

Model selection for mixture models

Anna Magdalena Kedzierska
joint work with Marta Casanellas and Jesus
Fernandez-Sanchez

Department de Matemàtica Aplicada I
Universitat Politècnica de Catalunya
Centre de Regulació Genòmica

Special Sessions on Advances in Algebraic Statistics
March 28, 2010

Content

- 1 Sequence evolution
- 2 Mixture models
- 3 Identifiability
- 4 Summary

Markov process along a tree

- T - unrooted tree on random variables with states $B = \{A, C, G, T\}$
- $L(T)$ - observed bases at the taxa, $n = |L(T)|$
- (M_e) - transition matrices associated to the edges $e \in E(T)$
- $p = p_{x_1 \dots x_n} := \mathbb{P}(X_1 = x_1, \dots, X_n = x_n)$ is the joint probability of the observed states at $L(T)$

Equivariant Evolutionary Models

Equivariant models are in agreement with certain symmetries of the transition matrices imposed by biological assumptions.

general Markov model (GM), general Kimura 3-parameter (K81), Kimura 2-parameter (K80), general Jukes-Cantor (JC69), strand symmetric model (CS05)

no restrictions

$$\begin{pmatrix} a & b & c & d \\ e & f & g & h \\ j & k & l & m \\ n & o & p & q \end{pmatrix}$$

...

$a+3b=1$

$$\begin{pmatrix} a & b & b & b \\ b & a & b & b \\ b & b & a & b \\ b & b & b & a \end{pmatrix}$$

Group structure

Let $\mathbf{G} \subseteq \mathfrak{S}_4$ be group of permutations on the set of nucleotides B and $W := \mathbb{C}^4 = \langle B \rangle_{\mathbb{C}}$ a vector space over \mathbb{C} generated by B .

Definition

An equivariant evolutionary model is a pair $\mathcal{M}^{\mathbf{G}} = (\mathbf{G}, W)$ such that the transition matrices preserve the action of \mathbf{G} on W .

Group action

Example (JC69, $\mathbf{G}^{JC69} = \mathfrak{S}_4 = \langle (ACGT), (AC) \rangle$)

$$\begin{pmatrix} a & b & b & b \\ b & a & b & b \\ b & b & a & b \\ b & b & b & a \end{pmatrix} \quad (AC) \mapsto \begin{pmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

Permutations of rows and columns leave the matrix unchanged.

JC69 transition matrices coincide with the matrices invariant under the group action \mathbf{G}^{JC69} . We call such set of matrices $\text{Hom}_{\mathbf{G}^{JC}}(W, W)$

$$\mathbf{G}^{GM} = id, \mathbf{G}^{K81} = \langle (AC)(GT), (AG)(CT) \rangle, \mathbf{G}^{K80} = \langle (ACGT), (AG) \rangle, \\ \mathbf{G}^{JC69} = \mathfrak{S}_4$$

The entries of M_e are the model parameters and depend on the chosen model of evolution.

For a given tree, we have a map

$$\Psi_T : \text{parameters} \mapsto \mathcal{L} = \{\text{joint probabilities } p \in \mathbb{R}^{4^n}\}$$

If we consider an equivariant model \mathcal{M}^G , then the parametrization map Ψ_T restricts to:

$$\Psi_T^G : \prod_{e \in E(T)} \text{Hom}_{\mathbf{G}}(W, W) \mapsto \mathcal{L}^G := \{p : p_{x_1 \dots x_n} = p_{g(x_1) \dots g(x_n)} \forall g \in \mathbf{G}\}$$

$$\mathcal{L}^G \subseteq \mathcal{L} \cong \mathbb{R}^{4^n}$$

\mathcal{L}^G depends on the model, NOT on a particular choice of T .

Description of \mathcal{L}^G

We will see that the inclusion

$$\mathcal{L}^G := \{p : p_{x_1 \dots x_n} = p_{g(x_1) \dots g(x_n)}\} \subseteq \mathcal{L}$$

is defined by a set of linear polynomial equations. Furthermore, those are the (linear) invariants for any T on \mathcal{M}^G .

