

InferNetwork_ML

Description

Infers a species network(s) with a specified number of reticulation nodes using maximum likelihood. The returned species network(s) will have inferred branch lengths and inheritance probabilities. During the search, branch lengths and inheritance probabilities of a proposed species network can be either sampled or optimized. For the first case, after the search, users can ask the program to further optimize those parameters of the inferred network. To optimize the branch lengths and inheritance probabilities to obtain the maximum likelihood for that species network, we use Richard Brent's algorithm (from his book "Algorithms for Minimization without Derivatives", p. 79). The species network and gene trees must be specified in the [Rich Newick Format](#).

The inference can be made using only topologies of gene trees, or using both topologies and branch lengths of gene trees. The latter one requires the input gene trees to be ultrametric.

The input gene trees can be gene tree distributions inferred from Bayesian methods like MrBayes. See the second example below.

Usage

```
InferNetwork_ML geneTreeList numReticulations [-a taxa map] [-bl] [-b threshold] [-s startingNetwork] [-n numNetReturned] [-h {s1 [,s2...]}] [-w (w1,...,w7)] [-x numRuns] [-m maxNetExamined] [-md moveDiameter] [-rd reticulationDiameter] [-f maxFailure] [-o] [-po] [-p (rel,abs)] [-r maxRounds] [-t maxTryPerBr] [-i improveThreshold] [-l maxBL] [-pl numProcessors] [-di] [resultOutputFile]
```

<i>geneTreeList</i>	Comma delimited list of gene tree identifiers or comma delimited list of sets of gene tree identifiers. See details .	mandatory
<i>numReticulations</i>	Maximum number of reticulations to added.	mandatory
<i>-b threshold</i>	Gene trees bootstrap threshold. Edges in the gene trees that have support lower than <i>threshold</i> will be contracted.	optional
<i>-a taxa map</i>	Gene tree / species tree taxa association .	optional
<i>-bl</i>	Use the branch lengths of the gene trees for the inference.	optional
<i>-s startingNetwork</i>	Specify the network to start search. Default value is the optimal MDC tree.	optional
<i>-n numNetReturned</i>	Number of optimal networks to return. Default value is 1.	optional
<i>-h {s1 [, s2...]}</i>	A set of specified hybrid species.	optional
<i>-w (w1, ..., w7)</i>	The weights of operations for network arrangement during the network search. Default value is (0.1,0.1,0.15,0.55,0.15,0.15,2.8).	optional
<i>-x numRuns</i>	The number of runs of the search. Default value is 5.	optional
<i>-m maxNetExamined</i>	Maximum number of network topologies to examined. Default value is infinity.	optional
<i>-md moveDiameter</i>	Maximum diameter to make an arrangement during network search. Default value is infinity.	optional
<i>-rd reticulationDiameter</i>	Maximum diameter for a reticulation event (the distance between two parents of a reticulation node). Default value is infinity.	optional
<i>-f maxFailure</i>	Maximum consecutive number of failures for hill climbing. Default value is 100.	optional
<i>-o</i>	If specified, during the search, for every proposed species network, its branch lengths and inheritance probabilities will be optimized to compute its likelihood. Default value is false.	optional
<i>-po</i>	If specified, after the search the returned species networks will be optimized for their branch lengths and inheritance probabilities. Default value is false.	optional
<i>-p (rel, abs)</i>	The original stopping criterion of Brent's algorithm. Default value is (0.01, 0.001).	optional

<code>-r maxRound</code>	Maximum number of rounds to optimize branch lengths for a network topology. Default value is 100.	optional
<code>-t maxTryPerBr</code>	Maximum number of trial per branch in one round to optimize branch lengths for a network topology. Default value is 100.	optional
<code>-i improveThreshold</code>	Minimum threshold of improvement to continue the next round of optimization of branch lengths. Default value is 0.001.	optional
<code>-l maxBL</code>	Maximum branch lengths considered. Default value is 6.	optional
<code>-pl numProcessors</code>	Number of processors if you want the computation to be done in parallel. Default value is 1.	optional
<code>-di</code>	Output the Rich Newick string of the inferred network that can be read by Dendroscope .	optional
<code>resultOutputFile</code>	Optional file destination for command output.	optional

It is mandatory to specify the number of reticulation nodes to added to the starting network. By default, the inference uses only the topologies of gene trees, however, users can also use both topologies and branch lengths of the gene trees to do the inference, by specifying option `-bl`. By default, it is assumed that only one individual is sampled per species in gene trees. However, the option `[-a taxa map]` allows multiple alleles to be sampled. If users have a prior knowledge of the hybrid species, they can specify them using option `-h`.

