

## Background

Conditional inequality statements encoded in graphical models are known to produce active constraints on the parameters of the models. Zwiernik and Smith (2011) fully characterised the inequality constraints for binary phylogenetic trees with manifest leaf variables and hidden interior variables.

## Aim

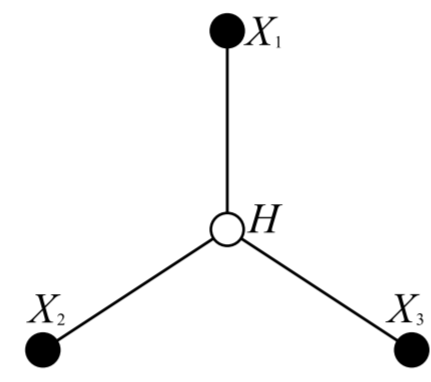
- To demonstrate two graphical diagnostics complementary to existing techniques:
  - ▶ A diagnostic which checks whether data appears to fit with any tree structure.
  - ▶ A diagnostic which utilises functions of associated statistics to guide candidate tree model selection.

## Constraints on a tripod tree

The inequality constraints associated with the tripod tree are of particular interest, as these inequalities must hold for any triple of manifest variables in a strictly trivalent tree (that is trees where all hidden vertices have degree 3). Using the notation of Zwiernik and Smith (2011, Proposition 2.5) we note that an observed joint probability table  $\mathbf{P}$  ( $2 \times 2 \times 2$ ) is consistent with the tripod tree structure if and only if:

$\mu_{123} = 0$  and at least two of  $\mu_{12}$ ,  $\mu_{13}$ ,  $\mu_{23}$  vanish *or*

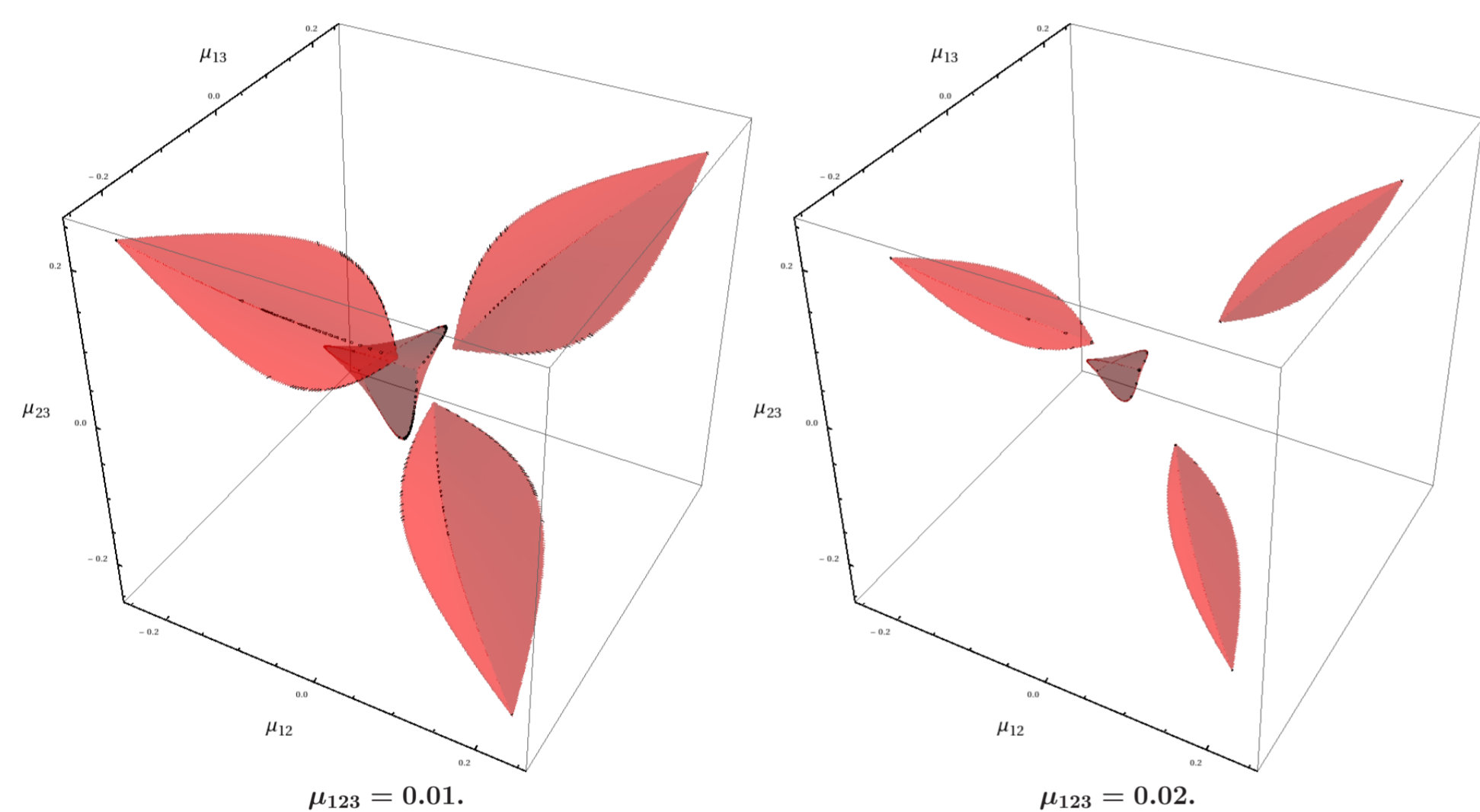
$$\begin{aligned} \mu_{12}\mu_{13}\mu_{23} &> 0 \\ |\mu_{jk}|\sqrt{\det(\mathbf{P})} + \mu_{123}\mu_{jk} &\leq (1 + \bar{\mu}_i)\mu_{jk}^2 \\ |\mu_{jk}|\sqrt{\det(\mathbf{P})} - \mu_{123}\mu_{jk} &\leq (1 - \bar{\mu}_i)\mu_{jk}^2 \end{aligned}$$