\mathcal{L}^G for G^{K81}

$G^{K81} = \{id, (AC)(GT), (AG)(CT), (AT)(CG)\}$. For each element $g \in G^{K81}$ consider $\{p_x : p_x = p_{g(x)}, x = (x_1, \dots, x_n)\}$.

$$S = 4^{n-1}$$

$$\begin{array}{ll}
 p_{AAA} = p_{CCC} = p_{GGG} = p_{TTT} & p_{AGA} = p_{CTC} = p_{GAG} = p_{TCT} \\
 p_{AAC} = p_{CCA} = p_{GGT} = p_{TTG} & p_{AGC} = p_{CTA} = p_{GAT} = p_{TCG} \\
 p_{AAG} = p_{CCT} = p_{GGA} = p_{TTC} & p_{ACC} = p_{CAA} = p_{GTT} = p_{TGG} \\
 p_{AAT} = p_{CCG} = p_{GGC} = p_{TTA} & p_{ACT} = p_{CAG} = p_{GTC} = p_{TGA} \\
 p_{ACA} = p_{CAC} = p_{GTG} = p_{TGT} & p_{ACG} = p_{CAT} = p_{GTA} = p_{TGC} \\
 p_{AGT} = p_{CTG} = p_{GAC} = p_{TCA} & p_{ATA} = p_{CGC} = p_{GCC} = p_{TAT} \\
 p_{ATC} = p_{CGA} = p_{GCT} = p_{TAG} & p_{ATG} = p_{CGT} = p_{GCA} = p_{TAC} \\
 p_{ATT} = p_{CGG} = p_{GCC} = p_{TAA} & p_{AGG} = p_{CTT} = p_{GAA} = p_{TCC}
 \end{array}$$

Denote by $O^G = (O_1, \dots, O_S)$ the sets of indices of p of the elements p assigned to the collection of these equivalence classes. We call O^G the `orbits` under the action of G .

Example: K81

$$S = 4^{n-1}$$

$$\begin{aligned} p_{AAA} &= p_{CCC} = p_{GGG} = p_{TTT} & p_{AGA} &= p_{CTC} = p_{GAG} = p_{TCT} \\ p_{AAC} &= p_{CCA} = p_{GGT} = p_{TTG} & p_{AGC} &= p_{CTA} = p_{GAT} = p_{TCG} \\ & \dots & & \end{aligned}$$

We select a representant from each orbit O_s , denoted by x^s . A set of such representants is called a `transversal set`,

e.g. $\{AAA, AAC, AAG, AAT, ACA, ACC, ACG, ACT, AGA, AGC, AGG, AGT, ATA, ATC, ATG, ATT\}$

Recursive procedure for generating orbits

Since $\mathbf{G}^{K81} \subset \mathbf{G}^{K80}$, can we use \mathcal{L}^{K81} to generate \mathcal{L}^{K80} ?

Note: $\mathbf{G}^{K80} = \langle \mathbf{G}^{K81}, (AG) \rangle$.

Consider the action of (AG) on the elements of the transversal set of $\mathcal{O}^{\mathbf{G}^{K81}}$ (**significant reduction in computation**),

e.g. $(AG)_{\mathcal{P}_{ACA}} = \mathcal{P}_{GCG}$, so the orbits of \mathbf{G}^{K81} containing \mathcal{P}_{ACA} and \mathcal{P}_{GCG} are combined into a new orbit of \mathbf{G}^{K80} .

Similarly, $\mathbf{G}^{JC69} = \langle \mathbf{G}^{K80}, (AC) \rangle$.

Motivation

We have been looking the spaces \mathcal{L}^G and provided a convenient description in terms of linear equations.

Can we say more? Do the spaces \mathcal{L}^G have a phylogenetic interpretation?

We will see that they are related to the phylogenetic mixture models.

Content

- 1 Sequence evolution
- 2 Mixture models**
- 3 Identifiability
- 4 Summary

Mixture models

Let n be the number of taxa, $[n]$ the set of taxa and

$\tau := \{\text{tree topologies on } [n] \text{ up to isomorphism}\}$.

Let Ψ_T be a parametrization of a model of evolution for any $T \in \tau$:

$$\Psi_T : \prod_{e \in E(T)} \text{Hom}_{\mathbf{G}}(W, W) \rightarrow \mathcal{L}, \quad W = \langle B \rangle_{\mathbb{R}}$$

Definition (in the notation of Matsen-Mosztel-Steel '07)

A phylogenetic mixture (on h classes) is any vector p in \mathcal{L} of the form:

$$p = \sum_{i=1}^h \alpha_i \Psi_{T_i}(\zeta_i),$$

where $T_i \in \tau$, ζ_i -edge parameters, $\alpha_i \in \mathbb{C}$ (or \mathbb{R}).

