





What about majority rule?

If $\frac{\text{speciation rate}}{\text{mutation rate}} < 4$, then *any* method loses *all* information about the ancestral state as *t* grows (we'll see why in 10 mins!).

Theorem [Mossel +S, 2014]

$$\Pr(\text{MR correct}) > \frac{1}{2} + \frac{1}{2} \left(1 - \frac{4m}{\lambda}\right)$$

All t

Part 2: Information loss on trees

• Probability Primer:

- Let X and Y be any two discrete random variables.
- The *mutual information* of (*X*,*Y*) is:

$$I(X;Y) = \sum_{x,y} p(x,y) \log\left(\frac{p(x,y)}{p(x)p(y)}\right)$$
$$= D_{\mathrm{KL}}(p_{XY}||p_Xp_Y)$$

 $I(X;Y) = 0 \iff X \text{ and } Y \text{ are independent}$

Fano's lemma

 $I(X;Y) << 1 \Rightarrow Y$ cannot accurately predict X by any method!

15

13

What can we say in general (for Yule trees)?



14

Assume any conservative GTR model on any number of states

- □ **Theorem 1:** The accuracy of any method in predicting the state at the **root node** of the tree vanishes (as *n* or *t* goes to infinity) when the mutation rate passes a threshold (dep. on the speciation rate).
- □ **Theorem 2:** There is a very simple method that can predict the state of a **randomly selected node** with an accuracy that does not vanish (as *n* or *t* grows) for any fixed mutation rate.

Markovian processes that destroy information (exponentially fast)

Theorem: For **any** finite state Markov chain, where a transition from any state to any other is possible in some fixed number of steps with probability $\ge p > 0$ then:

$$I(X_0; X_t) < Ce^{-ct}$$

- Transition rates can vary arbitrarily with time
- Theorem applies to a (1-dimensional) Markov **chain**, not for a Markov process on a (branching) tree (but applies to a path from root to leaf).



Data processing inequality



19

¹ S. Roch (Pers. Comm).

A caution...

- Throwing away data never increases information (e.g. deleting fast evolving sites for tree estimation) – same for MLE under true model.
 - But when don't know the 'true' model (i.e. always except in simulations!) it can still (sometimes) be a good thing to do (to avoid correcting under an incorrect model)...

18

Another property of I: subadditivity

If (Y_1, Y_2, \ldots, Y_k) are conditionally independent given X then:

$$I(X; (Y_1, Y_2, \dots, Y_k)) \le \sum_{i=1}^k I(X; Y_i)$$

Example: Data = k characters ($c_1, c_2,...$) generated i.i.d. by an unknown tree topology T, with some prior on branch lengths.

$$I(T; \text{Data}) \le kI(T; c_1)$$







Recent developments

Does 'rates across sites' help?

• For <u>long</u> edges it can: for certain distributions instead of $k = \Theta(\exp(cT))$ it suffices to have $k = \Theta(T^{\gamma})$

[Martyn and S. (2012). JTB 314: 157-163.]

What about additional noise from lineage sorting?

□ For <u>short</u> edges, the sequence length can still be kept at

 $k = \Theta\left(\frac{1}{\epsilon^2}\right)$

Dasarathy, Nowak, Roch (2014). [Data requirement for phylogenetic inference from multiple loci: A new distance method. ArXiv: April 28, 2014:]

Fine, but what about 'evolved' data

Suppose we evolve k characters independently on a tree under a 2-state symmetric model with



 $p(e) \in [p, P]$ for every edge e

Theorem 1 [Erdos, PL, Szkeley, S, Warnow (1999)] For some ('stringy') trees accurate tree reconstruction is possible with $k = \Theta(\log(n))$

But for other ('bushy') trees our approach required $k=\Theta(n^t)$

However, for almost all trees it suffices to have: $k = \Theta(\log(n)^s)$

Conjecture: Provided that $P < \frac{1}{2} \left(1 - \frac{1}{\sqrt{2}} \right)$ accurate tree reconstruction can be achieved for ALL trees with $k = \Theta(\log(n))$

Theorem 2 [Daskalakis, Mossel, Roch (2011)]

This conjecture holds (and is tight)

















Example 1



1338 rooted gene trees on variable taxon sets from the Actinobacteria phylum.

"Primordial tree" in Dendroscope

58

Further details

Gascuel, O. and Steel, M. (2010). Inferring ancestral sequences in taxon-rich phylogenies. *Mathematical Biosciences*, 227: 125-135.

Gascuel, O. and Steel, M. (2014). Predicting the ancestral character changes in a tree is typically easier than predicting the root state. *Systematic Biology*, 63(3): 421-435.

Mossel, E. and Steel, M. (2014). Majority rule has transition ratio 4 on Yule trees under a 2-state symmetric model (ArXiv, April 2014).

Steel, M., Linz, S., Huson, D. and Sanderson, M. (2013) Identifying a species tree subject to random lateral gene transfer. *Journal of Theoretical Biology*, 322: 81--93

Sand, A., and Steel, M. (2013). The standard lateral gene transfer model is statistically consistent for pectinate four-taxon trees. *Journal of Theoretical Biology*, 335: 295-298.