where  $\det(\mathbf{P}) = \mu_{123}^2 + 4\mu_{12}\mu_{13}\mu_{23}$ ,  $\bar{\mu}_i = 1 - 2\lambda_i$  for all  $i \in \{1, 2, 3\}$ ,  $\mu$  are central moments, and  $\lambda$  means.

## Admissible covariance regions

We can view these admissible regions (indicated below in red) by fixing  $\mu_{123}$  and  $\bar{\mu}_i$ , and then setting the covariances as the axes in three-dimensional space.



## Theorem 1

The constraints that can be used for a tree diagnostic can be derived simply in the following graphical way:

Given any strictly trivalent phylogenetic tree  $\mathcal{T}$ , for all triples  $\mathbf{X}_i, \mathbf{X}_j, \mathbf{X}_k$  there exists a unique hidden variable  $\mathbf{H}_{ijk}$  such that  $\mathbf{H}_{ijk}$  separates  $(\mathbf{X}_i, \mathbf{X}_j, \mathbf{X}_k)$  in  $\mathcal{T}$ . i.e.

$$\perp\!\!\!\perp (\mathbf{X}_i, \mathbf{X}_j, \mathbf{X}_k) | \mathbf{H}_{ijk}$$

This result allows us to construct a diagnostic test to check whether any tree could be consistent with a sample data set.

