Change-point Detection on a Tree to Study Evolutionary Adaptation from Present-day Species

Cécile Ané^{1,2}, <u>Paul Bastide^{3,4}</u>, Mahendra Mariadassou⁴, Stéphane Robin³

¹ Department of Statistics, University of Wisconsin–Madison, WI, 53706, USA
 ² Department of Botany, University of Wisconsin–Madison, WI, 53706, USA
 ³ UMR MIA-Paris, AgroParisTech, INRA, Université Paris-Saclay, 75005, Paris, France

⁴ MaIAGE, INRA, Université Paris-Saclay, 78352 Jouy-en-Josas, France

19 April 2016









Introduction





Dermochelys Coriacea



Homopus Areolatus

Turtles phylogenetic tree with habitats. (Jaffe et al., 2011).

- How can we explain the diversity, while accounting for the phylogenetic correlations ?
- Modelling: a shifted stochastic process on the phylogeny.

Outline



Stochastic Processes on Trees

- 2 Identifiability Problems and Counting Issues
- Statistical Inference
- Turtles Data Set

Principle of the Modeling Shifts Equivalency OU/BM

Stochastic Process on a Tree

(Felsenstein, 1985)



Only *tip* values are observed



Brownian Motion:

$$\mathbb{V}$$
ar $[A \mid R] = \sigma^2 t$
 \mathbb{C} ov $[A; B \mid R] = \sigma^2 t_{AB}$

BM vs OU



Principle of the Modeling

lentifiability Problems and Counting Issues Statistical Inference Turtles Data Set Principle of the Modeling Shifts Equivalency OU/BM

Shifts



BM Shifts in the mean:

$$m_{
m child} = m_{
m parent} + \delta$$

$$\beta_{\mathsf{child}} = \beta_{\mathsf{parent}} + \delta$$

lentifiability Problems and Counting Issues Statistical Inference Turtles Data Set Principle of the Modeling Shifts Equivalency OU/BM

Shifts



BM Shifts in the mean:

$$m_{
m child} = m_{
m parent} + \delta$$

$$\beta_{\mathsf{child}} = \beta_{\mathsf{parent}} + \delta$$

lentifiability Problems and Counting Issues Statistical Inference Turtles Data Set Principle of the Modeling Shifts Equivalency OU/BM

Shifts



BM Shifts in the mean:

$$m_{
m child} = m_{
m parent} + \delta$$

$$\beta_{\mathsf{child}} = \beta_{\mathsf{parent}} + \delta$$

lentifiability Problems and Counting Issues Statistical Inference Turtles Data Set Principle of the Modeling Shifts Equivalency OU/BM

Shifts



BM Shifts in the mean:

$$m_{
m child} = m_{
m parent} + \delta$$

$$\beta_{\mathsf{child}} = \beta_{\mathsf{parent}} + \delta$$

Principle of the Modeling Shifts Equivalency OU/BM

Linear Regression Model



Principle of the Modeling Shifts Equivalency OU/BM

Linear Regression Model



$$W(\alpha) = \text{Diag}(1 - e^{-\alpha(h - t_{\text{pa}(i)})}, 1 \le i \le m + n)$$

$$\lambda = \mu e^{-\alpha h} + \beta_0(1 - e^{-\alpha h})$$

$$OU: \quad Y = TW(\alpha)\Delta^{OU} + E^{OU}$$

Equivalencies

• Number of shifts K fixed, several equivalent solutions.



• Problem of over-parametrization: parsimonious configurations.

Equivalencies

• Number of shifts K fixed, several equivalent solutions.



• Problem of over-parametrization: parsimonious configurations.

