

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

Cécile Ané<sup>1,2</sup>, Paul Bastide<sup>3,4</sup>, Mahendra Mariadassou<sup>4</sup>,  
Stéphane Robin<sup>3</sup>

<sup>1</sup> Department of Statistics, University of Wisconsin–Madison, WI, 53706, USA

<sup>2</sup> Department of Botany, University of Wisconsin–Madison, WI, 53706, USA

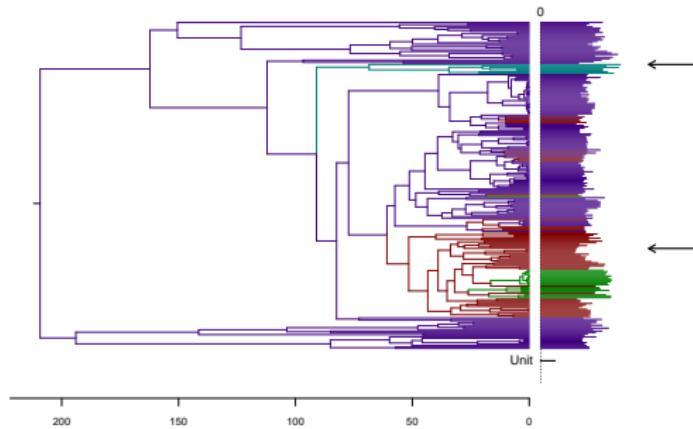
<sup>3</sup> UMR MIA-Paris, AgroParisTech, INRA, Université Paris-Saclay, 75005, Paris, France

<sup>4</sup> MaIAGE, INRA, Université Paris-Saclay, 78352 Jouy-en-Josas, France

23 June 2016



# Introduction



*Dermochelys Coriacea*



*Homopus Areolatus*

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

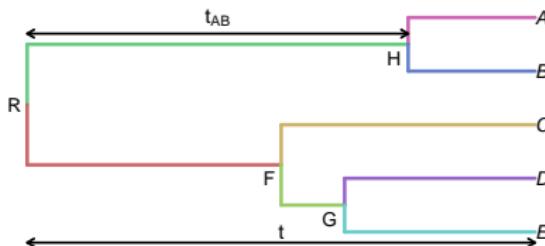
- How can we explain the diversity, while accounting for the phylogenetic correlations ?
- Modelling: a shifted stochastic process on the phylogeny.

# Outline

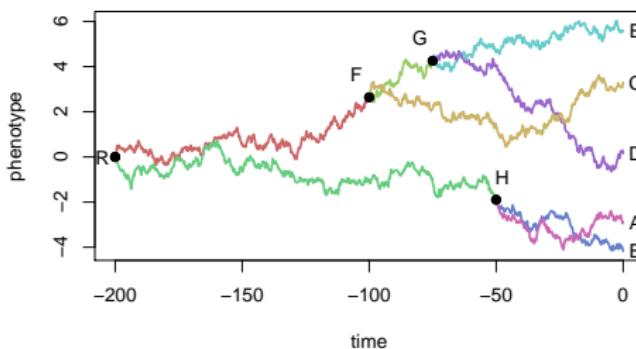
- ① Stochastic Processes on Trees
- ② Identifiability Problems and Counting Issues
- ③ Statistical Inference
- ④ Turtles Data Set

# Stochastic Process on a Tree

(Felsenstein, 1985)



The tree is known.  
 Only *tip* trait values  
 are observed.

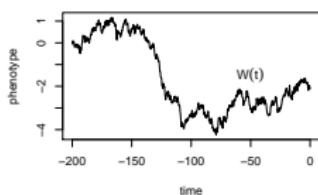
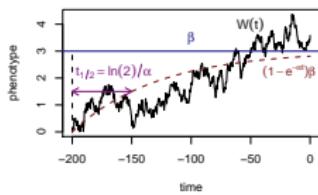


Brownian Motion:

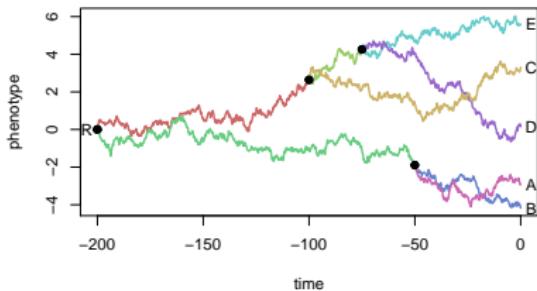
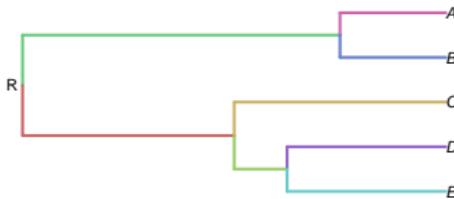
$$\text{Var}[A | R] = \sigma^2 t$$

$$\text{Cov}[A; B | R] = \sigma^2 t_{AB}$$

# Brownian Motion vs Ornstein-Uhlenbeck

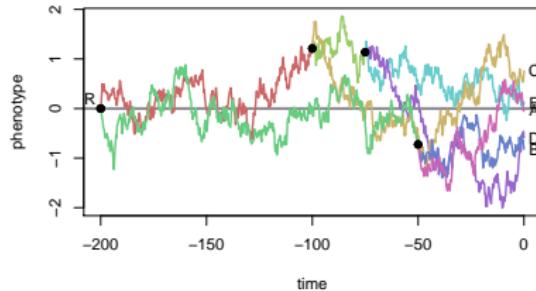
Equation	Stationary State	Variance
 <p><math>W(t)</math></p>	$dW(t) = \sigma dB(t)$	None. $\sigma_{ij} = \sigma^2 t_{ij}$
 <p><math>W(t)</math></p> <p><math>\beta</math></p> <p><math>t_1/t_2 = \ln(2)/\alpha</math></p> <p><math>(1-e^{-\alpha t})\beta</math></p>	$dW(t) = \sigma dB(t) + \alpha[\beta - W(t)]dt$	$\begin{cases} \mu = \beta_0 \\ \gamma^2 = \frac{\sigma^2}{2\alpha} \end{cases} \quad \sigma_{ij} = \gamma^2 e^{-\alpha(t_i+t_j)} \times (e^{2\alpha t_{ij}} - 1)$

# Shifts



**BM Shifts in the mean:**

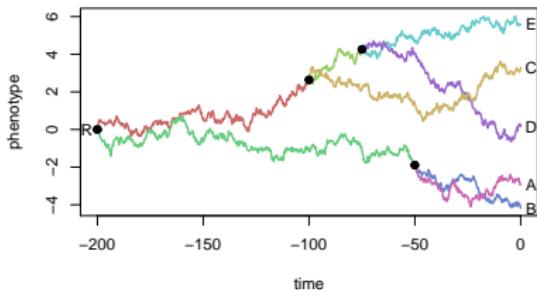
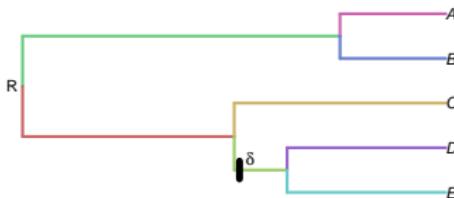
$$m_{\text{child}} = m_{\text{parent}} + \delta$$



**OU Shifts in the optimal value:**

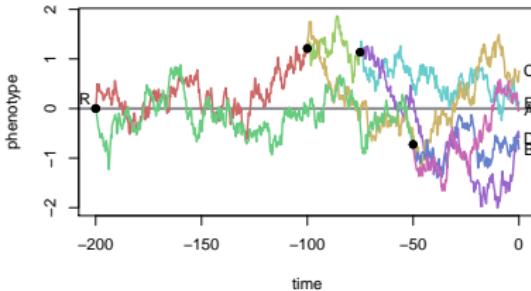
$$\beta_{\text{child}} = \beta_{\text{parent}} + \delta$$

# Shifts



**BM Shifts in the mean:**

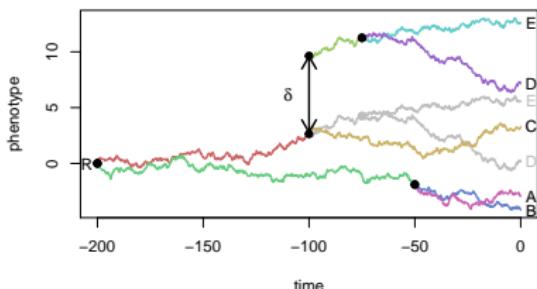
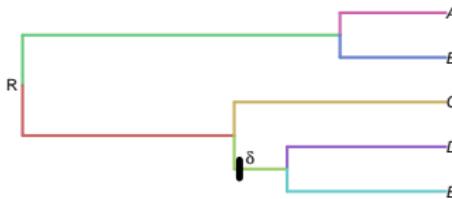
$$m_{\text{child}} = m_{\text{parent}} + \delta$$



**OU Shifts in the optimal value:**

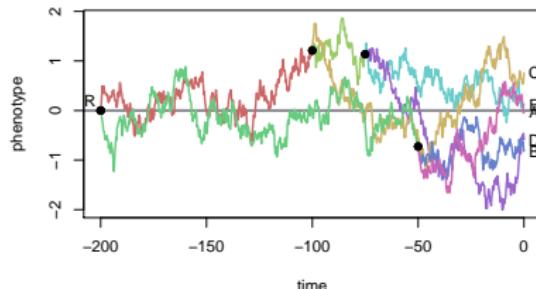
$$\beta_{\text{child}} = \beta_{\text{parent}} + \delta$$

# Shifts



**BM Shifts in the mean:**

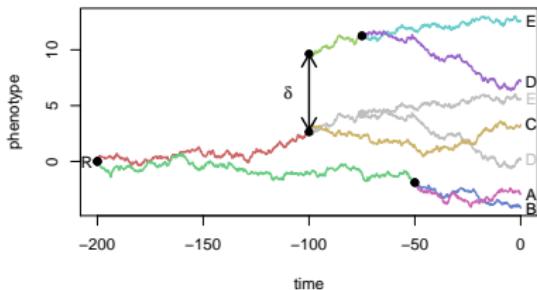
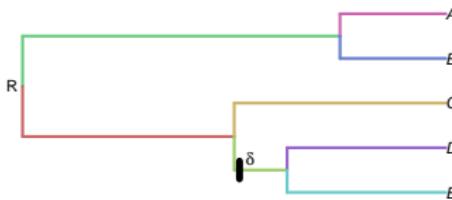
$$m_{\text{child}} = m_{\text{parent}} + \delta$$



**OU Shifts in the optimal value:**

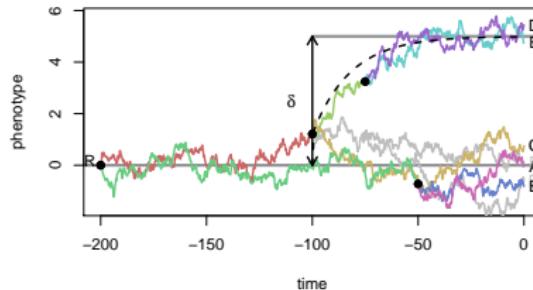
$$\beta_{\text{child}} = \beta_{\text{parent}} + \delta$$

# Shifts



**BM Shifts in the mean:**

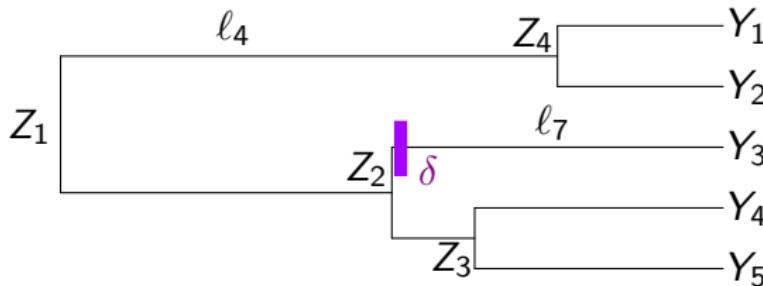
$$m_{\text{child}} = m_{\text{parent}} + \delta$$



