies.txt", and these values will be used to generate the root sequence during simulation.

10.2 Matrix scaling options

By convention, rate matrices are scaled such that the mean substitution rate is 1:

$$-\sum_{i=1} \pi_i q_{ii} = 1 \quad (7)$$

[3, 19]. Using this regime, branch lengths explicitly indicate the expected number of substitutions per unit (nucleotide, amino acid, or codon). By default, Pyvolve will scale rate matrices according to equation 7, as this approach remains conventional in the field.

Unfortunately, this scaling approach can lead to some unexpected results for modeling frameworks which contain explicit parameters for natural selection (mechanistic codon and MutSel models). For example, when multiple mechanistic codon (dN/dS) models are used, thus allowing for variable dN/dS values across sites, multiple rate matrices must be used – one matrix per dN/dS value. This scaling approach, then, would cause sites with $dN/dS = 0.05$ to experience the same average number of substitutions as sites with $dN/dS = 2.5$. From a biological perspective, this result is undesirable, as sites with low dN/dS values should evolve more slowly than sites with high dN/dS values. Figure 2 demonstrates the distinction between the conventional scaling approach [equation (7)], which we call "Yang", and our neutral scaling approach.

To overcome this issue, Pyvolve provides an option to scale matrices such that the mean *neutral* substitution rate is 1. For dN/dS codon models, this approach scales the matrix such that the mean number of substitutions when $dN/dS = 1$ is 1. For mutation-selection models (both nucleotide and codon), this approach scales the matrix such that the mean substitution rate is 1 when all states (nucleotides/codons) have equal fitness. Note that invoking neutral-scaling option has no effect on nucleotide or amino-acid models!

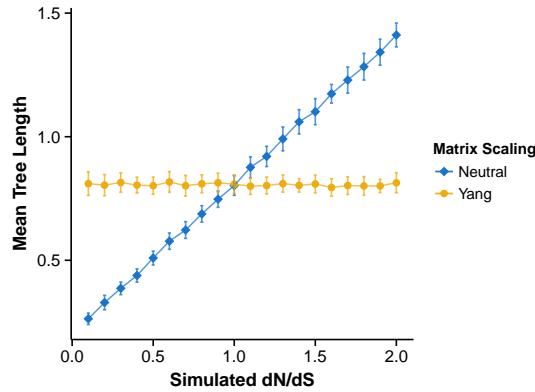


Figure 2: Neutral scaling approach yields more realistic simulations than the Yang (equation (7)) approach. Using Pyvolve, we simulated alignments, of 200 codons each, under the GY94 [3?] mechanistic codon model. All simulations were performed along the same, randomly generated (using the `rtree` function in the R package APE [12]) 25-taxon tree. We simulated alignments (50 replicates each) with a global dN/dS value ranging from 0.1 - 2.0, for each scaling approach. For each resulting alignment, we inferred a maximum-likelihood phylogeny with RAxML [14] under the GTRGAMMA model. We then used DendroPy [15] to calculate a tree length for each inferred phylogeny, indicating the average number of substitutions in the full tree. This figure demonstrates that the average number of substitutions remains constant across dN/dS values when the Yang scaling approach is applied, but this value increases linearly when Pyvolve's neutral scaling approach is applied. Further, the Yang and neutral scaling approaches yield the same number of average substitutions when $dN/dS = 1$, as expected. Error bars represent standard deviations.

To invoke the neutral scaling, provide the keyword argument `scale_matrix = "neutral"` argument "neutral" when defining a Model object:

```
1 # Construct a codon model with neutral-scaling
2 m = pyvolve.Model("GY94", {"omega":0.5}, scale_matrix = "neutral")
```

While we believe that this neutral scaling approach leads to more realistic simulated data, we urge caution when using this scaling approach. Most phylogenetic inference softwares and modeling frameworks (including HyPhy [8] and PAML [21]), scale matrices according to equation 7, and thus inferences on data simulated with neutral scaling may be confounded due to conventions in third-party softwares.

10.3 Perturbing branch lengths

Conventional sequence simulation algorithms apply a given branch length uniformly across a given branch. For example, if a given branch has a length of 0.1, then every site along that branch will evolve with a branch length of exactly 0.1. However, given that phylogenetic inference methods compute branch lengths effectively as an average value for all sites along that branch, there is no reasonable justification to assume all sites should have an identical branch length.

Therefore, Pyvolve allows you to specify some amount of noise in the branch lengths. Specifically, you can opt to draw *site-specific* branch lengths from one of three distributions: normal, gamma, or exponential. Each scheme will retain the given branch length in the newick phylogeny as the average branch length across sites. Given a distribution, Pyvolve will sample, at each branch, a certain number of new branch lengths (default is 10% of the sequence length) and apply these randomly to each evolving site.

