

Phyx: Phylogenetic tools for Unix

Joseph W. Brown[†], Joseph F. Walker[†], and Stephen A. Smith^{*}

Department of Ecology and Evolutionary Biology, University of Michigan, Ann Arbor, Michigan, 48109, USA

[†]Joint first authors.

^{*}Corresponding author

Emails: josephwb@umich.edu, jfwalker@umich.edu, and eebsmith@umich.edu.

What is Phyx?

Phyx (pronounced “fix”) is a set of data analysis programs modeled after POSIX-style command-line tools, to help them be easily incorporated in bioinformatic pipelines. The majority of Phyx programs focus on phylogenetic analyses, which includes a variety of programs to clean data matrices, simulate data, and perform basic phylogenetic analyses. Phyx is an ever expanding library of programs, and we welcome any feature requests through github issue submissions: <https://github.com/FePhyFoFum/phyx/issues>. All programs are open source and Phyx operates under the GPL3 licence: <https://www.gnu.org/licenses/gpl-3.0.html>

Installing Phyx

Installation instructions for both unix and mac may be found at the Phyx wiki: <https://github.com/FePhyFoFum/phyx/wiki/Installation>.

Pheatures of Phyx

The following is a list of the programs that are currently available. For all programs typing the argument `-h` will produce a list of program options. Below is a list of current programs with examples on how to run them. Programs are updated regularly with new options being added. The most up to date list of options for each program can be found in the help menu.

Data formats

Phyx supports the popular formats for sequence alignments (fasta, fastq, phylip, and Nexus) and trees (newick and Nexus) and uses format auto-detection on the fly. Therefore, for nearly all Phyx programs any of these formats (or a combination thereof, e.g. with `pxcat`, below) can be used without user specification.

Repeatability

Phyx will track the programs and commands input to a file called “phyx.logfile”. This allows users to see exactly what settings they used to run a program and can help others replicate experiments (say, but including the logfile as supplementary information to a paper).

Piping

By optionally reading from `stdin` and writing to `stdout`, Phyx provides the ability of programs to pipe the output of one into another, allowing for an efficient processing of data. An example of this would be to if someone wanted to perform a codon alignment, then remove all columns with missing data, and finally make a rough neighbor-joining tree.

```
pxaa2cdn -a amino_acid_alignment -n nucleotide_alignment || pxclsq -p 1.0  
|| pxnj -n 3 -o output_tree_file
```

Examples

A list of example files can be found in the “example_files” folder. Subfolder names correspond to individual program names.

pxaa2cdn

Often times a coding DNA alignment does not end up with the data divided into sets of three (codons) and as a result this may introduce bias into the analysis or make positive selection tests difficult. This program allows the user to first align the amino acid alignment, then using the alignment the user inputs the corresponding nucleotide sequences and the program will return the codon aligned sequence.

```
pxaa2cdn -a AA_Alignment.fa -n Unaligned_Nucleotide.fa -o CDN_aln.fa
```

pxbdfit

Diversification has becoming a rapidly expanding field and as a result tools to analyze the data are essential. This program will fit a diversification model to a tree. The model which is chosen with `-m` may be (the default) birth-death model (`bd`), a pure-birth yule model (`yule`), or the optimal of these two models (`best`) as determined by AIC. The program will return the model parameters (`b,d,r,e`), likelihood, aic and tree statistics.

```
pxbdfit -t bd.tre -m yule
```

pxbdsim

Birth death processes are an essential part to understanding diversification and simulation gives researchers the ability to study these processes using known birth and death parameters. Per-lineage-per-time birth and death rates are specified with the **-b** and **-d** arguments, respectively. The user chooses the termination condition, either specifying the final number of extant taxa (**-e**) or the simulation timeframe (**-t**). In addition, when running with non-zero extinction the user can choose to return a tree that includes all extinct lineages as well by providing the **-s** argument.

```
pxbdsim -e 100 -s -b 1 -d 0.5 -o output_tree_file
```

pxboot

Using a variety of statistical methods for evaluating certainty of phylogenetic trees is essential as all methods have both positives and negatives associated with them. This program will allow the user to create resampled datasets for two of the most commonly used methods: non-parametric bootstrapping and the jackknife. The proportion of data to be incorporated in a jackknife may be specified with **-f** and a random seed may be specified with **-x**.