**OU Shifts in the optimal value:**

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

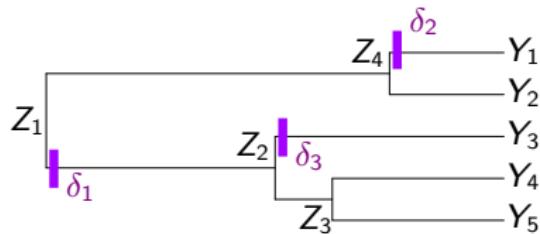
# Incomplete Data Model



$$\begin{aligned}
 BM \quad Z_4 | Z_1 &\sim \mathcal{N}\left(Z_1, \sigma^2 \ell_4\right) \\
 Y_3 | Z_2 &\sim \mathcal{N}\left(Z_2 + \delta, \sigma^2 \ell_7\right)
 \end{aligned}$$

$$OU \quad Y_3 | Z_2 \sim \mathcal{N}\left(Z_2 e^{-\alpha \ell_7} + (1 - e^{-\alpha \ell_7})(\beta_{Z_2} + \delta), \frac{\sigma^2}{2\alpha}(1 - e^{-2\alpha \ell_7})\right)$$

# Linear Regression Model



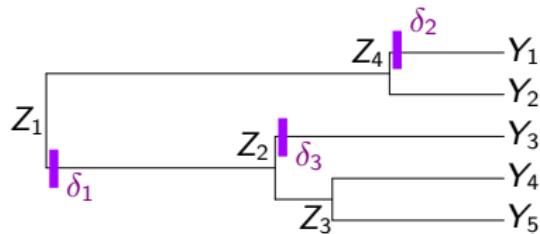
$$\Delta = \begin{pmatrix} \mu \\ \delta_1 \\ 0 \\ 0 \\ \delta_2 \\ 0 \\ \delta_3 \\ 0 \\ 0 \end{pmatrix}$$

$$T\Delta = \begin{pmatrix} \mu + \delta_2 \\ \mu \\ \mu + \delta_1 + \delta_3 \\ \mu + \delta_1 \\ \mu + \delta_1 \end{pmatrix}$$

$$T = \begin{matrix} & Z_1 & Z_2 & Z_3 & Z_4 & Y_1 & Y_2 & Y_3 & Y_4 & Y_5 \\ Y_1 & \left( \begin{array}{ccccccccc} 1 & 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 \end{array} \right) \\ Y_2 & \left( \begin{array}{ccccccccc} 1 & 0 & 0 & 1 & 0 & 1 & 0 & 0 & 0 \end{array} \right) \\ Y_3 & \left( \begin{array}{ccccccccc} 1 & 1 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \end{array} \right) \\ Y_4 & \left( \begin{array}{ccccccccc} 1 & 1 & 1 & 0 & 0 & 0 & 0 & 1 & 0 \end{array} \right) \\ Y_5 & \left( \begin{array}{ccccccccc} 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 1 \end{array} \right) \end{matrix}$$

$$BM : \quad Y = T\Delta^{BM} + E^{BM}$$

# Linear Regression Model



$$\Delta = \begin{pmatrix} \lambda \\ \delta_1 \\ 0 \\ 0 \\ \delta_2 \\ 0 \\ \delta_3 \\ 0 \\ 0 \end{pmatrix} \quad TW(\alpha)\Delta = \begin{pmatrix} \lambda + w_5\delta_2 \\ \lambda \\ \lambda + w_2\delta_1 + w_7\delta_3 \\ \lambda + w_2\delta_1 \\ \lambda + w_2\delta_1 \end{pmatrix}$$

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

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

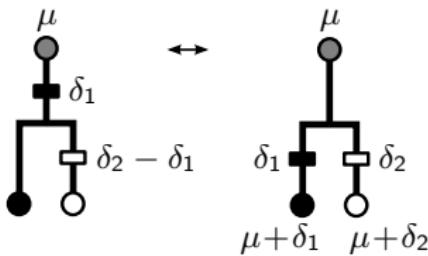
$$BM : \quad Y = T\Delta^{BM} + E^{BM}$$

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



# Equivalencies

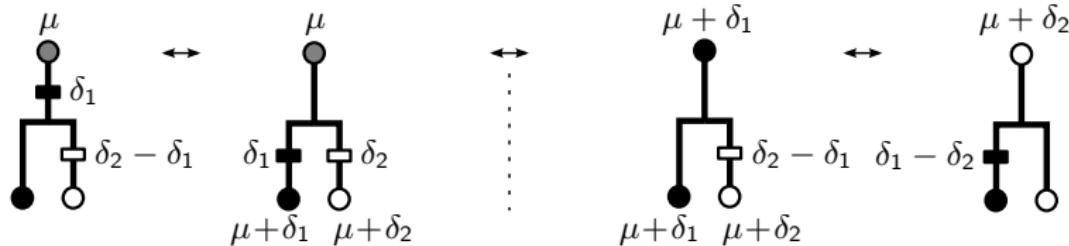
- Number of shifts  $K$  fixed, several equivalent solutions.



- Problem of over-parametrization: parsimonious configurations.

# Equivalencies

- Number of shifts  $K$  fixed, several equivalent solutions.



- Problem of over-parametrization: parsimonious configurations.

# Parsimonious Solution : Definition

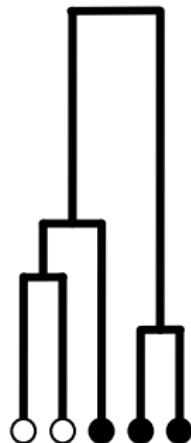
## Definition (Parsimonious Allocation)

A coloring of the tips being given, a *parsimonious* allocation of the shifts is such that it has a minimum number of shifts.

# Parsimonious Solution : Definition

## Definition (Parsimonious Allocation)

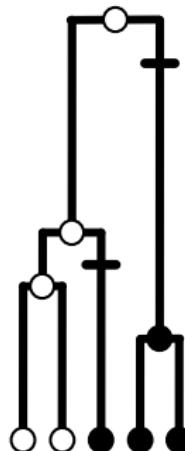
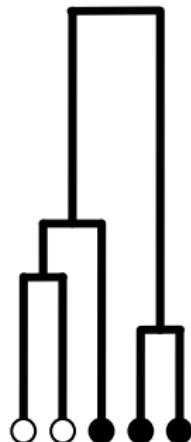
A coloring of the tips being given, a *parsimonious* allocation of the shifts is such that it has a minimum number of shifts.



# Parsimonious Solution : Definition

## Definition (Parsimonious Allocation)

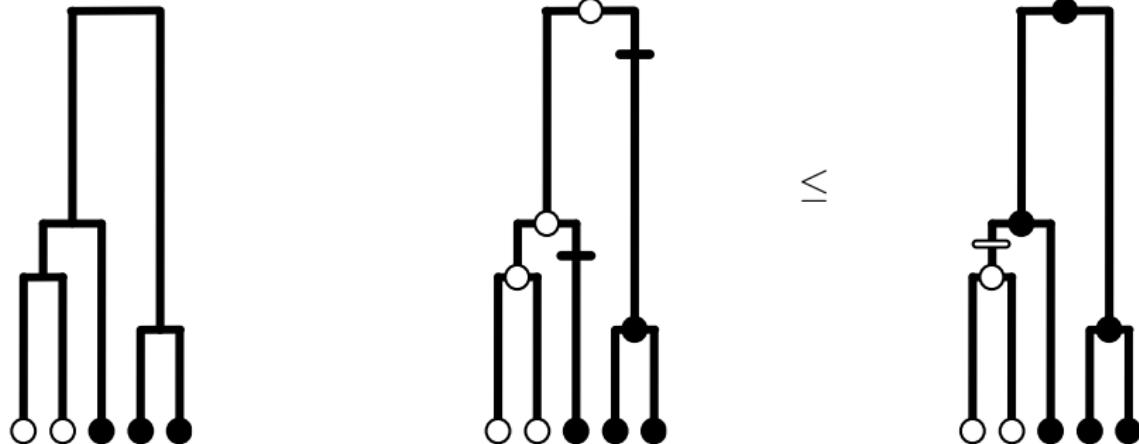
A coloring of the tips being given, a *parsimonious* allocation of the shifts is such that it has a minimum number of shifts.



# Parsimonious Solution : Definition

## Definition (Parsimonious Allocation)

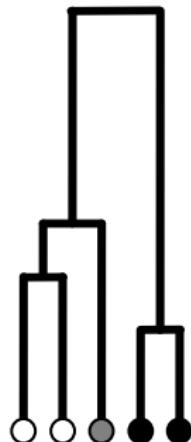
A coloring of the tips being given, a *parsimonious* allocation of the shifts is such that it has a minimum number of shifts.



# Parsimonious Solution : Definition

## Definition (Parsimonious Allocation)

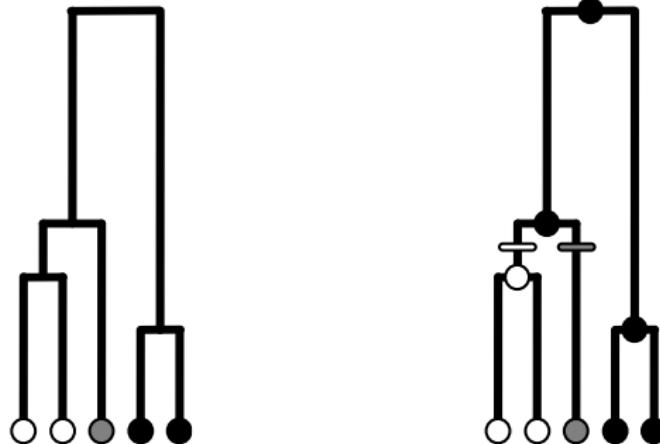
A coloring of the tips being given, a *parsimonious* allocation of the shifts is such that it has a minimum number of shifts.



# Parsimonious Solution : Definition

## Definition (Parsimonious Allocation)

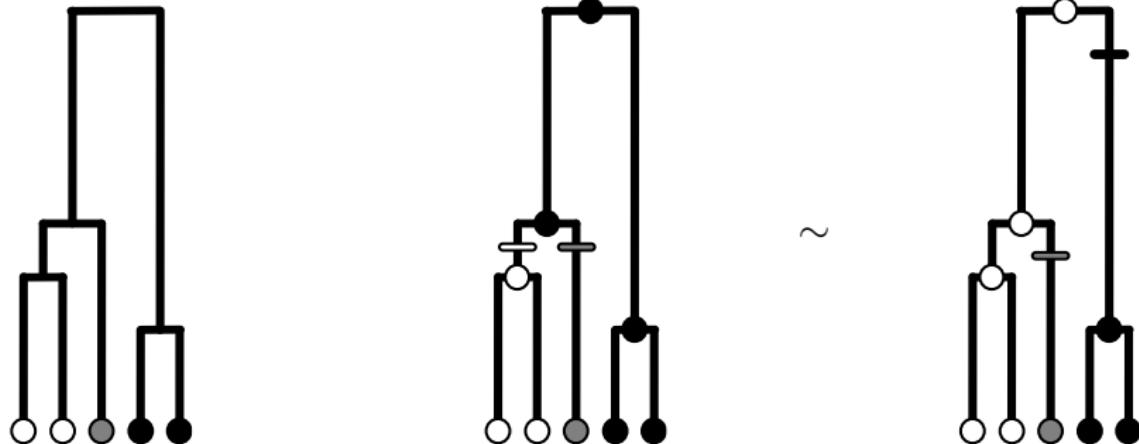
A coloring of the tips being given, a *parsimonious* allocation of the shifts is such that it has a minimum number of shifts.