Identifiability Problems Number of Parsimonious Solutions Number of Models with K Shifts

Parsimonious Solution : Definition

Definition (Parsimonious Allocation)

Identifiability Problems Number of Parsimonious Solutions Number of Models with K Shifts

Parsimonious Solution : Definition

Definition (Parsimonious Allocation)



Identifiability Problems Number of Parsimonious Solutions Number of Models with K Shifts

Parsimonious Solution : Definition

Definition (Parsimonious Allocation)



Identifiability Problems Number of Parsimonious Solutions Number of Models with K Shifts

Parsimonious Solution : Definition

Definition (Parsimonious Allocation)



Identifiability Problems Number of Parsimonious Solutions Number of Models with K Shifts

Parsimonious Solution : Definition

Definition (Parsimonious Allocation)



Identifiability Problems Number of Parsimonious Solutions Number of Models with K Shifts

Parsimonious Solution : Definition

Definition (Parsimonious Allocation)



Identifiability Problems Number of Parsimonious Solutions Number of Models with K Shifts

Parsimonious Solution : Definition

Definition (Parsimonious Allocation)



Identifiability Problems Number of Parsimonious Solutions Number of Models with K Shifts

Equivalent Parsimonious Allocations

Definition (Equivalency)

Two allocations are said to be *equivalent* (noted \sim) if they are both parsimonious and give the same colors at the tips.

Find one solution Several existing Dynamic Programming algorithms (Fitch, Sankoff, see Felsenstein, 2004).

Enumerate all solutions New recursive algorithm, adapted from previous ones (and implemented in R).

Identifiability Problems Number of Parsimonious Solutions Number of Models with K Shifts

Equivalent Parsimonious Solutions for an OU Model.



Equivalent allocations and values of the shifts - OU.

Identifiability Problems Number of Parsimonious Solutions Number of Models with *K* Shifts

Collection of Models

New Problem Number of Equivalence Classes: $|\mathcal{S}_{\mathcal{K}}^{PI}|$?

•
$$\left|\mathcal{S}_{K}^{PI}\right| \leq {m+n-1 \choose K} = {\# \text{ of edges} \\ \# \text{ of shifts}}$$

- A recursive algorithm to compute $|\mathcal{S}_{K}^{PI}|$ (implemented in R).
- $\mapsto\,$ Generally dependent on the topology of the tree.

• Binary tree:
$$|\mathcal{S}_{K}^{PI}| = {\binom{2n-2-K}{K}} = {\binom{\# \text{ of edges}-\# \text{ of shifts}}{\# \text{ of shifts}}}$$

EM Algorithm Model Selection

EM Algorithm: number of shifts K fixed



$$\begin{aligned} Y_3 \mid Z_2 &\sim \mathcal{N}\Big(Z_2 + \delta, \ \ell_7 \sigma^2\Big) \\ Z_4 \mid Z_1 &\sim \mathcal{N}\Big(Z_1, \ \ell_4 \sigma^2\Big) \end{aligned}$$

$$\log p_{\theta}(Y) = \mathbb{E}_{\theta}[\log p_{\theta}(Z, Y) \mid Y] - \mathbb{E}_{\theta}[\log p_{\theta}(Z) \mid Y]$$

$$p_{\theta}(Z,Y) = p_{\theta}(Z_1) \prod_{1 < j \le m} p_{\theta}(Z_j | Z_{\mathsf{parent}(j)}) \prod_{1 \le i \le n} p_{\theta}(Y_i | Z_{\mathsf{parent}(i)})$$

EM Algorithm Maximize $\mathbb{E}_{\theta}[\log p_{\theta}(Z, Y) \mid Y]$

E step Given
$$\theta^h$$
, compute $p_{\theta^h}(Z \mid Y)$
M step $\theta^{h+1} = \operatorname{argmax}_{\theta} \mathbb{E}_{\theta^h}[\log p_{\theta}(Z, Y) \mid Y]$

EM Algorithm Model Selection

Turtles Data

Model Selection on K



Simulated OU ($\alpha = 3$, $\gamma^2 = 0.1$)

EM Algorithm Model Selection



EM Algorithm Model Selection



EM Algorithm Model Selection



EM Algorithm Model Selection



EM Algorithm Model Selection



EM Algorithm Model Selection



EM Algorithm Model Selection



EM Algorithm Model Selection



EM Algorithm Model Selection



Stochastic Processes on Trees dentifiability Problems and Counting Issues Statistical Inference Curricia Data Statistical Inference

EM Algorithm Model Selection

Model Selection: Penalized Likelihood

Idea
$$\hat{K} = - \operatorname*{argmin}_{0 \le K \le p-1} \frac{n}{2} \log \left(\frac{\left\| Y - \hat{Y}_K \right\|_V^2}{n} \right) - \frac{1}{2} \operatorname{pen}'(K)$$