Main result

We denote by \mathcal{D} the set of all phylogenetic mixtures, $\mathcal{D} \subseteq \mathcal{L}$.

Lemma

The set of all phylogenetic mixtures \mathcal{D} is a vector subspace of \mathcal{L}

Proposition

If $\Psi_{\Gamma}^{\mathbf{G}}$ are equivariant maps for a certain equivariant model corresponding to a group $\mathbf{G} \subseteq \mathfrak{S}_4$, then \mathcal{D} coincides with $\mathcal{L}^{\mathbf{G}}$.

Model selection

Model selection is important: phylogenetic analysis and characterization of the evolutionary process (estimating \mathcal{M}_e , hypothesis testing).

- incorrect model affects the results: incorrect MLE, branch lengths, lower accuracy, inconsistency (Sullivan and Swofford, 1997; Cuningham et.al.1998; Kelsey et.al.1999; Yang et.al. 1994; Buckley et.al 2000; Felstenstein 1978; Bruno and Halpren, 1999; Penny et.al.1994; Huelsenbeck and Hillis 1993)
- Likelihood Ratio (LR) Test and Akaike Information Criterion rely on the candidate input tree to calculate the likelihood-approximation by a single 'best' tree.
- Asymptotic distribution of the LR statistics for mixture models is unknown.

Model selection: alternative outlook

Prior to performing phylogenetic inference, one should assess the fitness of a model to the data for *ALL* possible trees (including their mixtures!).

GOAL:

The linear equations of \mathcal{L}^G are the invariants that vanish on the joint distributions for any T . Therefore, they can be used to assess the goodness-of-fit of \mathcal{M}^G to the data.

Empirical data

- $D = (D_1, \dots, D_N)$ - multiple alignment on n taxa;

- $K = 4^n$;

- (n_1, \dots, n_K) - empirical frequencies, $\sum_{k=1}^K n_k = N$;

- $(p_{x_1}, \dots, p_{x_K})$ - joint probabilities, $\sum_{k=1}^K p_{x_k} = 1$;

Let \mathcal{M} and \mathcal{N} be two equivariant models, such that $\mathcal{M} \subseteq \mathcal{N}$:

- $O_{\mathcal{M}} = \{O_1, \dots, O_R\}$ and $O_{\mathcal{N}} = \{Q_1, \dots, Q_S\}$ - orbits

- $(x^s)_{s=1}^R$ and $(y^s)_{s=1}^S$ - transversal sets;

- $o^s = |O_s|$, $s = 1, \dots, R$ and $q^s = |Q_s|$, $s = 1, \dots, S$ - cardinalities of the orbits;

- $(n^s)_{s=1}^R$ and $(m^s)_{s=1}^S$ - cardinalities of the elements of O and Q .

Maximum Likelihood Estimate

Under \mathcal{M} the likelihood takes the form:

$$\mathbf{L}(D \mid \text{parameters}) = \prod_{s=1}^R (\mathbb{p}_{x^s})^{n^s}$$

$\hat{\mathbb{p}}_{x^s} = \frac{n^s}{o^s R}$ is the MLE of \mathbb{p}_{x^s} for all $s = 1, \dots, R$

Therefore if $x_k \in O_s$, $\hat{\mathbb{p}}_{x_k} = \frac{n^s}{o^s R}$.

Since \mathcal{L}^G are linear spaces, $\hat{\mathbb{p}}_{x_k}$ is the global MLE under \mathcal{M}^G .

In fact, $\hat{\mathbb{p}} = (\hat{\mathbb{p}}_{x_1}, \dots, \hat{\mathbb{p}}_{x_K})$ corresponds to the orthogonal projection of \mathbb{p} onto \mathcal{M}^G .

Model selection

$H_0 : p \in \mathcal{L}^{\mathbf{G}_1}$ vs. $H_1 : p \in \mathcal{L}^{\mathbf{G}_2}$

$$\lambda = -2 \log(\mathbf{L}_0 - \mathbf{L}_1) = -2 \sum_{k=1}^K n_k \log e_k \sim \chi^2(S - R)$$

where $e_k = \sum_{s=1}^R \sum_{l=1}^S \frac{n^s o^l}{m^l q^s} \delta_{sl}(k)$, $\delta_{sl}(k) = 1$ for $p_{x_k} \in \mathcal{O}_s \cap \mathcal{Q}_l$ and 0 otherwise

Alternatives: AIC, MCMC using Markov bases, parametric bootstrap

Content

- 1 Sequence evolution
- 2 Mixture models
- 3 Identifiability**
- 4 Summary

Identifiability

Definition

A statistical model is identifiable if the model parameters (tree topology and the continuous parameters) can be uniquely determined from the distribution of the observed random variables.