# Parsimonious Solution : Definition

## Definition (Parsimonious Allocation)

A coloring of the tips being given, a *parsimonious* allocation of the shifts is such that it has a minimum number of shifts.



# Equivalent Parsimonious Allocations

## Definition (Equivalency)

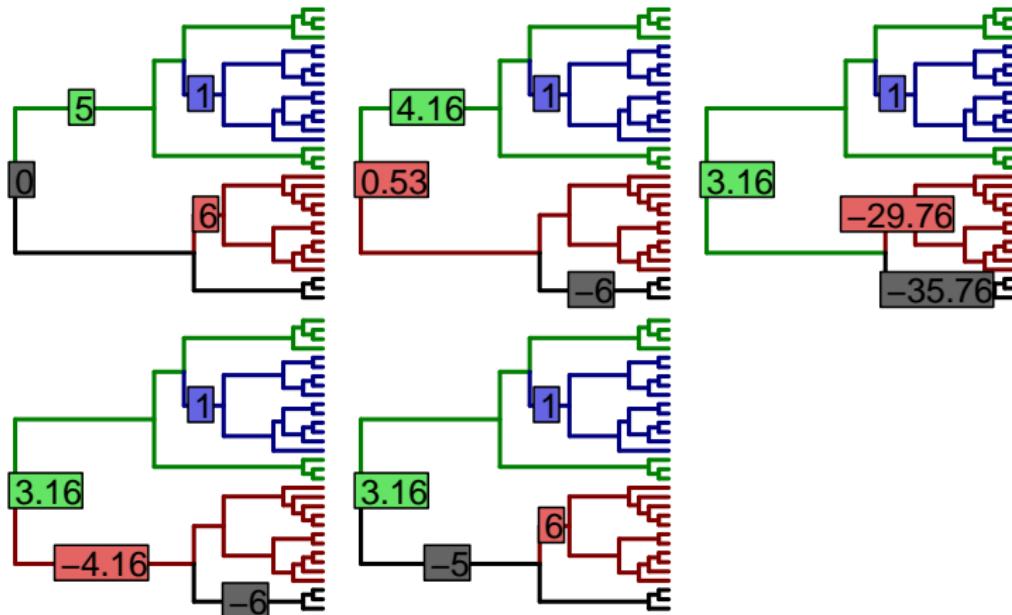
Two allocations are said to be *equivalent* (noted  $\sim$ ) if they are both parsimonious and give the same colors at the tips.

**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).



## Equivalent Parsimonious Solutions for an OU Model.



*Equivalent allocations and values of the shifts - OU.*

# Collection of Models

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

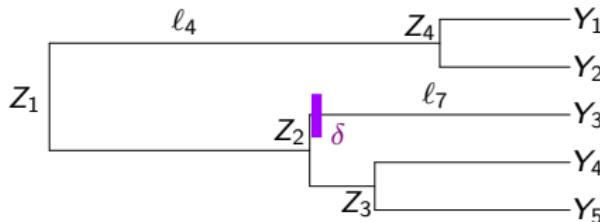
- $|\mathcal{S}_K^{PI}| \leq \binom{m+n-1}{K} = \binom{\text{\# of edges}}{\text{\# of shifts}}$
- A recursive algorithm to compute  $|\mathcal{S}_K^{PI}|$  (implemented in R).

→ 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: number of shifts K fixed



$$\text{BM } Z_4|Z_1 \sim \mathcal{N}\left(Z_1, \sigma^2 \ell_4\right)$$

$$Y_3|Z_2 \sim \mathcal{N}\left(Z_2 + \delta, \sigma^2 \ell_7\right)$$

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

**EM Algorithm**  $\log p_{\theta}(Y) = \mathbb{E}_{\theta}[\log p_{\theta}(Z, Y) | Y] - \mathbb{E}_{\theta}[\log p_{\theta}(Z) | Y]$

**E step** Given  $\theta^h$ , compute  $p_{\theta^h}(Z | Y)$

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

**Initialization** Lasso



# Model Selection on $K$

Assumption  $\alpha$  fixed

$$Y = TW(\alpha)\Delta + \gamma E \quad , \quad E \sim \mathcal{N}(0, V(\alpha))$$

## Models

$\eta \in \bigcup_{K=0}^{p-1} \mathcal{S}_K^{PI}$ : Identifiable parsimonious allocations of shifts

## EM Estimators

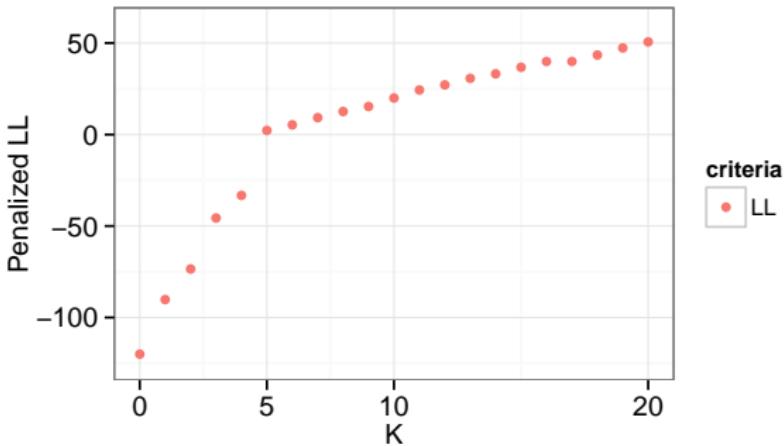
$$\hat{Y}_K = \operatorname{argmin}_{\eta \in \mathcal{S}_K^{PI}} \| Y - \hat{Y}_\eta \|_V^2$$

## Oracle

$$\inf_{\eta \in \bigcup_{K=0}^{p-1} \mathcal{S}_K^{PI}} \| \mathbb{E}[Y] - Y_\eta^* \|_V^2$$

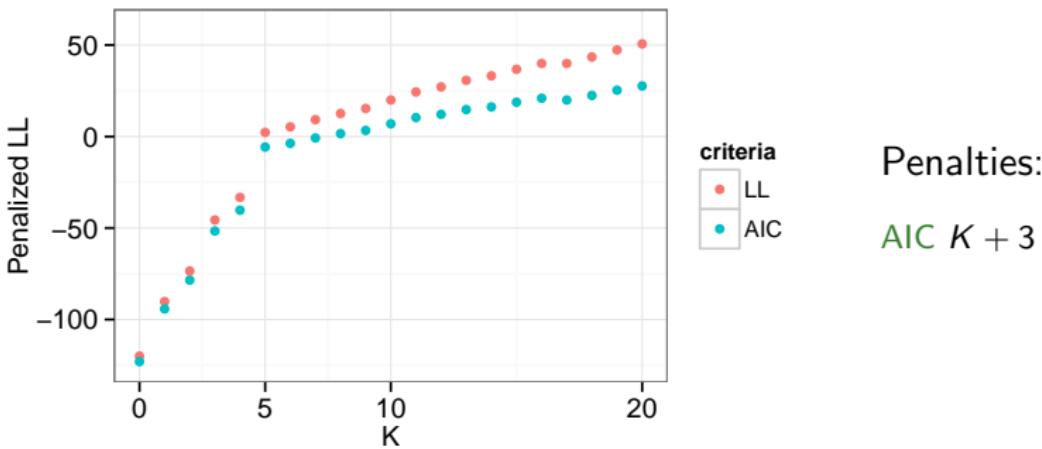
# Model Selection: Penalized Likelihood

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



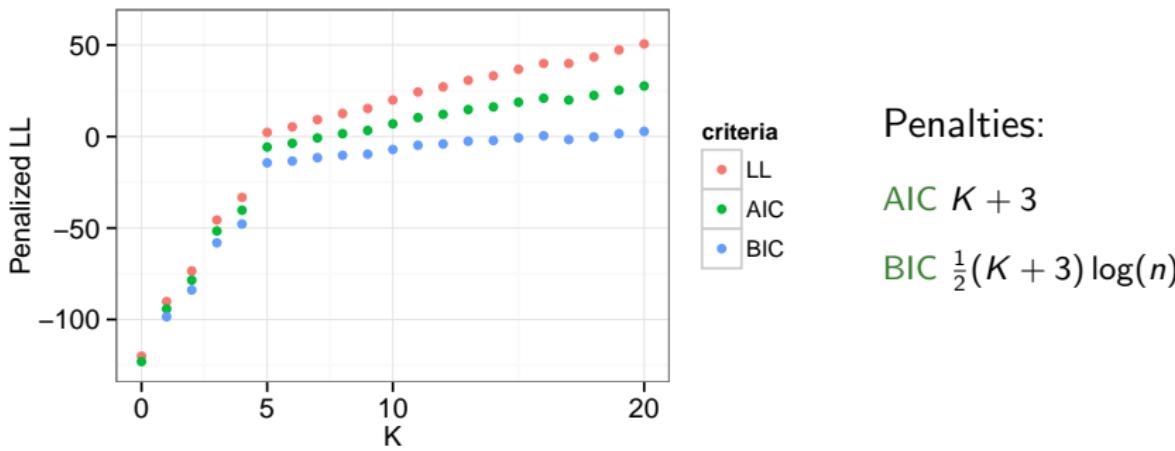
# Model Selection: Penalized Likelihood

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



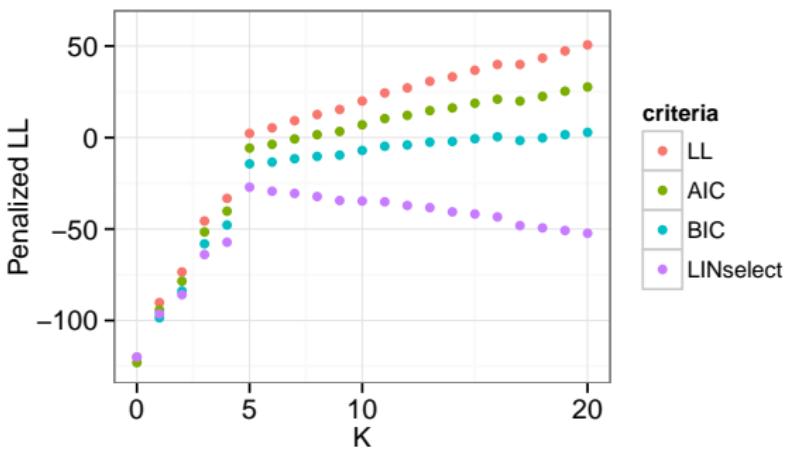
# Model Selection: Penalized Likelihood

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



# Model Selection: Penalized Likelihood

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



Penalties:

- AIC  $K + 3$
- BIC  $\frac{1}{2}(K + 3) \log(n)$
- LINselect  $\operatorname{pen}(n, K, |\mathcal{S}_K^{PI}|)$

# Proposition: LINselect Penalty

Proposition (Form of the Penalty and guarantees ( $\alpha$  known))

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

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

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

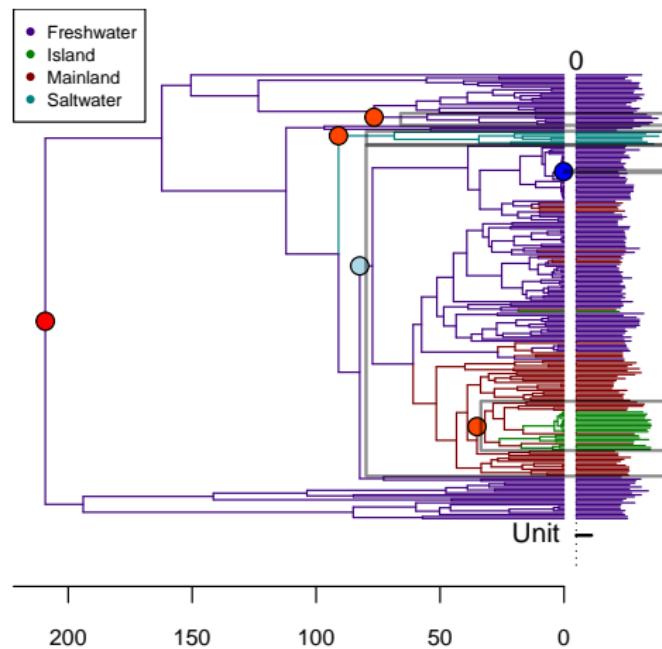
$$\mathbb{E} \left[ \frac{\|\mathbb{E}[Y] - \hat{Y}_{\hat{K}}\|_V^2}{\gamma^2} \right] \leq C(A, \kappa) \inf_{\eta \in \mathcal{M}} \left\{ \frac{\|\mathbb{E}[Y] - Y_{\eta}^*\|_V^2}{\gamma^2} + (K_{\eta} + 2)(3 + \log(n)) \right\}$$