Stochastic Processes on Trees dentifiability Problems and Counting Issues Statistical Inference Curricia Data Statistical Inference

EM Algorithm Model Selection

Model Selection: Penalized Likelihood

Idea
$$\hat{K} = - \operatorname*{argmin}_{0 \le K \le p-1} \frac{n}{2} \log \left(\frac{\left\| Y - \hat{Y}_K \right\|_V^2}{n} \right) - \frac{1}{2} \operatorname{pen}'(K)$$


Stochastic Processes on Trees dentifiability Problems and Counting Issues Statistical Inference Curricia Data Statistical Inference

EM Algorithm Model Selection

Model Selection: Penalized Likelihood

Idea
$$\hat{K} = - \operatorname*{argmin}_{0 \le K \le p-1} \frac{n}{2} \log \left(\frac{\left\| Y - \hat{Y}_K \right\|_V^2}{n} \right) - \frac{1}{2} \operatorname{pen}'(K)$$



Stochastic Processes on Trees dentifiability Problems and Counting Issues Statistical Inference Curricia Data Statistical Inference

EM Algorithm Model Selection

Model Selection: Penalized Likelihood

Idea
$$\hat{K} = - \operatorname*{argmin}_{0 \le K \le p-1} \frac{n}{2} \log \left(\frac{\left\| Y - \hat{Y}_K \right\|_V^2}{n} \right) - \frac{1}{2} \operatorname{pen}'(K)$$



EM Algorithm Model Selection

Proposition: LINselect Penalty

Proposition (Form of the Penalty and guarantees (α known))

Under our setting: $Y = TW(\alpha)\Delta + \gamma E$ with $E \sim \mathcal{N}(0, V)$, define the penalty:

$$\mathsf{pen}(\mathcal{K}) = A \frac{n-K-1}{n-K-2} \mathsf{EDkhi}\left[\mathcal{K}+2, n-K-2, \exp\left(-\log\left|\mathcal{S}_{\mathcal{K}}^{PI}\right| - 2\log(\mathcal{K}+2)\right)\right]$$

If
$$\kappa < 1$$
, and $p \le \min\left(\frac{\kappa n}{2 + \log(2) + \log(n)}, n - 7\right)$, we get:

$$\mathbb{E}\left[\frac{\left\|\mathbb{E}\left[Y\right]-\hat{Y}_{\hat{K}}\right\|_{V}^{2}}{\gamma^{2}}\right] \leq C(A,\kappa)\inf_{\eta\in\mathcal{M}}\left\{\frac{\left\|\mathbb{E}\left[Y\right]-Y_{\eta}^{*}\right\|_{V}^{2}}{\gamma^{2}}+\left(K_{\eta}+2\right)\left(3+\log(n)\right)\right\}$$

with $C(A, \kappa)$ a constant depending on A and κ only.

Based on Baraud et al. (2009) 🕕

Turtles Dataset



	Habitat	EM
No. of shifts	16	5
No. of regimes	4	6
InL	-133.86	-97.59
$\ln 2/lpha$ (%)	7.44	5.43
$\sigma^2/2\alpha$	0.33	0.22
CPU t (min)	65.25	134.49

(Jaffe et al., 2011)

Colors: habitats. Boxes: selected EM regimes.

CA, PB, MM, SR

Change-point Detection on a Tree

Turtles Dataset





Chelonia mydas

Colors: habitats. Boxes: selected EM regimes.

CA, PB, MM, SR

Turtles Dataset



Colors: habitats. Boxes: selected EM regimes.

CA, PB, MM, SR

Turtles Dataset





Chitra indica

Colors: habitats. Boxes: selected EM regimes.

CA, PB, MM, SR

Change-point Detection on a Tree

Conclusion and Perspectives

A general inference framework for trait evolution models.

Conclusions • Some problems of identifiability arise.

- An EM can be written to maximize likelihood.
- Adaptation of model selection results to non-iid framework.

R codes Available on GitHub:

https://github.com/pbastide/Phylogenetic-EM

- Perspectives Multivariate traits.
 - Deal with uncertainty (tree, data).
 - Use fossil records.