Example jackknife:

```
pxboot -s Alignment -x 112233 -f 0.50 -o output_of_50_jackknife
```

Example bootstrap:

```
pxboot -s Alignment -p parts -o output_of_bootstrap
```

pxbp

Analyzing similarities among phylogenetic trees has become a growing part of phylogenetics, especially in the field of phylogenomics to determine if a clade is found in a gene tree and in a species tree. This program allows the user to print out all the bipartitions that are in phylogenetic tree file (which may contain multiple trees).

```
pxbp -t Tree.tre -o bp_output
```

pxcat

When developing a supermatrix for an analysis concatenation of the genes is essential and manual programs that perform this for thousands of genes at once are capable of saving users a lot of time importing each gene into visualization software. This program allows the user to specify a variety of different file types to be concatenated together and can print partition information to a file with the **-p** argument.

An example where the sequences to be concatenated are in a variety of formats:

```
pxcat -s *.fas *.fa *.phy -p Parts.txt -o Supermatrix.fa
```

pxclsq

Having a large amount of missing data in a column of a supermatrix, may be due to errors in alignment or a variety of other factors. Therefore, removing highly ambiguous columns of data may help better estimate a model of evolution for a dataset. This program allows the user to specify a proportion of data that is required to be present (**-p**). The program attempts to detect the sequence type (DNA or protein); if it fails, a protein interpretation can be forced with the **-a** argument. If the verbose **-v** argument is used the program will print to screen the name of any sequences that are entirely removed (i.e. are left with only ambiguous characters after other columns have been removed).

An example to clean a nucleotide alignment down to only columns with a maximum of 40% data allowed to be missing:

```
pxclsq -s Alignment -p 0.6
```

pxconsq

This program will allow the user to get the consensus sequence

```
pxconsq -s Alignment
```

pxcontrates

Comparing continuous characters across phylogenies provides a valuable tool for understanding the evolution of such characters. Two of the most commonly used models are Brownian and OU models, and this program can be used to estimate the rate of character evolution. The input for this is a fasta file where instead of nucleotide data there is tab delimited character states and a tree file for this to be mapped onto. The program may then perform an ancestral state reconstruction (**-a 0** or default) or test for model fit between OU and Brownian motion (**-a 1**).

Example model test for a set of characters across a tree:

```
pxcontrates -c contrates_file.txt -t contrates_tree.tre -a 1
```

pxfqfilt

Filtering based on a certain quality score is essential for processing raw fastq reads from next generation sequencing data. This program allows the user to specify a mean quality score (`-m`) and filter based on that quality score.

```
pxfqfilt -s fqfilt_test.fastq -m 10
```

pxlog

This program is an MCMC log manipulator and concatenator. Resamples parameter or tree MCMC samples using some burnin and thinning across an arbitrary number of log files. Input files may be indicated using wildcards e.g. `*.trees`. NOTE: resampling parameters are in terms of number of samples, not number of generations. To determine the attributes of the log files, you can first use the `-i` flag:

```
pxlog -t *.trees -i
```

and then sample accordingly:

```
pxlog -t *.trees -b 75 -n 2 -o some_output_filename
```

pxlssq

Due to the high variability that is found in sequences and in data matrices it is often important to find out various aspects (e.g. amount of missing data, character state frequencies, etc.). This program will allow provide the user with a variety of these aspects of a the data and provide an easy way to summarize sequence data and concatenated matrices.

```
pxlssq -s Alignment
```