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

Based on Baraud et al. (2009)



# Turtles Dataset

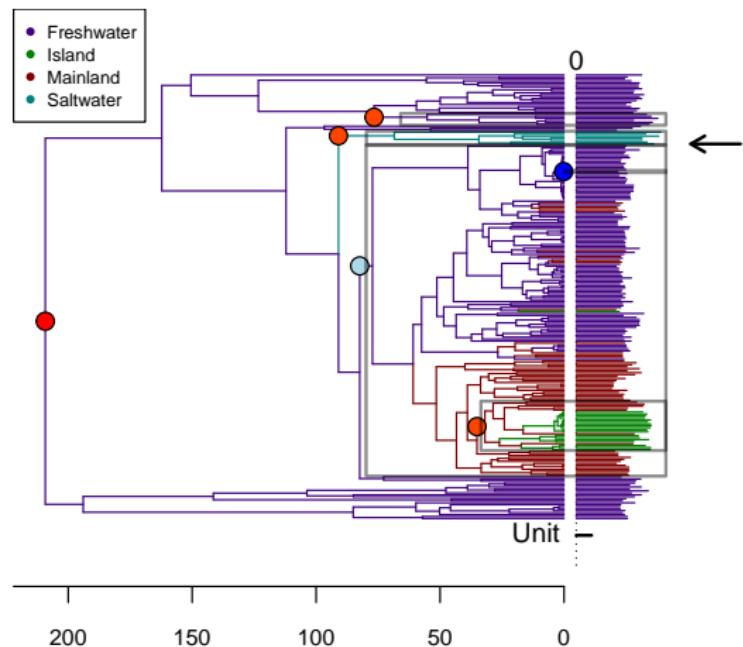


Colors: habitats.  
 Boxes: selected EM regimes.

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

(Jaffe et al., 2011)

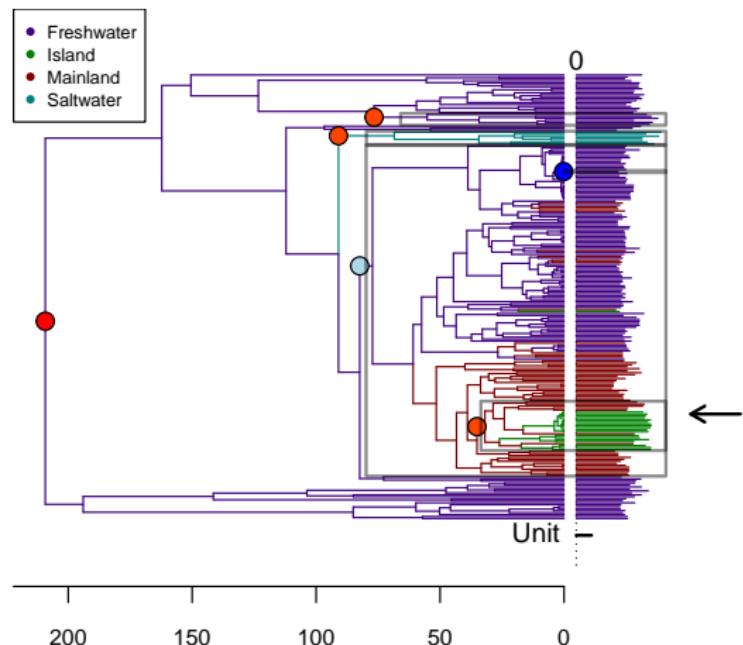
# Turtles Dataset



*Chelonia mydas*

Colors: habitats.  
Boxes: selected EM regimes.

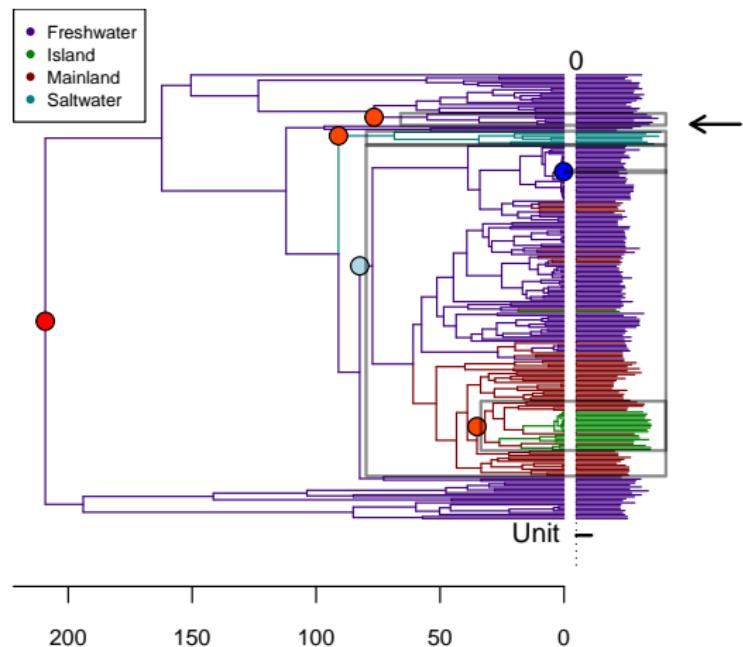
# Turtles Dataset



*Geochelone nigra abingdoni*

Colors: habitats.  
Boxes: selected EM regimes.

# Turtles Dataset



*Chitra indica*

Colors: habitats.  
Boxes: selected EM regimes.

# Conclusion and Perspectives

A general inference framework for trait evolution models.

## Conclusions

- Identifiability can be assessed.
- An EM can be written to maximize likelihood.
- Model selection for a non-iid framework.

R codes Available on GitHub:

<https://github.com/pbastide/PhylogeneticEM>

## Perspectives

- Multivariate traits.
- Deal with uncertainty (tree, data).
- Phylogenetic networks.

# Bibliography

- 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, Sunderland, 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, USVI Uploaded 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
  - "Dudhwali 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



[pbastide.github.io](https://pbastide.github.io)

# Appendices

## ⑤ Inference

- EM
- Model Selection

## ⑥ Identifiability Issues

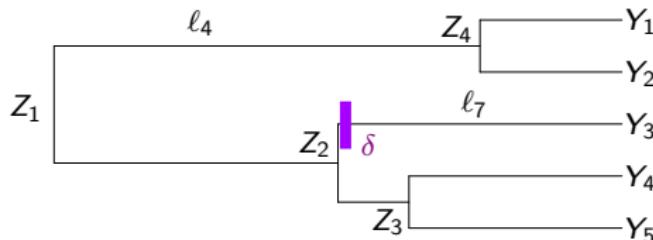
- OU  $\iff$  BM
- Cardinal of Equivalence Classes
- Number of Tree Compatible Clustering

## ⑦ Simulations Results

## ⑧ Multivariate

- Models
- Inference

# E step



Compute the following quantities:

$$\mathbb{E}^{(h)}[Z_j | Y], \text{Var}^{(h)}[Z_j | Y], \text{Cov}^{(h)}[Z_j, Z_{\text{parent}(j)} | Y]$$

- Using Gaussian properties. Need to invert matrices: complexity in  $O(n^3)$ .
- Using Gaussian properties **and** the tree structure: "Upward-Downward" algorithm. Complexity in  $O(n)$ .



# M Step

Maximize:

$$\mathbb{E} [\log p_\theta(X) \mid Y] = - \sum_{j=2}^{m+n} C_j(\alpha, \text{shifts}) + \mathcal{F}^{(h)} (\mu, \gamma^2, \sigma^2, \alpha)$$

- $\mu, \gamma^2, \sigma^2$ : simple maximization
- Discrete location of  $K$  shifts
  - ↳ Exact and fast for the BM
- $\alpha$ : numerical maximization and/or on a grid
  - ↳ Generalized EM



# Initialization

Shifts : Lasso regression.

$$\hat{\Delta} = \operatorname{argmin}_{\Delta} \left\{ \|Y - TW(\alpha)\Delta\|_{\Sigma_{YY}^{-1}}^2 + \lambda \|\Delta\|_1 \right\}$$

- Initialize  $\Sigma_{YY}(\alpha)$ , then estimate  $\Delta$  with a Gauss Lasso procedure, using a Cholesky decomposition.
- $\lambda$  chosen to get  $K$  shifts.



The selection strength  $\alpha$  : Initialization using couples of tips.

back

# Cholesky Decomposition

The problem is:

$$\hat{\Delta} = \operatorname{argmin}_{\Delta} \left\{ \|Y - R\Delta\|_{\Sigma_{YY}}^2 + \lambda |\Delta|_1 \right\}$$

Cholesky decomposition of  $\Sigma_{YY}$ :

$$\Sigma_{YY} = LL^T, \text{ } L \text{ a lower triangular matrix}$$

Then:

$$\|Y - R\Delta\|_{\Sigma_{YY}}^2 = \|L^{-1}Y - L^{-1}R\Delta\|^2$$

And if  $Y' = L^{-1}Y$  and  $R' = L^{-1}R$ , the problem becomes:

$$\hat{\Delta} = \operatorname{argmin}_{\Delta} \left\{ \|Y' - R'\Delta\|^2 + \lambda |\Delta|_1 \right\}$$

# Gauss Lasso

Let  $\hat{m}_\lambda$  be the set of selected variables (including the root). Then:

$$\hat{\Delta}^{\text{Gauss}} = \Pi_{\hat{F}_\lambda}(Y') \text{ with } \hat{F}_\lambda = \text{Span}\{R'_j : j \in \hat{m}_\lambda\}$$

back

# Goal and Notations

**Data** A process on a tree with the following structure:

$$\forall j > 1, \quad X_j | X_{\text{pa}(j)} \sim \mathcal{N} (m_j(X_{\text{pa}(j)}) = q_j X_{\text{pa}(j)} + r_j, \sigma_j^2)$$

$$\text{BM: } \begin{cases} q_j = 1 \\ r_j = \sum_k \mathbb{I}\{\tau_k = b_j\} \delta_k \\ \sigma_j^2 = \ell_j \sigma^2 \end{cases} \quad \text{OU: } \begin{cases} q_j = e^{-\alpha \ell_j} \\ r_j = \beta^{\text{pa}(j)} (1 - e^{-\alpha \ell_j}) + \sum_k \mathbb{I}\{\tau_k = b_j\} \delta_k (1 - e^{-\alpha(1-\nu_k) \ell_j}) \\ \sigma_j^2 = \frac{\sigma^2}{2\alpha} (1 - e^{-2\alpha \ell_j}) \end{cases}$$

**Goal** Compute the following quantities, at every node  $j$ :

$$\mathbb{V}\text{ar}^{(h)}[Z_j | Y], \mathbb{C}\text{ov}^{(h)}[Z_j, Z_{\text{pa}(j)} | Y], \mathbb{E}^{(h)}[Z_j | Y]$$