- Y. Baraud, C. Giraud, and S. Huet. Gaussian Model Selection with an Unknown Variance. The Annals of Statistics, 37(2):630–672, Apr. 2009.
- J.-P. Baudry, C. Maugis, and B. Michel. Slope Heuristics: Overview and Implementation. Statistics and Computing, 22(2):455–470, March 2012.
- V. Brault, J.-P. Baudry, C. Maugis, and B. Michel. capushe: Capushe, Data-Driven Slope Estimation and Dimension Jump. R package version 1.0, 2012.
- J. Felsenstein. Phylogenies and the Comparative Method. The American Naturalist, 125(1):1-15, Jan. 1985.
- J. Felsenstein. Inferring Phylogenies. Sinauer Associates, Suderland, USA, 2004.
- A. L. Jaffe, G. J. Slater, and M. E. Alfaro. The Evolution of Island Gigantism and Body Size Variation in Tortoises and Turtles. *Biology letters*, 11(11), November 2011.
- P. Massart. Concentration Inequalities and Model Selection, volume 1896 of Lecture Notes in Mathematics. Springer Berlin Heidelberg, 2007.
- J. C. Uyeda and L. J. Harmon. A Novel Bayesian Method for Inferring and Interpreting the Dynamics of Adaptive Landscapes from Phylogenetic Comparative Data. Systematic Biology, 63(6):902–918, July 2014.

Photo Credits :

- "Parrot-beaked Tortoise Homopus areolatus CapeTown 8" by Abu Shawka - Own work. Licensed under CC0 via Wikimedia Commons

 - "Leatherback sea turtle Tinglar, USVI (5839996547)" by U.S. Fish and Wildlife Service Southeast Region -Leatherback sea turtle/ Tinglar, USVIUploaded by AlbertHerring. Licensed under CC BY 2.0 via Wikimedia Commons

- "Hawaii turtle 2" by Brocken Inaglory. Licensed under CC BY-SA 3.0 via Wikimedia Commons
- "Dudhwalive chitra" by Krishna Kumar Mishra Own work. Licensed under CC BY 3.0 via Wikimedia Commons

- "Lonesome George in profile" by Mike Weston - Flickr: Lonesome George 2. Licensed under CC BY 2.0 via Wikimedia Commons

- "Florida Box Turtle Digon3a", "Jonathan Zander (Digon3)" derivative work: Materialscientist

Thank you for listening









Model Selection

Model Selection with Unknown Variance

Theorem (Baraud et al. (2009))

Under the following setting:

$$Y' = \mathbb{E}\left[Y'\right] + \gamma \mathsf{E}' \quad \textit{with} \quad \mathsf{E}' \sim \mathcal{N}(0, \mathit{I_n}) \quad \textit{and} \quad \mathcal{S}' = \{S'_\eta, \eta \in \mathcal{M}\}$$

If $D_{\eta} = \text{Dim}(S'_{\eta})$, $N_{\eta} = n - D_{\eta} \ge 7$, $\max(L_{\eta}, D_{\eta}) \le \kappa n$, with $\kappa < 1$, and:

$$\Omega' = \sum_{\eta \in \mathcal{M}} (D_\eta + 1) e^{-L_\eta} < +\infty$$

$$\text{lf:} \quad \hat{\eta} = \operatorname*{argmin}_{\eta \in \mathcal{M}} \left\| \mathbf{Y}' - \hat{Y}'_{\eta} \right\|^2 \left(1 + \frac{\operatorname{pen}(\eta)}{N_{\eta}} \right)$$

with:
$$pen(\eta) = pen_{\mathcal{A},\mathcal{L}}(\eta) = A \frac{N_{\eta}}{N_{\eta} - 1} EDkhi[D_{\eta} + 1, N_{\eta} - 1, e^{-L_{\eta}}]$$
, $A > 1$

Then:
$$\mathbb{E}\left[\frac{\left\|\mathbb{E}\left[Y'\right] - \hat{Y}'_{\hat{\eta}}\right\|^{2}}{\gamma^{2}}\right] \leq C(A,\kappa)\left[\inf_{\eta \in \mathcal{M}}\left\{\frac{\left\|\mathbb{E}\left[Y'\right] - Y'_{\eta}\right\|^{2}}{\gamma^{2}} + \max(L_{\eta}, D_{\eta})\right\} + \Omega'\right]$$

Model Selection

IID Framework ($\alpha = 0$)