The search: Option `-m` allows users to specify the maximum number of networks examined during the search. Users can specify the weights of seven operations for network arrangement through option `-w`. The seven weights correspond to adding a reticulation node, deleting a reticulation node, relocating the head of a reticulation edge, relocating the tail of an edge, reversing the direction of a reticulation edge, replacing a reticulation edge and changing branch lengths and inheritance probabilities, respectively. Furthermore, users can use option `-md` to specify the maximum move diameter of an operation for network rearrangement, like what local-SPR does. Also, users can use option `-rd` to specify the maximum reticulation diameter which is the distance (the number of branches) between the two parents of a reticulation node. In order to avoid getting stuck at some local optimum, it is recommended to performed the search multiple times, which users can specify by option `-x`. The `-s` option allows the users to specify a starting network (can be a tree) for network search. If the starting network is not specified, the optimal tree under MDC (command `infer_ST_MDC`) will be used. If it is not binary, a random resolution will be used. By default, only the first optimal species network will be returned. However, users can use `-n` option to ask for multiple optimal networks.

During the search, by default, simulated annealing is used (See *Salter and Pearl 2001* for details of settings), where the branch lengths and inheritance probabilities are sampled. In this case, through option `-po`, as a post-processing, users can optimize the branch lengths and inheritance probabilities of the species networks returned by the search. If the dataset is not large and a large amount of memory is available, users can use option `-o` to optimize the branch lengths and inheritance probabilities of every proposed network during the search. In this case, simple hill climbing will be used, and only the first 5 operations for network arrangement will be used. If branch lengths of the gene trees are used (option `-bl`), the latter case will be applied.

To optimize the branch lengths and inheritance probabilities of a species network, we use Richard Brent's algorithm (from his book "Algorithms for Minimization without Derivatives", p. 79). Users can use different options to control this process. Option `-p` allows users to specify the original stopping criterion of Brent's algorithm. More precisely, *abs* and *rel* define a tolerance $tol = rel |x| + abs$. We optimize the branch lengths one by one. For every branch, it terminates when either `maxTryPerBr` (option `-t`) trials have been made or the Brent's algorithm suggests so. Users can put an upper bound of the branch lengths through option `-l`. Optimization of all branch lengths consists of a round. After every round, if the improvement in terms of likelihood score is greater than that from last round by at least `improveThreshold` (option `-i`), we starts next round. A maximum of `maxRound` (option `-r`) rounds will be tried.

If users want to run the computation in parallel (in terms of the gene trees). Please specify the number of processors through option `-pl`.

Examples

```
#NEXUS

BEGIN TREES;

Tree geneTree1 = ((C,((B,D),A)),E);
Tree geneTree2 = (B,(D,(C,(A,E))));
Tree geneTree3 = (D,(B,((C,E),A)));
Tree geneTree4 = (D,((B,E),(C,A)));

END;

BEGIN PHYLONET;

InferNetwork_ML (geneTree1,geneTree2,geneTree3,geneTree4) 1;

END;
```

```
#NEXUS

BEGIN TREES;

Tree geneTree1 = [&W 0.9] ((C,((B,D),A)),E);
Tree geneTree2 = [&W 0.1] (B,(D,(C,(A,E))));
Tree geneTree3 = [&W 0.6] (D,(B,((C,E),A)));
Tree geneTree4 = [&W 0.4] (D,((B,E),(C,A)));

END;

BEGIN PHYLONET;

InferNetwork_ML (geneTree1,geneTree2,geneTree3,geneTree4) 1;

END;
```

```
#NEXUS

BEGIN TREES;

Tree geneTree1 = ((C:3,((B:1,D:1):1,A:2):1):1,E:4);
Tree geneTree2 = (B:4,(D:3,(C:2,(A:1,E:1):1):1):1);
Tree geneTree3 = (D:4,(B:3,((C:1,E:1):1,A:2):1):1);
Tree geneTree4 = (D:3,((B:1,E:1):1,(C:1,A:1):1):1);

END;

BEGIN PHYLONET;

InferNetwork_ML (geneTree1,geneTree2,geneTree3,geneTree4) 1 -bl;

END;
```

Command References

- Y. Yu, N. Ristic and L. Nakhleh. Fast algorithms and Heuristics for Phylogenomics under hybridization and incomplete lineage sorting. BMC Bioinformatics, vol. 14, no. Suppl 15, p. S6, 2013.
- Y. Yu, J. Dong, K. Liu, and L. Nakhleh, Maximum Likelihood Inference of Reticulate Evolutionary Histories, Proceedings of the National Academy of Sciences, vol. 111, no. 46, pp. 16448-16453, 2014.

See Also

- [List of PhyloNet Commands](#)