# Upward

**Goal** Compute for a vector of tips, given their common ancestor:

$$f_{\mathbf{Y}^j|X_j}(\mathbf{Y}^j; a) = A_j(\mathbf{Y}^j)\Phi_{M_j(\mathbf{Y}^j), S_j^2(\mathbf{Y}^j)}(a)$$

**Initialization** For tips:  $f_{Y_i|Y_i}(Y_i; a) = \Phi_{Y_i, 0}(a)$

**Propagation**

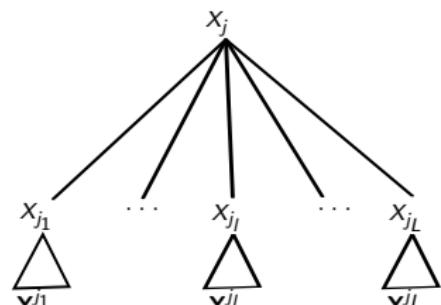
$$f_{\mathbf{Y}^j|X_j}(\mathbf{Y}^j; a) = \prod_{l=1}^L f_{\mathbf{Y}^{j_l}|X_j}(\mathbf{Y}^{j_l}; a)$$

$$f_{\mathbf{Y}^{j_l}|X_j}(\mathbf{Y}^{j_l}; a) = \int_{\mathbb{R}} f_{\mathbf{Y}^{j_l}|X_{j_l}}(\mathbf{Y}^{j_l}; b) f_{X_{j_l}|X_j}(b; a) db$$

**Root Node and Likelihood** At the root:

$$f_{X_1|\mathbf{Y}}(a; \mathbf{Y}) \propto f_{\mathbf{Y}|X_1}(\mathbf{Y}; a) f_{X_1}(a)$$

$$\left\{ \begin{array}{l} \text{Var}[X_1 | \mathbf{Y}] = \left( \frac{1}{\gamma^2} + \frac{1}{S_1^2(\mathbf{Y})} \right)^{-1} \\ \mathbb{E}[X_1 | \mathbf{Y}] = \text{Var}[X_1 | \mathbf{Y}] \left( \frac{\mu}{\gamma^2} + \frac{M_1(\mathbf{Y})}{S_1^2(\mathbf{Y})} \right) \end{array} \right.$$



# Downward

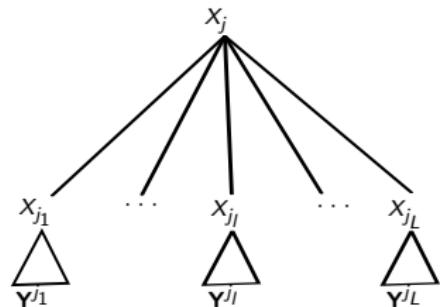
Compute  $E_j = \mathbb{E} [X_j | \mathbf{Y}]$  ,  $V_j^2 = \text{Var} [X_j | \mathbf{Y}]$  ,  $C_{j,\text{pa}(j)}^2 = \text{Cov} [X_j; X_{\text{pa}(j)} | \mathbf{Y}]$

**Initialization** Last step of Upward.

**Propagation**

$$f_{X_{\text{pa}(j)}, X_j | \mathbf{Y}}(a, b; \mathbf{Y}) = f_{X_{\text{pa}(j)} | \mathbf{Y}}(a; \mathbf{Y}) f_{X_j | X_{\text{pa}(j)}, \mathbf{Y}}(b; a, \mathbf{Y})$$

$$\begin{aligned} f_{X_j | X_{\text{pa}(j)}, \mathbf{Y}}(b; a, \mathbf{Y}) &= f_{X_j | X_{\text{pa}(j)}, \mathbf{Y}^j}(b; a, \mathbf{Y}^j) \\ &\propto f_{X_j | X_{\text{pa}(j)}}(b; a) f_{\mathbf{Y}^j | X_j}(\mathbf{Y}^j; b) \end{aligned}$$



# Formulas

Upward

$$\begin{cases} S_j^2(\mathbf{Y}^j) = \left( \sum_{l=1}^L \frac{q_{jl}^2}{S_{jl}^2(\mathbf{Y}^{ji}) + \sigma_{jl}^2} \right)^{-1} \\ M_j(\mathbf{Y}^j) = S_j^2(\mathbf{Y}^j) \sum_{l=1}^L q_{jl} \frac{M_{jl}(\mathbf{Y}^{ji}) - r_{jl}}{S_{jl}^2(\mathbf{Y}^{ji}) + \sigma_{jl}^2} \end{cases}$$

Downward

$$\begin{cases} C_{j,\text{pa}(j)}^2 = q_j \frac{S_j^2(\mathbf{Y}^j)}{S_j^2(\mathbf{Y}^j) + \sigma_j^2} V_{\text{pa}(j)}^2 \\ E_j = \frac{S_j^2(\mathbf{Y}^j)(q_j E_{\text{pa}(j)} + r_j) + \sigma_j^2 M_j(\mathbf{Y}^j)}{S_j^2(\mathbf{Y}^j) + \sigma_j^2} \\ V_j^2 = \frac{S_j^2(\mathbf{Y}^j)}{S_j^2(\mathbf{Y}^j) + \sigma_j^2} \left( \sigma_j^2 + p_j^2 \frac{S_j^2(\mathbf{Y}^j)}{S_j^2(\mathbf{Y}^j) + \sigma_j^2} V_{\text{pa}(j)}^2 \right) \end{cases}$$

back

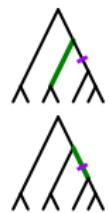
## M Step: Segmentation

$$C_j(\alpha, \tau, \delta) = \sigma_j^{-2} \left( \mathbb{E}[X_j | Y] - q_j \mathbb{E}[X_{\text{pa}(j)} | Y] - r_j - s_j \sum_k \mathbb{I}\{\tau_k = b_j\} \delta_k \right)^2$$

BM :  $r_j = 0$ , each cost is independent.

$$C_j^0(\alpha) = \sigma_j^{-2} \left( \mathbb{E}[X_j | Y] - q_j \mathbb{E}[X_{\text{pa}(j)} | Y] \right)^2$$

$$C_j^1(\alpha, \tau, \delta) = \sigma_j^{-2} \left( \mathbb{E}[X_j | Y] - q_j \mathbb{E}[X_{\text{pa}(j)} | Y] - s_j \sum_k \mathbb{I}\{\tau_k = b_j\} \delta_k \right)^2$$



Algorithm:

- ① Find the  $K$  branches  $j_1, \dots, j_K$  with largest  $C_j^0$ ;
- ② Allocate one change point in the first  $K$  branches;
- ③ For each of these branches, set  $\delta_{j_k}^{(h+1)}$  so that  $C_j^1(\tau, \delta) = 0$

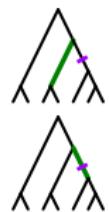
## M Step: Segmentation

$$C_j(\alpha, \tau, \delta) = \sigma_j^{-2} \left( \mathbb{E}[X_j | Y] - q_j \mathbb{E}[X_{\text{pa}(j)} | Y] - r_j - s_j \sum_k \mathbb{I}\{\tau_k = b_j\} \delta_k \right)^2$$

BM :  $r_j = 0$ , each cost is independent.

$$C_j^0(\alpha) = \sigma_j^{-2} \left( \mathbb{E}[X_j | Y] - q_j \mathbb{E}[X_{\text{pa}(j)} | Y] \right)^2$$

$$C_j^1(\alpha, \tau, \delta) = \sigma_j^{-2} \left( \mathbb{E}[X_j | Y] - q_j \mathbb{E}[X_{\text{pa}(j)} | Y] - s_j \sum_k \mathbb{I}\{\tau_k = b_j\} \delta_k \right)^2$$



Algorithm:

- ① Find the  $K$  branches  $j_1, \dots, j_K$  with largest  $C_j^0$ ;
- ② Allocate one change point in the first  $K$  branches;
- ③ For each of these branches, set  $\delta_{j_k}^{(h+1)}$  so that  $C_j^1(\tau, \delta) = 0$

## M Step: Segmentation

$$C_j(\alpha, \tau, \delta) = \sigma_j^{-2} \left( \mathbb{E}[X_j | Y] - q_j \mathbb{E}[X_{\text{pa}(j)} | Y] - r_j - s_j \sum_k \mathbb{I}\{\tau_k = b_j\} \delta_k \right)^2$$

BM :  $r_j = 0$ , each cost is independent.

$$C_j^0(\alpha) = \sigma_j^{-2} \left( \mathbb{E}[X_j | Y] - q_j \mathbb{E}[X_{\text{pa}(j)} | Y] \right)^2$$

$$C_j^1(\alpha, \tau, \delta) = \sigma_j^{-2} \left( \mathbb{E}[X_j | Y] - q_j \mathbb{E}[X_{\text{pa}(j)} | Y] - s_j \sum_k \mathbb{I}\{\tau_k = b_j\} \delta_k \right)^2$$

Algorithm:

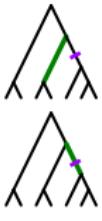
- ① Find the  $K$  branches  $j_1, \dots, j_K$  with largest  $C_j^0$ ;
- ② Allocate one change point in the first  $K$  branches;
- ③ For each of these branches, set  $\delta_{j_k}^{(h+1)}$  so that  $C_j^1(\tau, \delta) = 0$

## M Step: Segmentation

$$C_j(\alpha, \tau, \delta) = \sigma_j^{-2} \left( \mathbb{E}[X_j | Y] - q_j \mathbb{E}[X_{\text{pa}(j)} | Y] - r_j - s_j \sum_k \mathbb{I}\{\tau_k = b_j\} \delta_k \right)^2$$

BM :  $r_j = 0$ , each cost is independent.

$$C_j^0(\alpha) = \sigma_j^{-2} \left( \mathbb{E}[X_j | Y] - q_j \mathbb{E}[X_{\text{pa}(j)} | Y] \right)^2$$

$$C_j^1(\alpha, \tau, \delta) = \sigma_j^{-2} \left( \mathbb{E}[X_j | Y] - q_j \mathbb{E}[X_{\text{pa}(j)} | Y] - s_j \sum_k \mathbb{I}\{\tau_k = b_j\} \delta_k \right)^2$$


Algorithm:

- ① Find the  $K$  branches  $j_1, \dots, j_K$  with largest  $C_j^0$ ;
- ② Allocate one change point in the first  $K$  branches;
- ③ For each of these branches, set  $\delta_{j_k}^{(h+1)}$  so that  $C_j^1(\tau, \delta) = 0$

# M Step: Segmentation

$$C_j(\alpha, \tau, \delta) = \sigma_j^{-2} \left( \mathbb{E}[X_j | Y] - q_j \mathbb{E}[X_{\text{pa}(j)} | Y] - r_j - s_j \sum_k \mathbb{I}\{\tau_k = b_j\} \delta_k \right)^2$$

OU :  $r_j = \beta^{\text{pa}(j)}$ , a cost depends on all its parents.

- Exact minimization: too costly.
- Need of an heuristic.
- Idea: rewrite as a least square:

$$\|D - AU\Delta\|^2$$

with  $D$  a vector of size  $n + m$ ,  $A$  a diagonal matrix of size  $n + m$ ,  $\Delta$  the vector of shifts and  $U$  the incidence matrix of the tree.