## Application of Theorem 1 - Testing for consistency with a tree

We demonstrate the use of the inequality constraints in a phylogenetic context using mitochondrial DNA of length 883 base pairs from BOLD Systems 3 (<http://boldsystems.org>) for six species from the class Mammalia:

- $\mathbf{X}_1$  – *Ailurus fulgens* (red panda)
- $\mathbf{X}_2$  – *Procyon lotor* (raccoon)
- $\mathbf{X}_3$  – *Ailuropoda melanoleuca* (giant panda)
- $\mathbf{X}_4$  – *Ursus maritimus* (polar bear)
- $\mathbf{X}_5$  – *Tremarctos ornatus* (spectacled bear)
- $\mathbf{X}_6$  – *Ursus malayanus* (sun bear)

The genetic data did violate some constraints (e.g.  $(\mathbf{X}_3, \mathbf{X}_4, \mathbf{X}_6)$ ). The two figures below illustrate a covariance triple consistent with a tree structure, and a violation respectively.

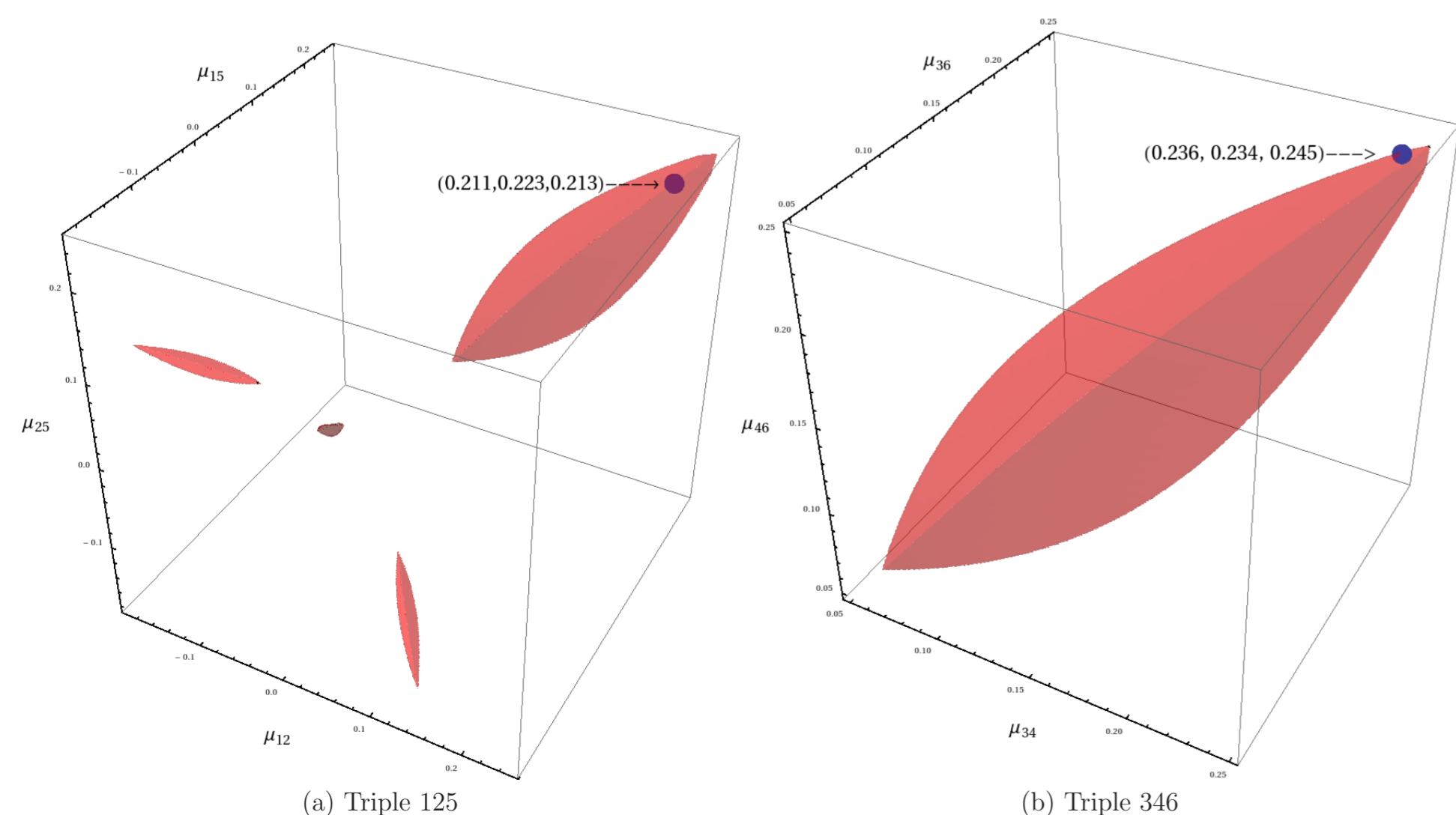


Figure: Point estimates of covariances.

A single violation does not automatically reject a tree structure, as sample errors can cause false positives. A discussion of sample size can be found in Shiers and Smith (2012).

## Theorem 2

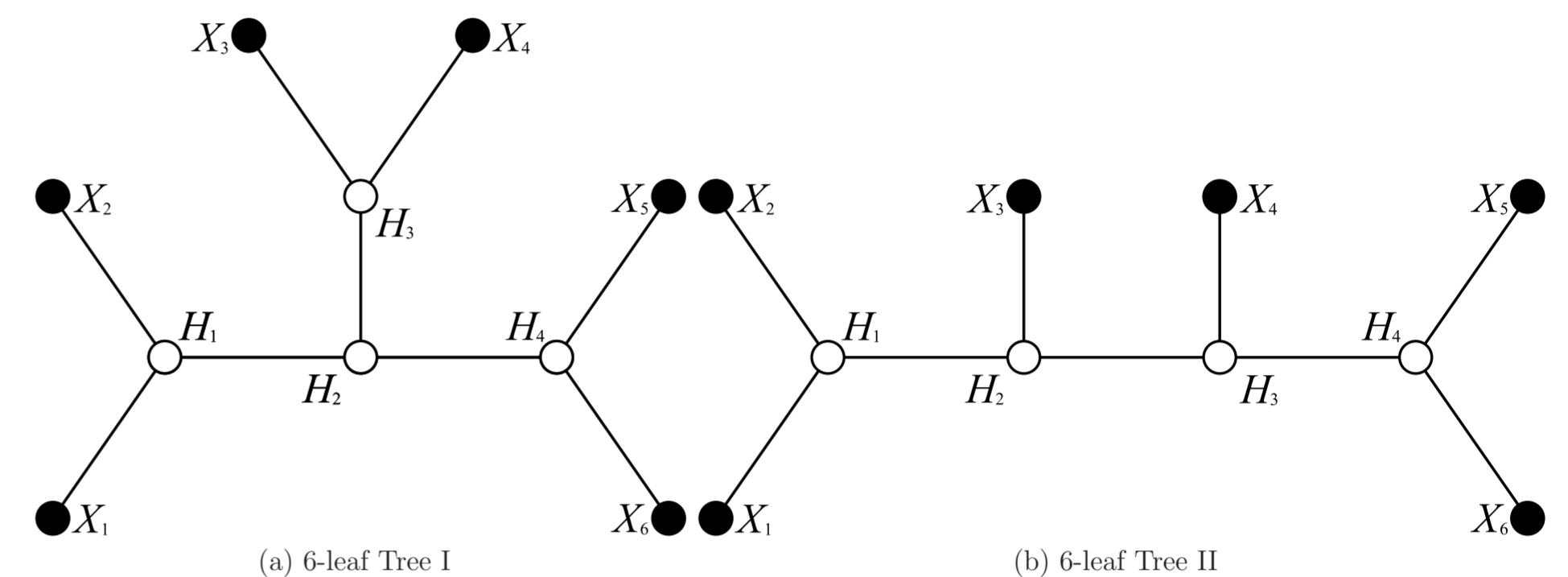
Note for every hidden variable  $\mathbf{H} \in \mathcal{H}$  of a strictly trivalent tree  $\mathcal{T}$ , there is a partition  $\Lambda(\mathbf{H}, \mathcal{T})$  of the manifest variables into 3 subsets, each subset being the leaves of a subtree rooted at  $\mathbf{H}$ . Usefully, these partitions uniquely define a tree  $\mathcal{T}$ . Thus we have:

Each strictly trivalent tree  $\mathcal{T}$  is uniquely identified by its set of partitions and  $\mathcal{X}(\mathcal{T}) \triangleq \{\Lambda(\mathbf{H}, \mathcal{T}) : \mathbf{H} \in \mathcal{H}\}$  acts as an identifier, under the assumption of faithfulness (see Spirtes et al. (2001)).

Thus, if there exists a unique phylogenetic tree, then the estimated moments of the triple will identify it.

## Application of Theorem 2 - Inferring trees from moments

We simulated binary data from the two non-isomorphic strictly trivalent trees with 6 leaves. The probability distributions of the graphs were simulated so expected means for each  $\mathbf{X}_i$  were the same across graphs.



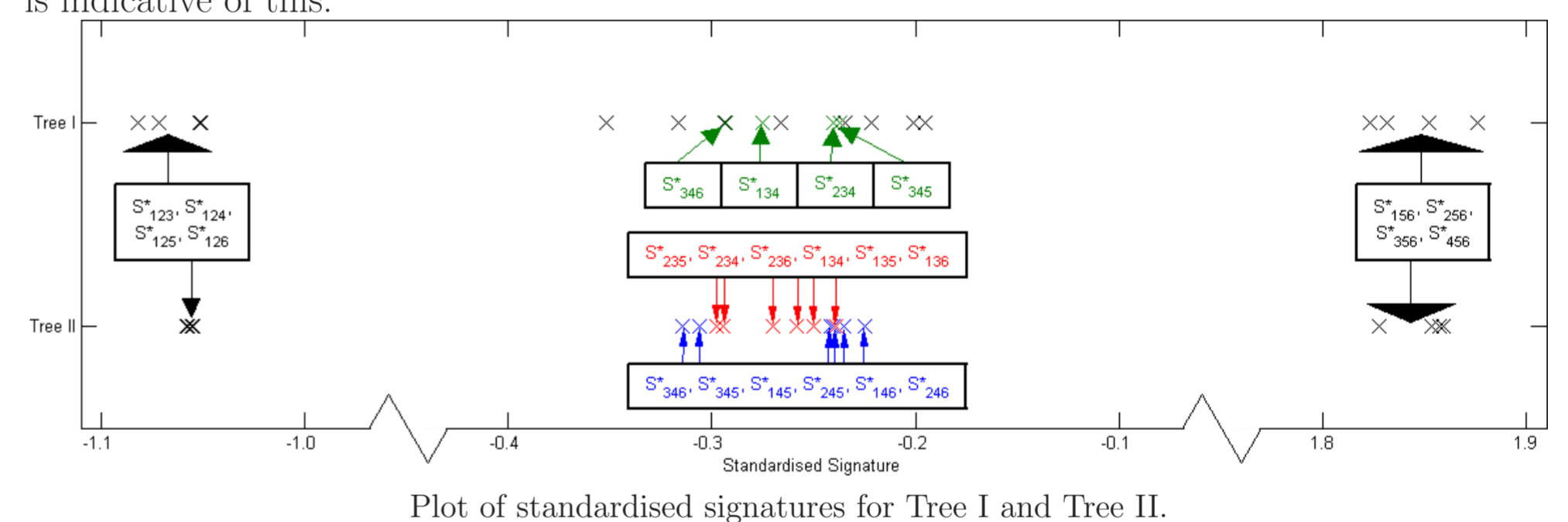
The first 3 moments provide us with consistent but inefficient estimates of a tree. In Settini and Smith (1998) it was shown that (provided none of the terms is zero) for any triple  $\mathbf{X}_i, \mathbf{X}_j, \mathbf{X}_k$  with  $(\mathbf{X}_i \perp\!\!\!\perp \mathbf{X}_j \perp\!\!\!\perp \mathbf{X}_k) | \mathbf{H}_{ijk}$