From now on, we will think of the **generic** identifiability.

Given a phylogenetic mixture, $p = \sum_{i=1}^h \alpha_i \Psi_{T_i}^G(\zeta_i)$, are we always guaranteed to recover (T_i, ζ_i) ? NOT

Identifiability is essential for the consistency of statistical inference

Proving consistency of the ML estimate is based on the argument of the identifiability of the parameters (Wald, 1949)

Recovering a tree from the leaf colourations it generates under a Markov model, M. Steel, 1994

Identifiability of parameters in MCMC Bayesian inference of phylogeny B. Rannala, 2002

Phylogenetic MCMC algorithms are misleading on mixtures of trees, E. Mossel and E. Vigoda, 2005

Performance of maximum parsimony and likelihood phylogenetics when evolution is heterogeneous. B. Kolaczkowski and J.W. Thornton, 2004

Identifiability of two-tree mixtures under group-based models, E. Allman, S. Petrovic, J. Rhodes, and S. Sullivant, 2009

Identifiability

Question: Is an h -mixture identifiable for any h ?

Let \mathcal{D}^h denote the space of all h -mixtures.

if $\mathcal{D}^h = \mathcal{L}^G$, then m -mixtures are not identifiable for $m \geq h$.

[MMS'07] show that

any $p \in \mathcal{D}$ is a phylogenetic mixture of class $h \leq d + 1$,
 $d := \dim \mathcal{D}$.

In other words, $\mathcal{D}^{d+1} = \mathcal{D} = \mathcal{L}^G$.

Since now we know that $\mathcal{D} = \mathcal{L}^G$, we can compute
 $d = \dim \mathcal{D} = \dim \mathcal{L}^G$ for any equivariant model.

Dimension of \mathcal{L}^G

Tools from representation theory allows us to calculate the dimension of the space of mixtures \mathcal{L}^G for any number of taxa under the equivariant models.

Proposition

Dimension of the space of all phylogenetic mixtures for K81, K80, JC69 for n taxa is:

$$\dim_{K81} = 4^{n-1}$$

$$\dim_{K80} = 2^{2n-3} + 2^{n-2}$$

$$\dim_{JC69} = \frac{1}{3}2^{2n-3} + 2^{n-2} + \frac{1}{3}$$

Upper bound on the number of mixtures

Proposition

Let (\mathbf{G}, W) be an equivariant model, n be a number of taxa and let d be the dimension of the algebraic varieties associated to trivalent trees on n taxa under this model, $d = \text{Im} \Psi_{\mathbf{T}}^{\mathbf{G}}$.

For $h \geq \frac{\dim \mathcal{L}^{\mathbf{G}} + 1}{d + 1}$, h -tree mixtures are not identifiable.

Can this bound be lowered?

Upper bound on the number of mixtures: example

Example

K81

$$n = 4: \dim \mathcal{L}^G = 4^{n-1} = 64, d = 3(2n - 3) + 1 = 16,$$

4-tree mixtures are non-identifiable

$$n = 5: \dim \mathcal{L}^G = 256, d = 22,$$

12-tree mixtures are non-identifiable.

JC69

$$n = 4: \dim \mathcal{L}^G = \frac{1}{3}2^{2n-3} + 2^{n-2} + \frac{1}{3} = 15, d = 2n - 2 = 6,$$

3-tree mixtures are non-identifiable

$$n = 5: \dim \mathcal{L}^G = 51, d = 8,$$

7-tree mixtures are non-identifiable

Content

- 1 Sequence evolution
- 2 Mixture models
- 3 Identifiability
- 4 Summary**

Ongoing and future work

Development of a web-based tool for model selection on the real-life data.

We showed that there exists a bound on h , the number of mixtures, for which model parameters are non-identifiable.

Are the $h - 1$ mixtures identifiable?

Fix n and h . Let $\mathcal{M}^0 \subseteq \mathcal{M}^1$.

Assume the topology (continuous parameters) to be non-identifiable under \mathcal{M}^0 for h mixtures. Is it non-identifiable under \mathcal{M}^1 ?

Fix \mathcal{M}^G and h . What is the smallest n such that for $k > n$ \mathcal{M}^G is identifiable?

Acknowledgements:

Sonja Petrović

Mathias Drton

Thank you