- Then use Stepwise selection or LASSO.

back

# Model Selection on $K$

$$\begin{aligned}\hat{K} &= \underset{0 \leq K \leq p-1}{\operatorname{argmin}} \frac{n}{2} \log \left( \frac{\|Y - \hat{s}_K\|_V^2}{n} \right) + \frac{1}{2} \operatorname{pen}'(K) \\ &= \underset{0 \leq K \leq p-1}{\operatorname{argmin}} \|Y - \hat{s}_K\|_V^2 \left( 1 + \frac{\operatorname{pen}(K)}{n - K - 1} \right)\end{aligned}$$

Definition (Baraud et al. (2009))

Let  $D, N > 0$ , and  $X_D \sim \chi^2(D)$ ,  $X_N \sim \chi^2(N)$ ,  $X_D \perp X_N$ .

$$\text{Dkhi}[D, N, x] = \frac{1}{\mathbb{E}[X_D]} \mathbb{E} \left[ \left( X_D - x \frac{X_N}{N} \right)_+ \right], \quad \forall x > 0$$

$$\text{Dkhi}[D, N, \text{EDkhi}[D, N, q]] = q, \quad \forall 0 < q \leq 1$$

# Model Selection with Unknown Variance

Theorem (Baraud et al. (2009))

*Under the following setting:*

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

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

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

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

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

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

# IID Framework ( $\alpha = 0$ )

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

Then:

$$\begin{aligned}
 \Omega' &= \sum_{\eta \in \mathcal{M}} (D_\eta + 1)e^{-L_\eta} = \sum_{\eta \in \mathcal{M}} (K_\eta + 2)e^{-L_\eta} \\
 &= \sum_{K=0}^{p-1} |S_K^{PI}| (K+2)e^{-L_K} = \sum_{K=0}^{p-1} |S_K^{PI}| (K+2)e^{-(\log|S_K^{PI}| + 2\log(K+2))} \\
 &= \sum_{K=0}^{p-1} \frac{1}{K+2} \leq \log(p) \leq \log(n)
 \end{aligned}$$

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  $\eta \in \mathcal{M}$ .

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

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

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

$$S'_\eta = L^{-1}S_\eta, \quad \hat{Y}'_\eta = \text{Proj}_{S'_\eta} Y' = \underset{a' \in S'_\eta}{\operatorname{argmin}} \|Y - La'\|_V^2 = L^{-1}\hat{Y}_\eta$$

$$\|\mathbb{E}[Y] - \hat{Y}_{\hat{\eta}}\|_V^2 = \|\mathbb{E}[Y'] - \hat{Y}'_{\hat{\eta}}\|^2, \quad \|Y - \hat{Y}_\eta\|_V^2 = \|Y' - \hat{Y}'_\eta\|^2$$

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

back

$$OU \iff BM$$

## Expectations

$$\mathbb{E}[Y | X_1 = \mu] = T \underbrace{W(\alpha) \Delta^{OU}}_{\Delta^{BM}}$$

**Remark:**  $\mu^{BM} = \lambda^{OU} = \mu e^{-\alpha h} + \beta_0(1 - e^{-\alpha h})$

## Variance

$$\text{Cov}[Y_i; Y_j | X_1 = \mu] = \sigma^2 \times \underbrace{\frac{1}{2\alpha} e^{-2\alpha h} (e^{2\alpha t'_{ij}} - 1)}_{t'_{ij}}$$

$OU \iff BM$  on a re-scaled tree with  $t' = e^{-2\alpha h}(e^{2\alpha t} - 1)$

# OU $\iff$ BM

OU  $\iff$  BM on a re-scaled tree with  $t' = e^{-2\alpha h}(e^{2\alpha t} - 1)$

## Remarks:

- This only works for an *ultrametric* tree.
- The laws of the internal nodes is changed.
- This is *not* the following standart time transformation:

## Lemma (Brownian Solution for the OU)

The stochastic process defined by:

$$X_t = X_0 e^{-\alpha t} + \beta(1 - e^{-\alpha t}) + \frac{\sigma}{\sqrt{2\alpha}} e^{-\alpha t} B_{e^{2\alpha t}-1}$$

is an OU, solution of the EDS  $dX_t = \alpha(\beta - X_t) + \sigma dB_t$ .

back

# Cardinal of Equivalence Classes

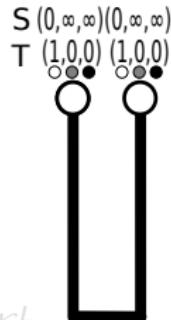
Initialization For tips

Propagation

$$\mathcal{K}_k^l = \operatorname{argmin}_{1 \leq p \leq K} \{ S_{ij}(p) + \mathbb{I}\{p \neq k\} \}$$

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

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



Termination Sum on the root vector

back

# Cardinal of Equivalence Classes

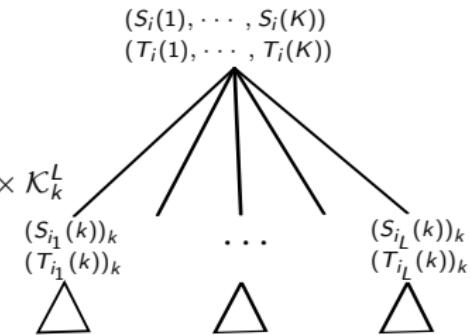
Initialization For tips

Propagation

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

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

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



Termination Sum on the root vector

back

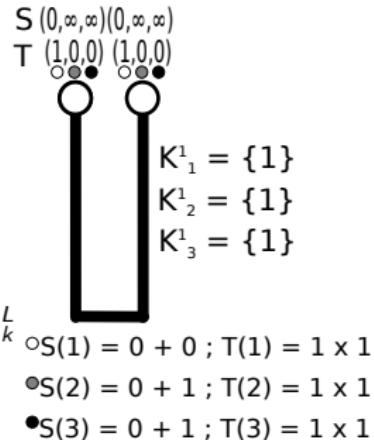
# Cardinal of Equivalence Classes

Initialization For tips  
 Propagation

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

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

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



Termination Sum on the root vector

back

# Cardinal of Equivalence Classes

Initialization For tips  
 Propagation

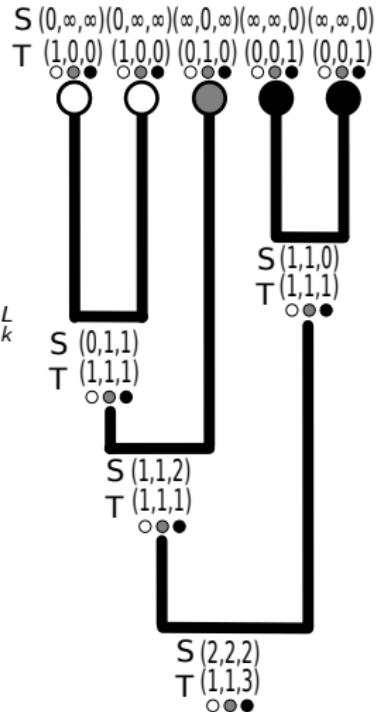
$$\mathcal{K}_k^l = \operatorname{argmin}_{1 \leq p \leq K} \{ S_{i_l}(p) + \mathbb{I}\{p \neq k\} \}$$

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

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

Termination Sum on the root vector

back



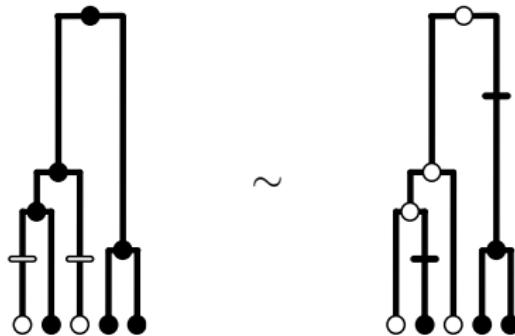
# Linking Shifts and Clustering

Assumption “No Homoplasy”: 1 shift = 1 new color

Proposition “ $K$  shifts  $\iff K + 1$  clusters”

# Linking Shifts and Clustering

Assumption “No Homoplasy”: 1 shift = 1 new color

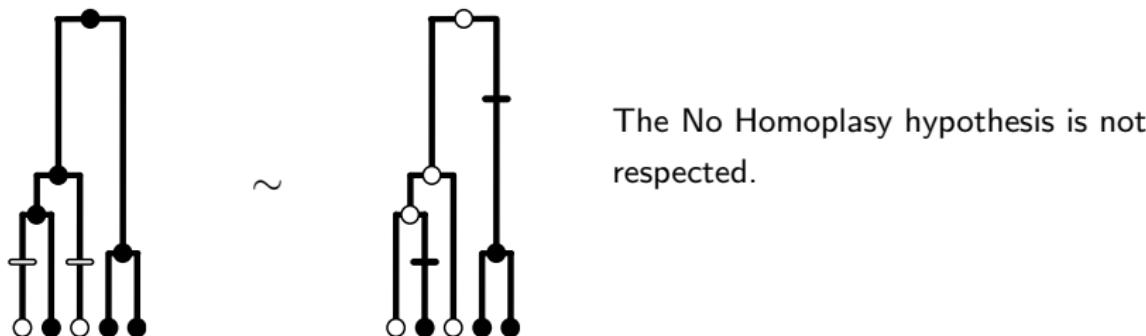


The No Homoplasy hypothesis is not respected.

Proposition “ $K$  shifts  $\iff K + 1$  clusters”

# Linking Shifts and Clustering

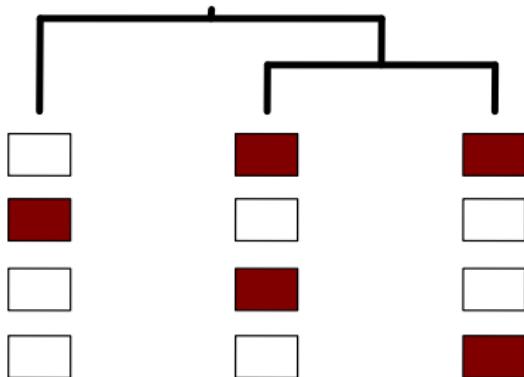
Assumption “No Homoplasy”: 1 shift = 1 new color



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

# Definitions

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



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

- $N_2^{(\mathcal{T}_3)} = 3$ : partitions 1 and 2 are equivalent
- $A_2^{(\mathcal{T}_3)} = 4$ : one marked color ("white = ancestral state")

*Partitions in two groups for a binary tree with 3 tips*

## General Formula (Binary Case)

If  $\mathcal{T}$  is a binary tree, consider  $T_\ell$  and  $T_r$  the left and right sub-trees of  $\mathcal{T}$ . Then:

$$\begin{cases} N_K^{(\mathcal{T})} = \sum_{k_1+k_2=K} N_{k_1}^{(T_\ell)} N_{k_2}^{(T_r)} + \sum_{k_1+k_2=K+1} A_{k_1}^{(T_\ell)} A_{k_2}^{(T_r)} \\ A_K^{(\mathcal{T})} = \sum_{k_1+k_2=K} A_{k_1}^{(T_\ell)} N_{k_2}^{(T_r)} + N_{k_1}^{(T_\ell)} A_{k_2}^{(T_r)} + \sum_{k_1+k_2=K+1} A_{k_1}^{(T_\ell)} A_{k_2}^{(T_r)} \end{cases}$$

We get:

$$N_{K+1}^{(\mathcal{T})} = N_{K+1}^{(n)} = \binom{2n - 2 - K}{K} \quad \text{and} \quad A_{K+1}^{(\mathcal{T})} = A_{K+1}^{(n)} = \binom{2n - 1 - K}{K}$$

# Recursion Formula (General Case)

If we are at a node defining a tree  $\mathcal{T}$  that has  $p$  daughters, with sub-trees  $\mathcal{T}_1, \dots, \mathcal{T}_p$ , then we get the following recursion formulas:

$$\left\{ \begin{array}{l} N_K^{(\mathcal{T})} = \sum_{\substack{k_1 + \dots + k_p = K \\ k_1, \dots, k_p \geq 1}} \prod_{i=1}^p N_{k_i}^{(\mathcal{T}_i)} + \sum_{\substack{I \subset [1, p] \\ |I| \geq 2}} \sum_{\substack{k_1 + \dots + k_p = K + |I| - 1 \\ k_1, \dots, k_p \geq 1}} \prod_{i \in I} A_{k_i}^{(\mathcal{T}_i)} \prod_{i \notin I} N_{k_i}^{(\mathcal{T}_i)} \\ A_K^{(\mathcal{T})} = \sum_{\substack{I \subset [1, p] \\ |I| \geq 1}} \sum_{\substack{k_1 + \dots + k_p = K + |I| - 1 \\ k_1, \dots, k_p \geq 1}} \prod_{i \in I} A_{k_i}^{(\mathcal{T}_i)} \prod_{i \notin I} N_{k_i}^{(\mathcal{T}_i)} \end{array} \right.$$