$$\mathbf{S}_{ijk} = \ln(|\mu_{ij}|) + \ln(|\mu_{ik}|) + \ln(|\mu_{jk}|) - 2\ln(|\mu_{ijk}|)$$

where  $\mathbf{S}_{ijk}$  (the **signature**) depends only on the margin distributions of the triple. From Theorem 2, for large enough datasets the signatures of the corresponding sample quantities will indicate candidate tree partitions  $\mathcal{X}(\mathcal{T})$  and hence, from the theorem, candidates  $\mathcal{T}$ . Note that triples  $(i, j, k)$  and  $(i', j', k')$  share the same separator  $\mathbf{H}$  in  $\mathcal{T}$  if the pairs  $(i, i')$ ,  $(j, j')$  and  $(k, k')$  all lie in different subsets in  $\Lambda(\mathbf{H}, \mathcal{T})$ . So these  $(n - 2)$  partitions can be calculated by first clustering the signatures by magnitude into up to  $(n - 2)$  clusters and from this deducing  $\Lambda(\mathbf{H}, \mathcal{T})$ .

## Standardised signature plot

The signature plot below shows the standardised signatures  $\mathbf{S}_{ijk}^*$  for both non-isomorphic trees 6-leafed tree with the signatures of interest labelled. The clustering of the signatures demonstrates the power of the diagnostic - there is remarkably clear separation of all  $\mathbf{S}_{12k}^*$  and  $\mathbf{S}_{156}^*$  which supports the topologies of the trees. The signatures involving 3 and 4 are less distinct, but this may be in part unavoidable overlapping of true clusters. However, trees like Tree I have an interior node and it can be shown that this leads to signatures with a higher variance. The wide spread in Tree I of the middle cluster (overlapping) with the  $\mathbf{X}_3 \mathbf{X}_4$  cluster is indicative of this.



## Summary & remarks

- ▶ Two graphical diagnostics based on inequality violations and low order sample moments.
- ▶ Simple to calculate, and scale-up to larger trees.
- ▶ Complementary to existing techniques.
- ▶ Useful in the preliminary stages of phylogenetic analysis to check consistency with a tree.
- ▶ Can be used to identify starting points for a tree-model searches.
- ▶ If non-tree structure detected, possible to identify subset of variables consistent with tree.
- ▶ There exist (albeit, in general not yet determined) equivalent inequalities for non-binary data, and in a similar manner alternative 'signatures'.

## Future work

- ▶ Development of a general theory associated with sample sizes and with interpretation of signature plots for first and second diagnostics respectively.
- ▶ Development of analogous inequality diagnostics using the same graphical properties but for differently distributed variables.
- ▶ Application to functional acoustic language data in a language phylogenetic tree context.

## References

- Settini, R and J Q Smith. 1998. On the geometry of bayesian graphical models with hidden variables. In *Proc. 14th Conf. Uncertain. Artif. Intell.* Morgan Kaufmann, pages 472–479.
  - Shiers, N and J Q Smith. 2012. Graphical inequality diagnostics for phylogenetic trees. In *Proc. 6th European Workshop Probabilistic Graph. Model. (to appear)*.
  - Spirtes, P, C Glymour, and R Scheines. 2001. *Causation, Prediction, and Search*. MIT Press, 2nd edition.
  - Zwiernik, P and J Q Smith. 2011. Implicit inequality constraints in a binary tree model. *Electron. J. Stat.* 5:1276–1312.
- We gratefully acknowledge the referees of the corresponding paper (Shiers and Smith (2012)) for their thorough and valuable feedback.