pxlstr

Aspects of trees often provide a large amount of information regarding the behavior of the data that was used to create the tree. This program conveniently allows the user to uncover many of these aspects from the command line, such as tree length, whether a tree is rooted, number of terminals, etc.

```
pxlstr -t Tree.tre
```

pxmrca

This program will provide the information regarding the most recent common ancestor, giving number of tips in the tree and number of tips for each clade specified. The clade that will be analyzed is the smallest clade containing the tips specified. Specifically the user provides the species in a clade of interest using an MRCA file formatted as follows: **MRCANAME = tip1 tip2**

```
pxmrca -t mrca_test.tre -m mrca.txt
```

pxmrcacut

With extremely large trees becoming more common place (species level, gene families etc.) it is useful to focus on certain clades. This program allows the user to specify tips of a clade (-m), only two are required and will remove a newick for the smallest clade that encompasses both species specified. mrca file format: **MRCANAME = tip1 tip2**

```
pxmrcacut -t tree -m mrca_file
```

pxmrca_name

This program allows the user to label the internal nodes with clade names. The program takes in an mrca file in the same format as (pxmrca and pxmrcacut)

```
pxmrca_name -t tree -m mrca_file
```

pxnj

This program will create a basic neighbor joining tree from an alignment matrix.

```
pxnj -s Alignment.aln
```

pxnw

This program will do pairwise alignments using the Needleman-Wunsch algorithm. It also allows alignment scores to be analyzed and various scoring matrices to be used by specifying the -m argument.

```
pxnw -s Alignment.aln
```

pxrecode

This program will recode a DNA alignment to specify only transitions/tranversions (RY-coding).

```
pxrecode -s Nucleotide.fa
```

pxrevcomp

This program will provide the reverse complement of DNA sequences from an alignment file.

```
pxrevcomp -s Nucleotide.fa
```

pxrls

This program allows the user to rename taxa by giving a sequence file and specifying files listing current **-c** and new names **-n**; name ordering in the files must be identical, with one taxon per line.

```
pxrls -s SeqFile -c CurrentNames -n NewNames
```

pxrlt

This program provides a way to relabel the tips of trees by giving a tree file and specifying files listing current **-c** and new names **-n**; name ordering in the files must be identical, with one taxon per line.

```
pxrlt -t kingdoms.tre -c kingdoms.oldnames.txt -n kingdoms.newnames.txt
```

pxrms

This program will remove sequences from a sequence file, either by typing them on the command line using **-n** (comma-delimited) or by specifying a file using **-f** (one taxon per line). Use the **-c** flag to remove the complementary taxa.

```
pxrms -s Nucleotide.fa -f taxa_to_delete.txt
```

pxrmt

This program will remove tips from a tree file, either by typing them on the command line using **-n** (comma-delimited) or by specifying a file using **-f** (one taxon per line). Use the **-c** flag to remove the complementary taxa.

Example to remove tips s1, s6, and s8:

```
pxrmt -t rmt_test.tre -n s1,s6,s8
```

pxrr

This program will re-root trees in a tree file based on specified outgroups (-g), or the program can unroot a tree (-u). For re-rooting, if not all the outgroups are found in the tree the program will print an error. However, **pxrr** can re-root the tree based on the outgroups that *are* available by using the silent option (-s). Alternatively, if the outgroups are ranked in preference but not all necessarily present in a given tree **pxrr** can root on the first outgroup present by using the -r option. It provides a useful tool for re-rooting thousands of trees which can then be used for analyzing gene discordance across phylogenies.

Example to root on the outgroups s1 and s2:

```
pxrr -t rr_test.tre -g s1,s2
```

pxs2fa and pxs2phy and pxs2nex

This programs are all designed in a similar vain, with the ability to convert a file from its current format to fasta, phylip or nexus, respectively. You may also specify if you would like to have the output in uppercase with the option -u.

```
pxs2* -s Alignment
```

pxseqgen

This is a sequence simulator that allows the user to give a tree and specify a model of evolution and sequences will be generated for that tree under the model. Some features are that it allows for the model of evolution to change at nodes along the tree using the -m option. The program also allows the user to specify rate variation through a value for the shape of the gamma distribution with the -g option and the user is able to specify the proportion of invariable sites the would like to include using the -i option. Other options can be found from the help menu by typing -h after the program.

