Methods to reconstruct phylogenetic networks accounting for ILS

Céline Scornavacca some slides have been kindly provided by Fabio Pardi

ISE-M, Equipe Phylogénie & Evolution Moléculaires Montpellier, France

February the 2nd, 2017

In a phylogenetic network, a reticulate event is represented as a reticulation, where branches converge to give rise to a new lineage:



In a phylogenetic network, a reticulate event is represented as a reticulation, where branches converge to give rise to a new lineage:



The genome(s) at the start of the new lineage is (are) a composition of those of the parent lineages.

In a phylogenetic network, a reticulate event is represented as a reticulation, where branches converge to give rise to a new lineage:



The genome(s) at the start of the new lineage is (are) a composition of those of the parent lineages.

The evolution of each part independently inherited is described by a *gene tree*

In a phylogenetic network, a reticulate event is represented as a reticulation, where branches converge to give rise to a new lineage:



The genome(s) at the start of the new lineage is (are) a composition of those of the parent lineages.

The evolution of each part independently inherited is described by a *gene tree*

In a phylogenetic network, a reticulate event is represented as a reticulation, where branches converge to give rise to a new lineage:



In a phylogenetic network, a reticulate event is represented as a reticulation, where branches converge to give rise to a new lineage:







An optimization problem where a candidate network is evaluated *on the basis of how well the trees it displays fit the data:*



Many possible formulations:

Data:

Trees with 3 taxa: (inferred from other data)

$$\bigwedge$$

$$\bigwedge_{c \in f} d$$

 $\bigwedge_{l \in f} \qquad \bigwedge_{d \neq b}$

 $\bigwedge_{f \ b} \qquad \bigwedge_{a \ c \ d}$

Goal:

Find the network N with the lower hybridization number such that the triplets are `consistent' with one of the trees displayed by N

subject to constraints on the complexity of ${\cal N}$

An optimization problem where a candidate network is evaluated *on the basis of how well the trees it displays fit the data:*



Many possible formulations:

Data:

Any trees on the same taxa: (inferred from other data)

Goal:

Find the network N with the lower hybridization number such that the input trees are `consistent' with one of the trees displayed by N

subject to constraints on the complexity of ${\cal N}$

An optimization problem where a candidate network is evaluated *on the basis of how well the trees it displays fit the data:*



Many possible formulations:

Data:

Clusters of taxa: $\{a, b\}, \{d, e\}, \{d, e, f\}, \{a, b, c, d, e, f\}, \{e, f\}, \{c, d, e, f\}, \dots$

Goal:

Find the network N with the lower hybridization number such that the input clusters are `explained' by one of the trees displayed by N

subject to constraints on the complexity of ${\cal N}$

An optimization problem where a candidate network is evaluated on the basis of how well the trees it displays fit the data:



Many possible formulations:

Data:

Sequence alignments: (typically given in blocks)



$$A_1 \qquad A_2 \qquad \cdots \qquad A_m$$

Goal:

Find N that minimizes $F(N|A_1, A_2, \dots, A_m) = \sum_{i=1}^{N} \min_{T \in \mathcal{T}(N)} F(T|A_i)$

subject to constraints on the complexity of N. F() is the parsimony score.

Jin et al. Parsimony Score of Phylogenetic Networks: Hardness Results and a Linear-Time Heuristic. TCCB. 2009.

An optimization problem where a candidate network is evaluated *on the basis of how well the trees it displays fit the data:*



Many possible formulations:

Data: Sequence alignments: (typically given in blocks) Goal: Find N that maximises $\Pr(A_1, A_2, \dots, A_m | N) = \prod_{i=1}^{m} \Pr(A_i | N) = \prod_{i=1$

Jin et al. Maximum likelihood of phylogenetic networks. Bioinformatics 2006.

N

Phylogenetic network inference

An optimization problem where a candidate network is evaluated *on the basis of how well the trees it displays fit the data:*

https://phylomemetic.wordpress.com/2015/04/17/phylodag/

http://old-bioinfo.cs.rice.edu/nepal/

Many possible formulations:



Sequence alignments: (typically given in blocks)



Goal: Find N that maximises $\mathbf{Pr}(A_1, A_2, \dots, A_m | N) = \prod_{i=1}^m \mathbf{Pr}(A_i | N) = \prod_{i=1}^m \left(\sum_{T \in \mathcal{T}(N)} \mathbf{Pr}(A_i | T) \mathbf{Pr}(T | N) \right)$

 $\mathcal{T}(N)$:

Jin et al. Maximum likelihood of phylogenetic networks. Bioinformatics 2006.

Some issues

- Searching the space of phylogenetic networks The space of networks with k reticulations is infinite.
- Controlling for Model Complexity
 Because any network with k reticulations provides a more complex model than
 any network with (k-1) reticulations, we must handle the model selection problem
 (AIC, BIC, K-fold cross-validation, ...).
- Identifiability issues

$$\mathbf{Pr}(A_1, A_2, \dots, A_m | N) = \prod_{i=1}^m \mathbf{Pr}(A_i | N) = \prod_{i=1}^m \left(\sum_{T \in \mathcal{T}(N)} \mathbf{Pr}(A_i | T) \mathbf{Pr}(T | N) \right)$$

• Not accounting for ILS and allopolyploidy

Different networks can display the same trees

Some networks display exactly the same trees:



Different networks can display the same trees

Some networks display exactly the same trees:

Because N_1 and N_2 display the same trees, they are equally good to any of the inference methods we saw – no matter the input data

(Recall that a network is evaluated on the basis of how well the trees it displays fit the data)