Assume
$$K_{\eta} = D_{\eta} - 1 \le p - 1 \le n - 8$$
, $\forall \eta \in \mathcal{M}$

Then:

$$\begin{split} \Omega' &= \sum_{\eta \in \mathcal{M}} (D_{\eta} + 1) e^{-L_{\eta}} = \sum_{\eta \in \mathcal{M}} (K_{\eta} + 2) e^{-L_{\eta}} \\ &= \sum_{K=0}^{p-1} \left| \mathcal{S}_{K}^{PI} \right| (K+2) e^{-L_{K}} = \sum_{K=0}^{p-1} \left| \mathcal{S}_{K}^{PI} \right| (K+2) e^{-(\log \left| \mathcal{S}_{K}^{PI} \right| + 2\log(K+2))} \\ &= \sum_{K=0}^{p-1} \frac{1}{K+2} \le \log(p) \le \log(n) \end{split}$$

And:

$$L_{K} \leq \log {\binom{n+m-1}{K}} + 2\log(K+2) \leq K\log(n+m-1) + 2(K+1) \leq p(2+\log(2n-2))$$

Hence, if $p \leq \min\left(\frac{\kappa n}{2 + \log(2) + \log(n)}, n - 7\right)$, then $\max(L_{\eta}, D_{\eta}) \leq \kappa n$ for any $\eta \in \mathcal{M}$. CA, PB, MM, SR Change-point Detection on a Tree Identifiability Issues Simulations Results Multivariate

Model Selection

Non-IID Framework ($\alpha \neq 0$)

Cholesky decomposition: $V = LL^T$ $Y' = L^{-1}Y$ $s' = L^{-1}s$ $E' = L^{-1}E$

$$Y' = \mathbb{E}\left[Y'\right] + \gamma E'$$
, with: $E' \sim \mathcal{N}(0, I_n)$

$$S'_{\eta} = L^{-1}S_{\eta}, \quad \hat{Y}'_{\eta} = \operatorname{Proj}_{S'_{\eta}} Y' = \operatorname*{argmin}_{a' \in S'_{\eta}} \|Y - La'\|_{V}^{2} = L^{-1}\hat{Y}_{\eta}$$
$$\left\|\mathbb{E}[Y] - \hat{Y}_{\hat{\eta}}\right\|_{V}^{2} = \left\|\mathbb{E}[Y'] - \hat{Y}'_{\hat{\eta}}\right\|^{2}, \quad \left\|Y - \hat{Y}_{\eta}\right\|_{V}^{2} = \left\|Y' - \hat{Y}'_{\eta}\right\|^{2}$$

$$\operatorname{Crit}_{MC}(\eta) = \left\| Y' - \hat{Y}'_{\eta} \right\|^{2} \left(1 + \frac{\operatorname{pen}_{A,\mathcal{L}}(\eta)}{N_{\eta}} \right) = \left\| Y - \hat{Y}_{\eta} \right\|_{V}^{2} \left(1 + \frac{\operatorname{pen}_{A,\mathcal{L}}(\eta)}{N_{\eta}} \right)$$

back

Cardinal of Equivalence Classes Number of Tree Compatible Clustering

 $S(0,\infty,\infty)(0,\infty,\infty)$

Cardinal of Equivalence Classes

Initialization For tips

Propagation

$$\mathcal{K}_{k}^{l} = \operatorname*{argmin}_{1 \le p \le K} \left\{ S_{i_{l}}(p) + \mathbb{I} \{ p \ne k \} \right\}$$

$$S_{i}(k) = \sum_{l=1}^{L} S_{i_{l}}(p_{l}) + \mathbb{I} \{ p_{l} \ne k \} , \ \forall (p_{1}, \dots, p_{L}) \in \mathcal{K}_{k}^{1} \times \dots \times \mathcal{K}$$

$$T_i(k) = \sum_{(p_1,\dots,p_L)\in\mathcal{K}_k^1\times\dots\times\mathcal{K}_k^L}\prod_{l=1}^{L}T_{i_l}(p_l) = \prod_{l=1}^{L}\sum_{p_l\in\mathcal{K}_k^l}T_{i_l}(p_l)$$