To invoke perturbed branch lengths, use the keyword argument `branch_lengths` when *defining* (not calling!) an `Evolver` object. This argument should be a *dictionary* with keys indicating the distribution and any distribution-specific parameters. The distribution should be specified with the key `"dist"` as either "normal", "gamma", or "exp" (for exponential). Normal distributions require the additional key `"sd"` for stan-

dard deviation, and gamma distributions require the additional key of either "shape" or "alpha" (treated equivalently) to specify the shape. Exponential distributions require no other keys.

```
1 # Specify perturbed branch lengths according to different distributions
2
3 # Normal distribution with a standard deviation of 0.1
4 my_evolver = Evolver(partitions = my_partitions, tree = my_tree, branch_lengths = {"
    dist": "normal", "sd": 0.1})
5
6 # Gamma distribution with a shape of 0.5
7 my_evolver = Evolver(partitions = my_partitions, tree = my_tree, branch_lengths = {"
    dist": "gamma", "shape": 0.5}) # can use "alpha" key instead of "shape", if desired
8
9 # Exponential distribution
10 my_evolver = Evolver(partitions = my_partitions, tree = my_tree, branch_lengths = {"
    dist": "exp"})
```

Finally, you can change the number of site-wise branch lengths, for each branch, to sample with the key "num_categories". Specify a number in the range 1 - sequence length, or alternatively to specify a completely unique branch length at every site, you can set the value as "full". (Remember - the default number of branch lengths to draw is 10% of the sequence length).

```
1 # Specify 20 branch length categories per branch (uses exponential distribution as
    example)
2 my_evolver = Evolver(partitions = my_partitions, tree = my_tree, branch_lengths = {"
    dist": "exp", "num_categories": 20})
3
4 # Specify that every site should have a unique branch length at every branch
5 my_evolver = Evolver(partitions = my_partitions, tree = my_tree, branch_lengths = {"
    dist": "exp", "num_categories": "full"})
```

References

- [1] M Arenas. Simulation of molecular data under diverse evolutionary scenarios. *PLoS Comp Biol*, 8:e1002495, 2012.
- [2] W Fletcher and Z Yang. INDELible: A flexible simulator of biological sequence evolution. *Mol Biol Evol*, 26(8):1879–1888, 2009.
- [3] N Goldman and Z Yang. A codon-based model of nucleotide substitution for protein-coding DNA sequences. *Mol Biol Evol*, 11:725–736, 1994.
- [4] AL Halpern and WJ Bruno. Evolutionary distances for protein-coding sequences: modeling site-specific residue frequencies. *Mol Biol Evol*, 15:910–917, 1998.
- [5] M Hasegawa, H Kishino, and T Yano. Dating of human-ape splitting by a molecular clock of mitochondrial DNA. *J Mol Evol*, 22(2):160–174, 1985.
- [6] DT Jones, WR Taylor, and JM Thornton. The rapid generation of mutation data matrices from protein sequences. *CABIOS*, 8:275–282, 1992.
- [7] TH Jukes and CR Cantor. Evolution of protein molecules. In HN Munro, editor, *Mammalian protein metabolism*. Academic Press, New York, 1969.

- [8] SL Kosakovsky Pond, SDW Frost, and SV Muse. HyPhy: hypothesis testing using phylogenies. *Bioinformatics*, 12:676–679, 2005.
- [9] C. Kosiol, I. Holmes, and N. Goldman. An empirical codon model for protein sequence evolution. *Mol Biol Evol*, 24:1464 – 1479, 2007.
- [10] SQ Le and O Gascuel. An improved general amino acid replacement matrix. *Mol Biol Evol*, 25:1307–1320, 2008.
- [11] SV Muse and BS Gaut. A likelihood approach for comparing synonymous and nonsynonymous nucleotide substitution rates, with application to the chloroplast genome. *Mol Biol Evol*, 11:715–724, 1994.
- [12] E. Paradis, J. Claude, and K. Strimmer. APE: analyses of phylogenetics and evolution in R language. *Bioinformatics*, 20:289–290, 2004.
- [13] SJ Spielman and CO Wilke. The relationship between dn/ds and scaled selection coefficients. *Mol Biol Evol*, 32:1097 – 1108, 2015.
- [14] A Stamatakis. RAxML Version 8: A tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics*, 30:1312 – 1313, 2014.
- [15] J. Sukumaran and Mark T. Holder. DendroPy: A Python library for phylogenetic computing. *Bioinformatics*, 26:1569–1571, 2010.
- [16] S Tavaré. Lines of descent and genealogical processes, and their applications in population genetics models. *Theor Popul Biol*, 26:119–164, 1984.
- [17] S Whelan, P de Bakker, E Quevillion, N Rodriguez, and N Goldman. PANDIT: an evolution-centric database of protein and associated nucleotide domains with inferred trees. *Nucleic Acids Res*, 34:D327 – D331.
- [18] S Whelan and N Goldman. A general empirical model of protein evolution derived from multiple protein families using a maximum likelihood approach. *Mol Biol Evol*, 18:691–699, 2001.
- [19] Z Yang. Maximum likelihood phylogenetic estimation from DNA sequences with variable rates over sites: approximate methods. *J Mol Evol*, 39:306 – 314, 1994.
- [20] Z. Yang. *Computational Molecular Evolution*. Oxford University Press, 2006.
- [21] Z. Yang. PAML 4: Phylogenetic analysis by maximum likelihood. *Molecular Biology and Evolution*, 24:1586–1591, 2007.