The sequence simulator features have been thoroughly tested except the multimodel simulation which is still under active development and has not been thoroughly tested to the developers comfort!

For multimodel simulations it is easiest to print out the node labels on your tree originally using the -p option. Once you know the nodes that you would like the model to change at you can specify these nodes on the input using the -m option. An example if you wanted two models of evolution on your tree one for the tree and one where it changes at node two, you would enter the command as follows.

example_file example (which uses a simple JC69 model):

```
pxseqgen -t seqgen_test.tre
```


Substitution model parameters are always given in the following order (with no spaces):
A<->C,A<->G,A<->T,C<->G,C<->T,G<->T.

Multi-model commands are given as the following: pxseqgen -t tree_file -m
A<->C,A<->G,A<->T,C<->G,C<->T,G<->T,Node#,A<->C,A<->G,A<->T,C<->G,C<->T,G<->T

If, for example, the model you want for the tree is (.33,.33,.33,.33,.33) and you want the model to change at node two to (.30,.30,.20,.50,.40), the command would be as follows:

```
pxseqgen -t tree_file -o output_alignment  
-m .33,.33,.33,.33,.33,.33,2,.3,.3,.2,.5,.4,.2
```

pxsstat

This program calculates multinomial alignment test statistics that can be used for assessing model adequacy. Currently limited to the test statistic of Bollback (2002) MBE, but more are coming.

```
pxsstat -s Alignment.fa
```

pxstrec

This is a program that does some ancestral state reconstruction and stochastic mapping of categorical characters. There are a number of options and the requirement for a control file. The control file can be as simple as ancstates = all which designates that you want ancestral states calculated for each node. The can then be output on a tree in a file given by an -o FILE option. If you only want to look at particular nodes, these can be designated in the control with the mrca = MRCANAME tipid1 tipid2. Then the MRCANAME can be given at the ancstates = MRCANAME. If you would like stochastic mapping with the time in the state mapped you can use the same format but instead of ancstates you would put stochtime. For stochastic number of events stochnumber or 'stochnumber any. For the stochastic mapping, you will need to designate an MRCA or MRCA's (not all). Multiple can be separated by commas or spaces. You can output these to a file with -n for number of events, -a for the total number of events, and -m for the duration.

```
pxstrec -d test.data.narrow -t test.tre -c config_stmap
```

pxsw

This program will do pairwise alignments using the Smith-Waterman algorithm. It also allows alignment scores to be analyzed and various scoring matrices to be used (-m).

```
pxsw -s Alignment.fa
```

pxt2new

This will convert a tree file to newick format.

```
pxt2new -t Tree.nex
```

pxtlate

This program will take an input coding DNA sequence and translate it to the associated amino acid alignment. The `-t` argument specifies which translation table to use; the standard code is used by default. Use `-h` to see which tables are currently available, and feel free to suggest additions (see alternative translation tables described on genbank).

```
pxtlate -s Sequence.fa
```

pxtscale

This program will rescale a tree when the user inputs either a scaling factor (`-s`) or a root height (`-r`; this option requires an ultrametric tree).

```
pxtscale -t Tree -s 2.0
```

pxupgma

Provides publication quality Unweighted Paired group Method with Arithmetic Mean (UP-GMA) tree, just kidding don't use this for a final tree mainly designed as a teaching tool.

```
pxupgma -s drosophila.aln
```

pxvcf2fa

Convert vcf file to fasta, and can force upper case with `-u`. Currently only handles haploid data; phased data will come soon.

```
pxvcf2fa -s vcf_file
```