Termination Sum on the root vector

back

Cardinal of Equivalence Classes Number of Tree Compatible Clustering

Cardinal of Equivalence Classes

Initialization For tips Propagation



Termination Sum on the root vector

Cardinal of Equivalence Classes Number of Tree Compatible Clustering

Cardinal of Equivalence Classes

Initialization For tips
Propagation

$$\mathcal{K}_{k}^{l} = \underset{1 \le p \le K}{\operatorname{argmin}} \{S_{i_{l}}(p) + \mathbb{I}\{p \ne k\}\}$$

$$S_{i}(k) = \sum_{l=1}^{L} S_{i_{l}}(p_{l}) + \mathbb{I}\{p_{l} \ne k\}, \quad \forall (p_{1}, \dots, p_{L}) \in \mathcal{K}_{k}^{1} \times \dots \times \mathcal{K}_{k}^{L} \xrightarrow{0} \{1, 0, 0\} \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad (1, 0, 0)$$

$$\int_{1 \le p \le K}^{K} (1, 0) \quad ($$

Termination Sum on the root vector

Cardinal of Equivalence Classes Number of Tree Compatible Clustering

Cardinal of Equivalence Classes

Initialization For tips Propagation

1

$$\mathcal{K}_{k}^{l} = \operatorname*{argmin}_{1 \leq p \leq K} \left\{ S_{i_{l}}(p) + \mathbb{I}\{p \neq k\} \right\}$$

$$S_i(k) = \sum_{l=1}^{L} S_{i_l}(p_l) + \mathbb{I}\{p_l \neq k\} , \ \forall (p_1, \dots p_L) \in \mathcal{K}_k^1 \times \dots \times \mathcal{K}_k^n$$

$$T_i(k) = \sum_{(p_1,\dots,p_L)\in\mathcal{K}_k^1\times\dots\times\mathcal{K}_k^L} \prod_{l=1}^L T_{i_l}(p_l) = \prod_{l=1}^L \sum_{p_l\in\mathcal{K}_k^l} T_{i_l}(p_l)$$

Termination Sum on the root vector



Cardinal of Equivalence Classes Number of Tree Compatible Clustering

Linking Shifts and Clustering

Assumption "No Homoplasy": 1 shift = 1 new color

Proposition "K shifts $\iff K + 1$ clusters"

Cardinal of Equivalence Classes Number of Tree Compatible Clustering

Linking Shifts and Clustering

Assumption "No Homoplasy": 1 shift = 1 new color



The No Homoplasy hypothesis is not respected.

Proposition "K shifts $\iff K + 1$ clusters"

Cardinal of Equivalence Classes Number of Tree Compatible Clustering

Linking Shifts and Clustering

Assumption "No Homoplasy": 1 shift = 1 new color



The No Homoplasy hypothesis is not respected.

Proposition "K shifts $\iff K + 1$ clusters"

Cardinal of Equivalence Classes Number of Tree Compatible Clustering

Definitions

- \mathcal{T} a rooted tree with *n* tips
- $N_{K}^{(T)} = |\mathcal{C}_{K}|$ the number of possible partitions of the tips in K clusters
- $A_{K}^{(\mathcal{T})}$ the number of possible *marked* partitions



Partitions in two groups for a binary tree with 3 tips

Difference between $N_2^{(\mathcal{T}_3)}$ and $A_2^{(\mathcal{T}_3)}$:

- $N_2^{(T_3)} = 3$: partitions 1 and 2 are equivalent
- A₂^(T₃) = 4: one marked color ("white = ancestral state")

Cardinal of Equivalence Classes Number of Tree Compatible Clustering

General Formula (Binary Case)

If \mathcal{T} is a binary tree, consider \mathcal{T}_{ℓ} and \mathcal{T}_{r} the left and right sub-trees of \mathcal{T} . Then:

$$\begin{cases} \mathsf{N}_{\mathsf{K}}^{(\mathcal{T})} = \sum_{k_1+k_2=\mathsf{K}} \mathsf{N}_{k_1}^{(\mathcal{T}_\ell)} \mathsf{N}_{k_2}^{(\mathcal{T}_r)} + \sum_{k_1+k_2=\mathsf{K}+1} \mathsf{A}_{k_1}^{(\mathcal{T}_\ell)} \mathsf{A}_{k_2}^{(\mathcal{T}_r)} \\ \mathsf{A}_{\mathsf{K}}^{(\mathcal{T})} = \sum_{k_1+k_2=\mathsf{K}} \mathsf{A}_{k_1}^{(\mathcal{T}_\ell)} \mathsf{N}_{k_2}^{(\mathcal{T}_r)} + \mathsf{N}_{k_1}^{(\mathcal{T}_\ell)} \mathsf{A}_{k_2}^{(\mathcal{T}_r)} + \sum_{k_1+k_2=\mathsf{K}+1} \mathsf{A}_{k_1}^{(\mathcal{T}_\ell)} \mathsf{A}_{k_2}^{(\mathcal{T}_r)} \end{cases}$$

We get:

$$\mathcal{N}_{K+1}^{(\mathcal{T})} = \mathcal{N}_{K+1}^{(n)} = egin{pmatrix} 2n-2-K\ K \end{pmatrix}$$
 and $\mathcal{A}_{K+1}^{(\mathcal{T})} = \mathcal{A}_{K+1}^{(n)} = egin{pmatrix} 2n-1-K\ K \end{pmatrix}$

Cardinal of Equivalence Classes Number of Tree Compatible Clustering

Recursion Formula (General Case)

If we are at a node defining a tree T that has p daughters, with sub-trees T_1, \ldots, T_p , then we get the following recursion formulas:

$$\begin{cases} \mathsf{N}_{\mathsf{K}}^{(\mathcal{T})} = \sum_{\substack{k_1 + \dots + k_p = \mathsf{K} \\ k_1, \dots, k_p \ge 1}} \prod_{i=1}^{p} \mathsf{N}_{k_i}^{(\mathcal{T}_i)} + \sum_{\substack{l \subset \llbracket 1, p \rrbracket \\ |l| \ge 2}} \sum_{\substack{k_1 + \dots + k_p = \mathsf{K} + |l| - 1 \\ k_1, \dots, k_p \ge 1}} \prod_{i \in I} \mathsf{A}_{k_i}^{(\mathcal{T}_i)} \prod_{i \notin I} \mathsf{N}_{k_i}^{(\mathcal{T}_i)} \\ \mathsf{A}_{\mathsf{K}}^{(\mathcal{T})} = \sum_{\substack{l \subset \llbracket 1, p \rrbracket \\ |l| \ge 1}} \sum_{\substack{k_1 + \dots + k_p = \mathsf{K} + |l| - 1 \\ k_1, \dots, k_p \ge 1}} \prod_{i \in I} \mathsf{A}_{k_i}^{(\mathcal{T}_i)} \prod_{i \notin I} \mathsf{N}_{k_i}^{(\mathcal{T}_i)} \end{cases}$$

No general formula. The result depends on the topology of the tree.

back

Simulations Design

- Topology of the tree fixed (unit height, $\lambda = 0.1$, with 64, 128, 256 taxa).
- Initial optimal value fixed: $\beta_0 = 0$
- One "base" scenario $\alpha_b = 3$, $\gamma_b^2 = 0.5$, $K_b = 5$.
- $\alpha \in \log(2)/\{0.01, 0.05, 0.1, 0.2, 0.23, 0.3, 0.5, 0.75, 1, 2, 10\}.$
- $\gamma^2 \in \{0.3, 0.6, 3, 6, 12, 18, 30, 60, 150\}/(2\alpha_b).$
- $K \in \{0, 1, 2, 3, 4, 5, 8, 11, 16\}.$
- Shifts values $\sim rac{1}{2}\mathcal{N}(4,1) + rac{1}{2}\mathcal{N}(-4,1)$
- Shifts randomly placed at regular intervals separated by 0.1 unit length.
- n = 200 repetitions : 16200 configurations.

CPU time on cluster MIGALE (Jouy-en-Josas):

- α known: 6 minutes per estimation (66 days in total).
- α unknown: 52 minutes per estimation (570 days in total).

Log-Likelihood



Log likelihood for a tree with 256 tips. Solid black dots are the median of the log likelihood for the true parameters.