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

back

# Simulations Design

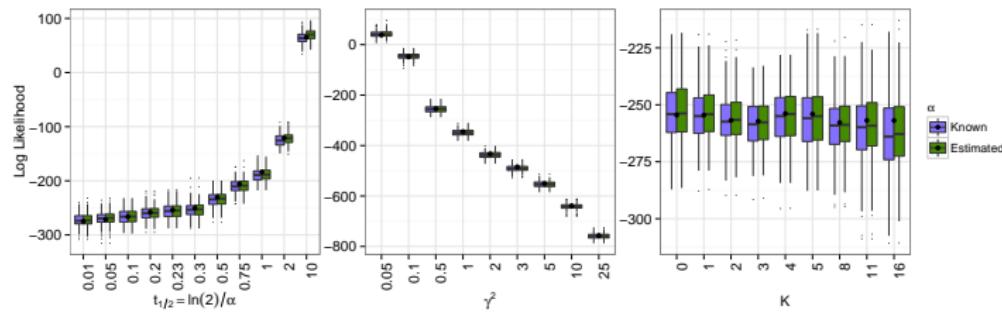
(Uyeda and Harmon, 2014)

- 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  $\sim \frac{1}{2}\mathcal{N}(4, 1) + \frac{1}{2}\mathcal{N}(-4, 1)$
- 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):

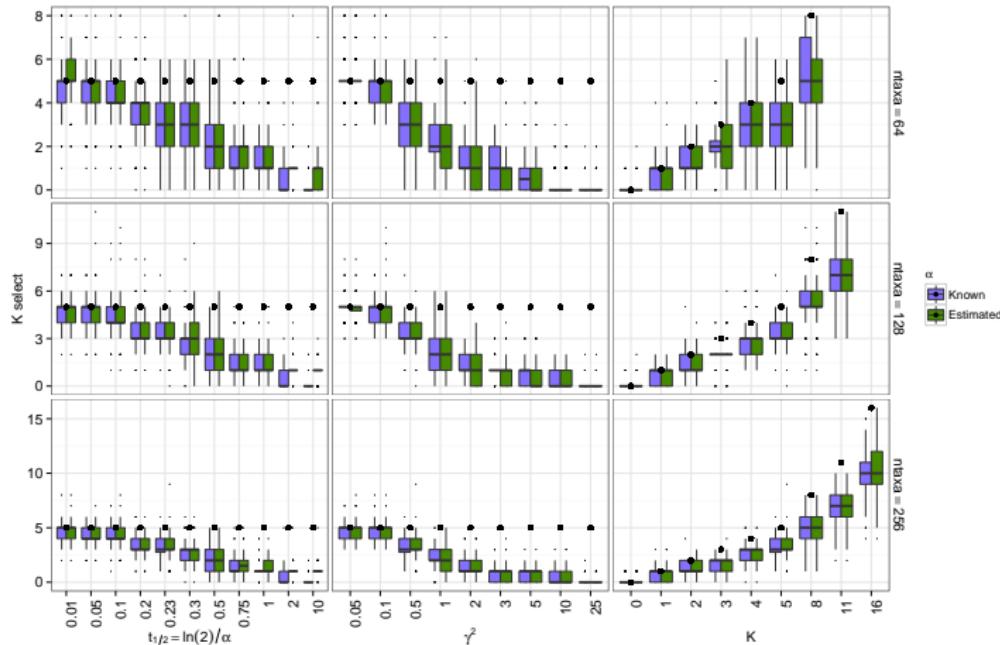
- $\alpha$  known: 6 minutes per estimation (66 days in total).
- $\alpha$  unknown: 52 minutes per estimation (570 days in total).

# Log-Likelihood

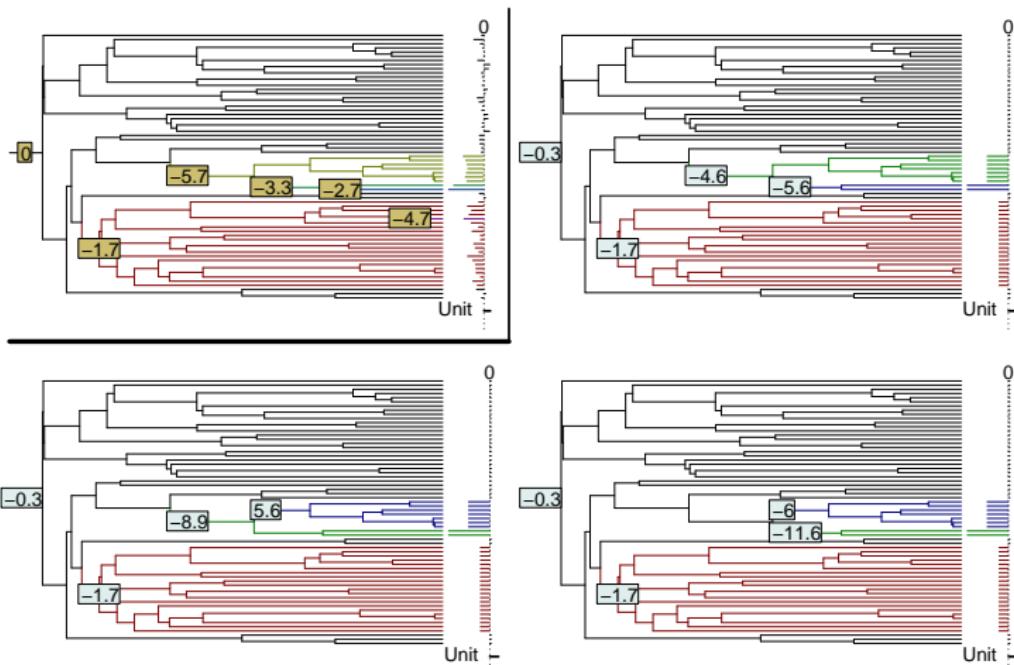


*Log likelihood for a tree with 256 tips. Solid black dots are the median of the log likelihood for the true parameters.*