Different networks can display the same trees

Some networks display exactly the same trees:

Because N_1 and N_2 display the same trees, they are equally good to any of the inference methods we saw – no matter the input data

UNIDENTIFIABILITY

Indistinguishable networks



 N_1 and N_2 display the same trees (i.e. including branch lengths) and are thus *indistinguishable* even to methods accounting for lengths

Indistinguishable networks



 N_1 and N_2 display the same trees (i.e. including branch lengths) and are thus *indistinguishable* even to methods accounting for lengths

Key observation: we can move reticulations up or down (until they hit a speciation node) and the trees displayed by a network remain the same:



Key observation: we can move reticulations up or down (until they hit a speciation node) and the trees displayed by a network remain the same:



Key observation: we can move reticulations up or down (until they hit a speciation node) and the trees displayed by a network remain the same:



Key observation: we can move reticulations up or down (until they hit a speciation node) and the trees displayed by a network remain the same:



By moving reticulations always down, N_1 and N_2 both end up becoming the same network.

True in general: indistinguishable networks always transform into the same network – the *canonical form* of N_1 and N_2 .



What it means for the evolutionary biologist

If N is reconstructed by a "classic" inference method, then even assuming perfect and unlimited data, the best you can hope is that the true phylogenetic network is just one of the many that are indistinguishable from N...



The canonical form of N is a unique representative of the networks indistinguishable from N, that excludes their unrecoverable aspects...

Take home message for the computational biologist

If N is reconstructed by a "classic" inference method, then even assuming perfect and unlimited data, the best you can hope is that the true phylogenetic network is just one of the many that are indistinguishable from N...



Classes of indistinguishable networks

"Classic" network inference methods should only attempt to reconstruct canonical networks .

Are gene trees always displayed by the network?

The approaches above assume that **deep coalescence** cannot occur, so gene trees are necessarily displayed trees, but **this is not always the case**



An example of incomplete lineage sorting (ILS)



ILS and the probability of gene trees



Rannal & Yang 2003 (Genealogies) Degnan & Salter 2005 (Topologies)

ILS and the probability of gene trees



ILS and the probability of gene trees



Estimating the species tree under the coalescent model

Since standard phylogenetic methods can be **inconsistent under ILS** (Degnan & Rosenberg 2006), new methods have been developed to cope for this bias, eg:

- STEM (Kubatko, Carstens & Knowles 2009)
- STAR (Liu, Yu, Pearl & Edwards 2009)
- MP-EST (Liu, Yu & Edwards 2010)
- ...

They estimate a phylogenetic species tree given a sample of gene trees (with branch lengths or not) under the coalescent model.

Estimating the species tree under the coalescent model

Since standard phylogenetic methods can be **inconsistent under ILS** (Degnan & Rosenberg 2006), new methods have been developed to cope for this bias, eg:

- STEM (Kubatko, Carstens & Knowles 2009)
- STAR (Liu, Yu, Pearl & Edwards 2009)
- MP-EST (Liu, Yu & Edwards 2010)

The same can be done with networks!



Phylogenetic network inference under ILS

$$\mathbf{Pr}(A_1, A_2, \dots, A_m | N) = \prod_{i=1}^m \mathbf{Pr}(A_i | N) = \prod_{i=1}^m \sum_{T \in \mathcal{T}(N)} \mathbf{Pr}(A_i | T) \mathbf{Pr}(T | N).$$
?

Phylogenetic network inference under ILS





Phylogenetic network inference under ILS

Y Yu & L Nakhleh's group :

- Yu et al. BMC Bioinformatics 2013 (without branch lengths)
- Yu et al. PNAS 2014 (with branch lengths)
- Wen el al. PLOS Genetics 2016 (Bayesian method)

PhyloNet <u>http://bioinfo.cs.rice.edu/phylonet</u>

But also other groups contributed significantly, e.g. L Kubabtko's group and C Ané's group.

$$\mathbf{Pr}(A_1, A_2, \dots, A_m | N) = \prod_{i=1}^m p(G_i | N).$$

Are gene trees always displayed by the network?

The approaches above assume that deep coalescence cannot occur, so gene trees are necessarily displayed trees, but this is not always the case:

Currently accepted introgression scenario among modern humans, Neanderthals and Denisovans, based on sequenced (ancient) nuclear genomes:

Modern

Kitunesi Rices

Are gene trees always displayed by the network?

The approaches above assume that deep coalescence cannot occur, so gene trees are necessarily displayed trees, but this is not always the case:













Network Multi-Species Coalescent (NMSC)

Given a network with branch lengths (in coalescent units) and inheritance probabilities, we can calculate the probability of every possible gene tree:



Network Multi-Species Coalescent (NMSC)

Given a network with branch lengths (in coalescent units) and inheritance probabilities, we can calculate the probability of every possible gene tree:



Network Multi-Species Coalescent (NMSC)

Given a network with branch lengths (in coalescent units) and inheritance probabilities, we can calculate the probability of every possible gene tree:



network also because of allopolyploidy.

Some issues

- Searching the space of phylogenetic networks The space of networks with k reticulations is infinite.
- Controlling for Model Complexity
 Because any network with k reticulations provides a more complex model than
 any network with (k-1) reticulations, we must handle the model selection problem
 (AIC, BIC, K-fold cross-validation, ...).
- Identifiability issues
- Not accounting for ILS and allopolyploidy

<u>http://phylnet.univ-mlv.fr/</u> <u>http://phylonetworks.blogspot.fr</u>

THANKS FOR YOUR ATTENTION