Number of Shifts



CA, PB, MM, SR Change-point Detection on a Tree

One Example



CA, PB, MM, SR Change-point Detection on a Tree

Adjusted Rand Index



CA, PB, MM, SR

Parameters: β_0



CA, PB, MM, SR Change-point Detection on a Tree

Parameters: α



CA, PB, MM, SR Change-point Detection on a Tree

Parameters: γ^2



CA, PB, MM, SR Chan

Models Inference

BM Model

Data *n* vectors of *p* traits at the tips:
$$\mathbf{Y}_{i} = \begin{pmatrix} Y_{i1} \\ \vdots \\ Y_{ip} \end{pmatrix}$$

SDE $d\mathbf{W}(t) = \mathbf{\Sigma} d\mathbf{B}_{t}$, rate matrix $\mathbf{R} = \mathbf{\Sigma} \mathbf{\Sigma}^{T} (p \times p)$

Covariances \mathbb{C} ov $[Y_{il}; Y_{jq}] = t_{ij}R_{lq}$ for i, j tips, and l, q characters

$$\mathbb{V}$$
ar [vec(\mathbf{Y})] = $\mathbf{C}_n \otimes \mathbf{R}$

Models Inference

BM Model

Linear Model Representation

$$\operatorname{vec}(\mathbf{Y}) = \operatorname{vec}(\mathbf{\Delta T}^{T}) + \mathbf{E}$$
 with $\mathbf{E} \sim \mathcal{N}(0, \mathbf{V} = \mathbf{C}_{n} \otimes \mathbf{R})$

Incomplete Data Representation

$$\mathbf{Y}_3 \mid \mathbf{Z}_2 \sim \mathcal{N} \Big(\mathbf{Z}_2 + oldsymbol{\delta}, \ \ell_7 \mathbf{R} \Big)$$



Models Inference

OU Model: General Case

Data *n* vectors of *p* traits at the tips:
$$\mathbf{Y}_i = \begin{pmatrix} Y_{i1} \\ \vdots \\ Y_{ip} \end{pmatrix}$$

SDE **A** $(p \times p)$ "selection strength"

$$d\mathbf{W}(t) = -\mathbf{A}(\mathbf{W}(t) - \beta(t))dt + \mathbf{\Sigma}d\mathbf{B}_t$$

Covariances

$$\mathbb{C} \text{ov} \left[\mathbf{X}_{i}; \mathbf{X}_{j} \right] = e^{-\mathbf{A}t_{i}} \mathbf{\Gamma} e^{-\mathbf{A}^{T}t_{j}} \\ + e^{-\mathbf{A}(t_{i}-t_{ij})} \left(\int_{0}^{t_{ij}} e^{-\mathbf{A}v} \mathbf{\Sigma} \mathbf{\Sigma}^{T} e^{-\mathbf{A}^{T}v} dv \right) e^{-\mathbf{A}^{T}(t_{j}-t_{ij})}$$
Shifts K shifts $\delta_{1}, \cdots, \delta_{K}$ vectors size p

 \mapsto On the optimal values

Models Inference

OU Model: A scalar

Assumption
$$\mathbf{A} = \alpha \mathbf{I}_{p}$$
 "scalar"

Stationnary State $\mathbf{S} = \frac{1}{2\alpha} \mathbf{R}$

Fixed Root For i, j tips and l, q characters:

$$\mathbb{C}\mathrm{ov}\left[Y_{il};Y_{jq}\right] = \frac{1}{2\alpha} e^{-2\alpha h} \left(e^{2\alpha t_{ij}} - 1\right) R_{lq}$$

 $\mapsto\,$ Can be reduced to a BM on a re-scaled tree
Models Inference

EM algorithm

BM Natural generalization of the univariate case.

OU M step intractable in general.

Incomplete Data Model: Can readily handle missing data.

Models Inference

Model Selection

- Previous criterion cannot be applied
- Solution: "Slope Heuristic"-based method
 - Massart (2007)
 - oracle inequality with known variance
 - penalty up to a multiplicative constant
 - Baudry et al. (2012)
 - Slope-heuristic method to calibrate the constant
 - Implemented in capushe (Brault et al., 2012)

Models Inference

Model Selection: Toy Example



Figure: Simulated Process.

Models Inference

Model Selection: Toy Example



Figure: capushe output for penalized log-likelihood.

Models Inference

Model Selection: Toy Example



Figure: Reconstructed Process.