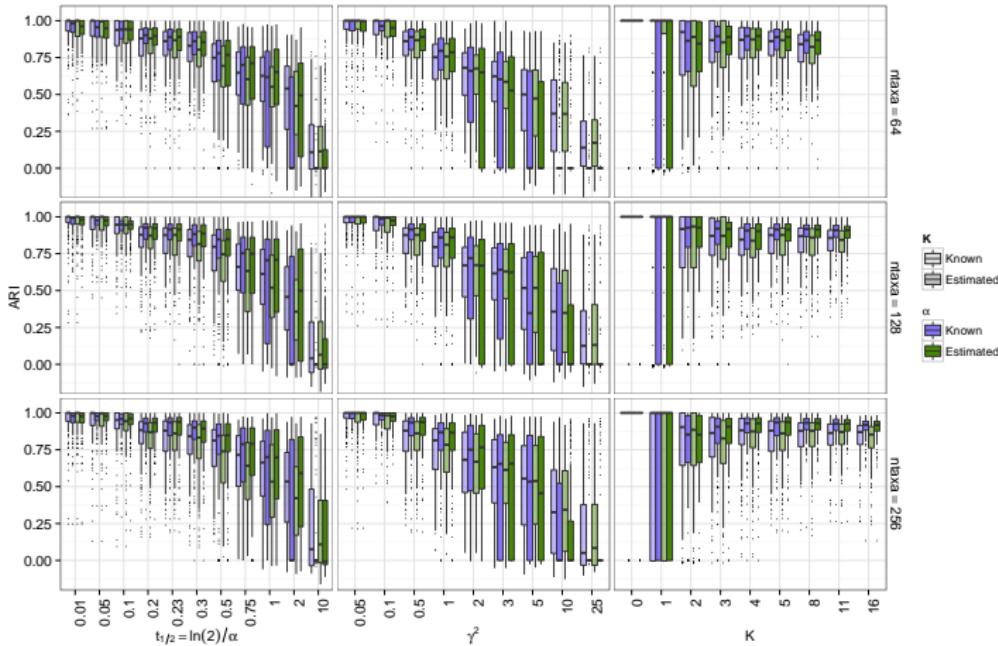
# Number of Shifts



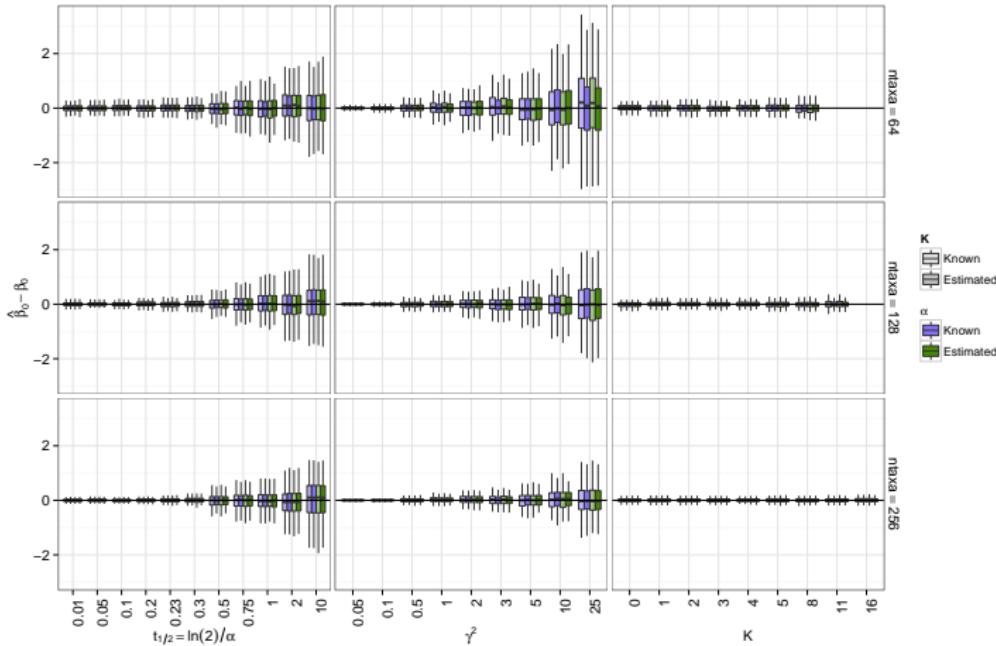
# One Example



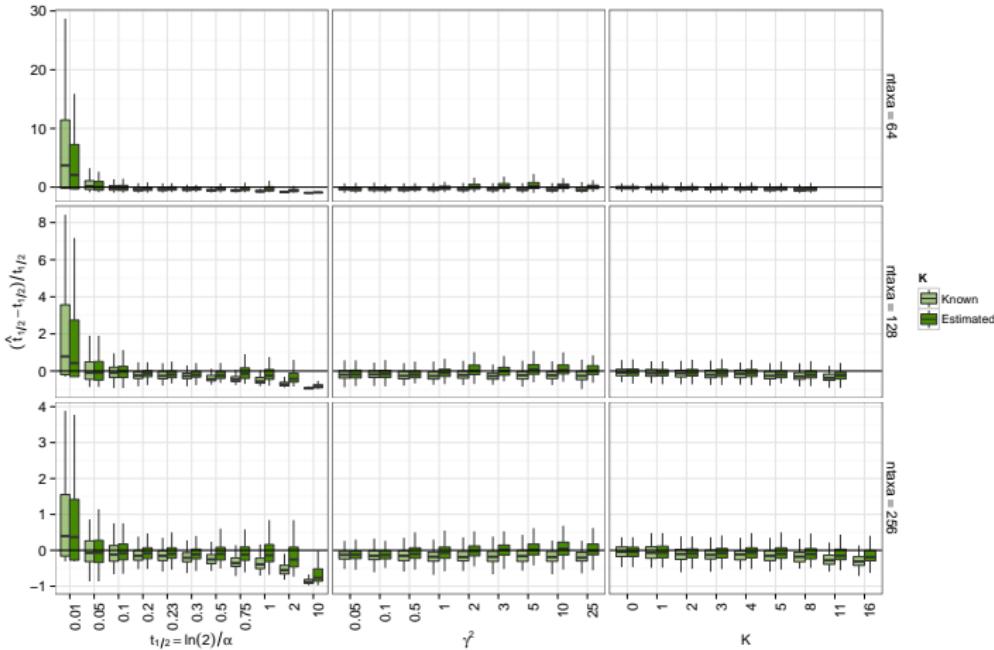
# Adjusted Rand Index



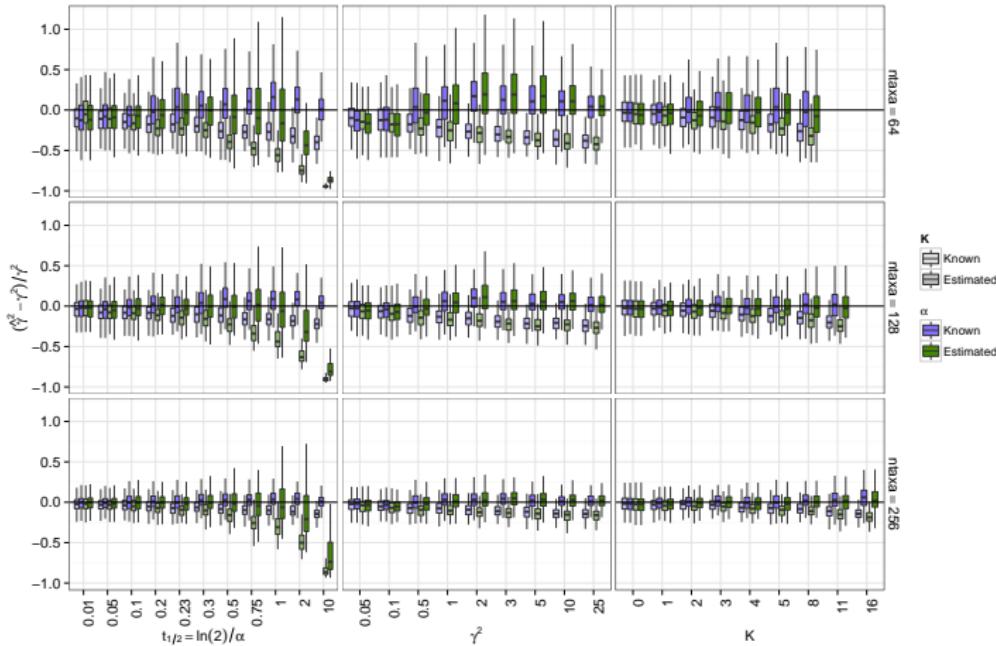
# Parameters: $\beta_0$



# Parameters: $\alpha$



# Parameters: $\gamma^2$



# BM Model

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

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

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

$$\text{Var}[\text{vec}(\mathbf{Y})] = \mathbf{C}_n \otimes \mathbf{R}$$

Shifts  $K$  shifts  $\delta_1, \dots, \delta_K$  vectors size  $p$

→ All the characters shift at the same time

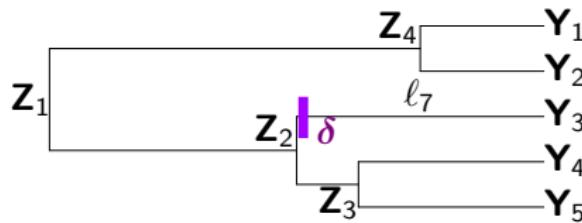
# BM Model

## Linear Model Representation

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

## Incomplete Data Representation

$$\mathbf{Y}_3 | \mathbf{Z}_2 \sim \mathcal{N}\left(\mathbf{Z}_2 + \boldsymbol{\delta}, \ell_7 \mathbf{R}\right)$$



## 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  $\mathbf{A}$  ( $p \times p$ ) “selection strength”

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

### Covariances

$$\begin{aligned} \text{Cov} [\mathbf{X}_i; \mathbf{X}_j] &= e^{-\mathbf{A}t_i} \boldsymbol{\Gamma} e^{-\mathbf{A}^T t_j} \\ &+ e^{-\mathbf{A}(t_i - t_{ij})} \left( \int_0^{t_{ij}} e^{-\mathbf{A}\nu} \boldsymbol{\Sigma} \boldsymbol{\Sigma}^T e^{-\mathbf{A}^T \nu} d\nu \right) e^{-\mathbf{A}^T (t_j - t_{ij})} \end{aligned}$$

Shifts  $K$  shifts  $\delta_1, \dots, \delta_K$  vectors size  $p$   
 ↳ On the optimal values

# OU Model: $\mathbf{A}$ 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:

$$\text{Cov}[Y_{il}; Y_{jq}] = \frac{1}{2\alpha} e^{-2\alpha h} (e^{2\alpha t_{ij}} - 1) R_{lq}$$

↪ Can be reduced to a BM on a re-scaled tree

# EM algorithm

BM Natural generalization of the univariate case.

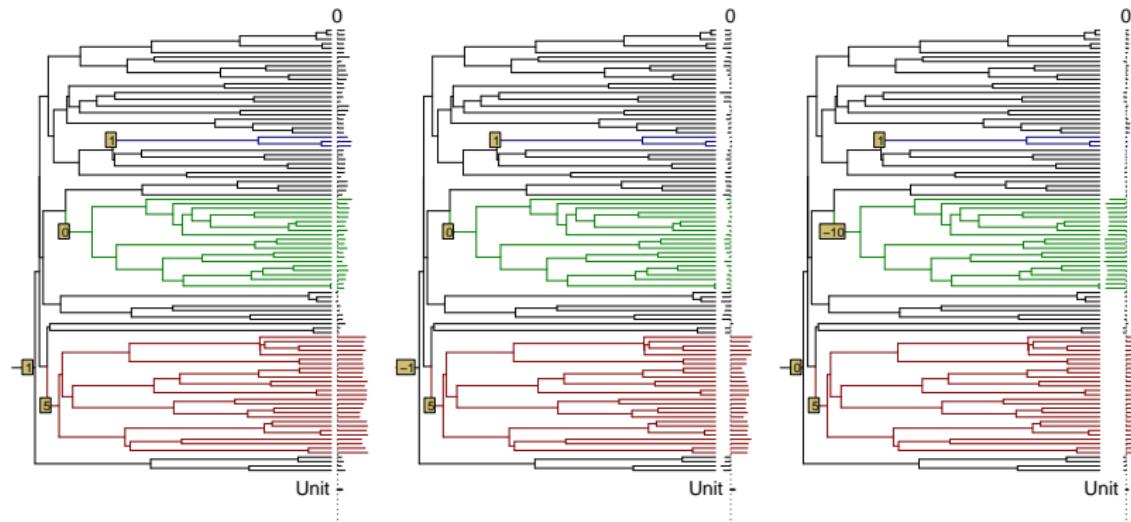
OU M step intractable in general.

Incomplete Data Model: Can readily handle missing data.

# Model Selection

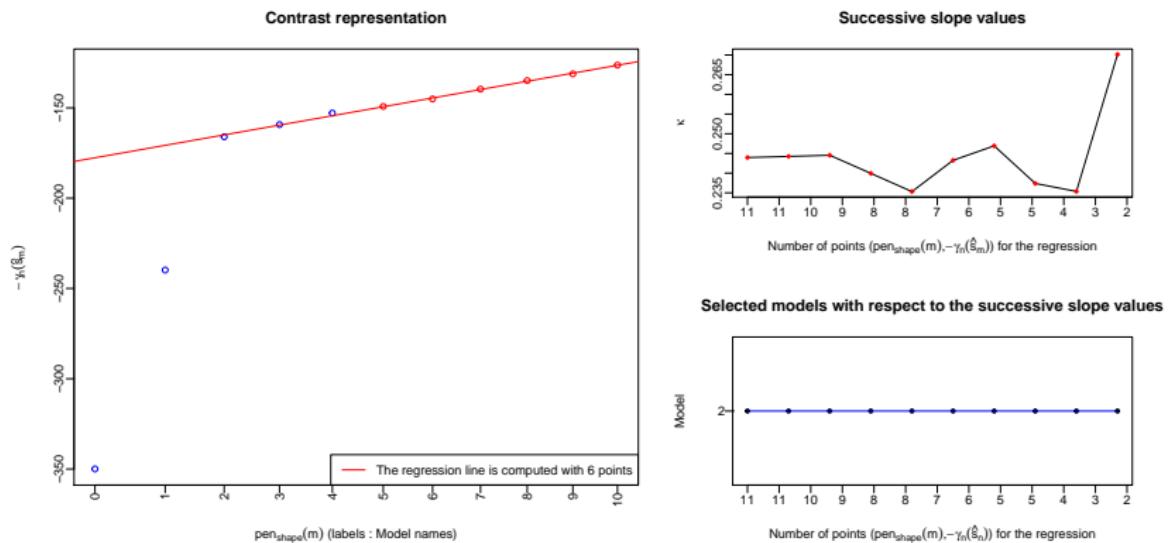
- Previous criterion cannot be applied
- Solution: “Slope Heuristic”-based method
  - Massart (2007)
    - oracle inequality with known variance
    - penalty up to a multiplicative constant
  - Baudry et al. (2012)
    - Slope-heuristic method to calibrate the constant
    - Implemented in capushe (Brault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

# Model Selection: Toy Example



*Figure: Simulated Process.*

# Model Selection: Toy Example



*Figure:* capushe output for penalized log-likelihood.

# Model Selection: Toy Example

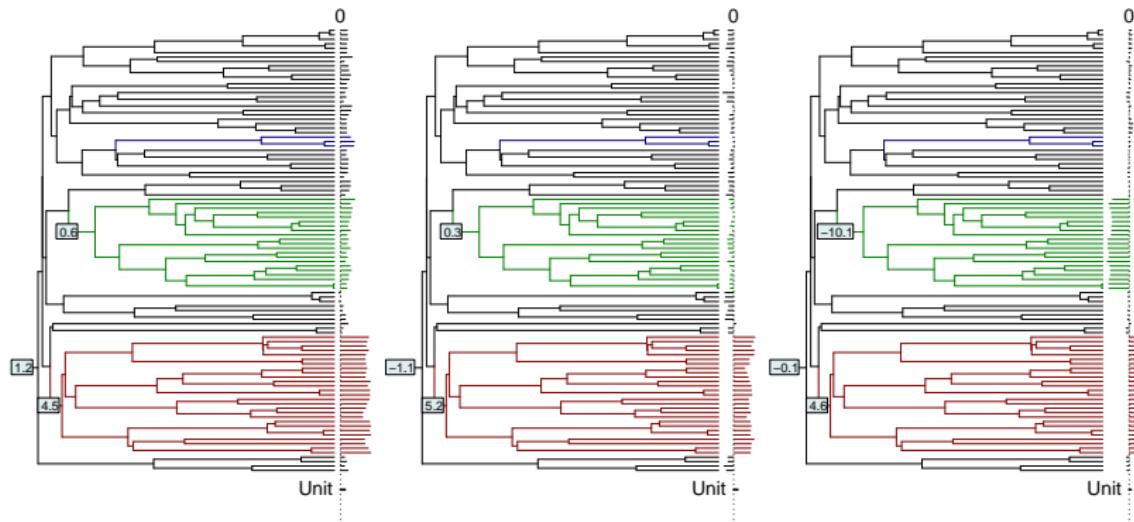


Figure: Reconstructed Process.