

SERIK SAGITOV

25 slides

Chalmers University and University of Gothenburg

A PHYLOGENETIC CONFIDENCE INTERVAL FOR THE OPTIMAL TRAIT VALUE

<http://arxiv.org/abs/1207.6488>

joint work with

Krzysztof Bartoszek



Comparative phylogenetics

Interspecies correlation

Yule-Brownian-Motion model

Yule-Ornstein-Uhlenbeck model

Adaptation to an optimal value

Sample mean and variance of the trait values

Phylogenetic confidence intervals

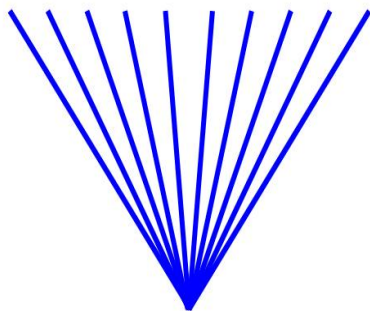
Exact and approximate formulae

Related projects

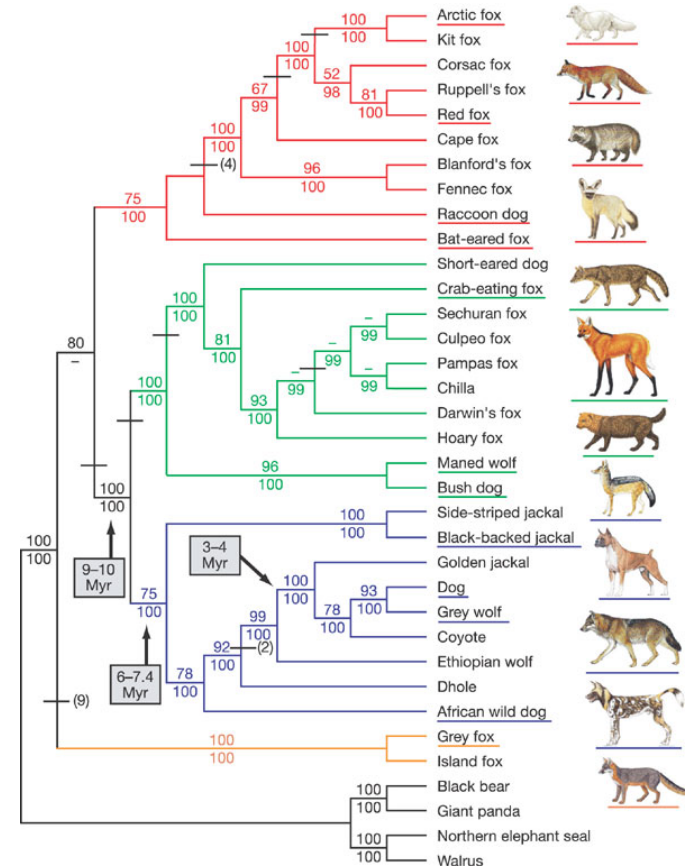
COMPARATIVE PHYLOGENETICS

Comparative phylogenetics: studies various trait values like log-bodysize (X_1, \dots, X_n) for a family of related species

**What
Conventional
Statistical
Methods
Assume**



**What
Evolution
Provides**



We will assume throughout that the species tree is unknown.

INTERSPECIES CORRELATION

In our previous paper [Interspecies correlation for neutrally evolving traits](#). *J. Theor. Biol.* 309 (2012) 11-19 we obtain explicit formulae for

$$\rho_n = \frac{1}{\binom{n}{2} \text{Var} [X]} \sum_{1 \leq i < j \leq n} \text{Cov} [X_i, X_j]$$

assuming that the trait evolves along a lineage as a **B**rownian **M**otion.

Aldous-Popovic (2005), Stadler (2008): [unknown species trees are modeled by branching processes conditioned on \$n\$ tips.](#)

In particular, we found for the **Y**ule-**B**M-model

$$\rho_n = \frac{2}{n-1} \left(\frac{n}{a_n} - 1 \right),$$

where $a_n = \sum_{i=1}^n (1/i)$ are the harmonic numbers.



YBM-MODEL FOR THE TRAIT EVOLUTION



G.U. Yule (1924):

let us model speciation by a pure birth process with parameter λ .

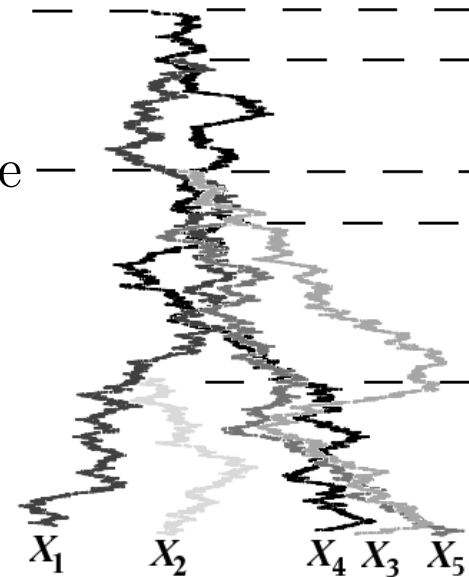


J. Felsenstein (1985):

the trait value evolves along

a lineage neutrally as a Brownian $B(t)$ motion

with variance σ^2 and ancestral state $B(0) = X_0$.



Coalescent processes model the gene trees.

Conditioned branching processes model the species trees.

YULE-ORNSTEIN-UHLENBECK MODEL

In this paper we work with the **YOU-model**:
(conditioned **Y**ule tree)+(OU trait evolution)

L.Ornstein and
G.Uhlenbeck (1930)

$$dX(t) = -\alpha(X(t) - \theta)dt + \sigma dB(t)$$

$$X(0) = X_0$$

Just six parameters fully describe the YOU-model:

- number of species n
- speciation rate λ , wlog we assume $\lambda = 1$
- optimal trait value θ ,
- **adaptation** rate $\alpha > 0$
- ancestral state X_0

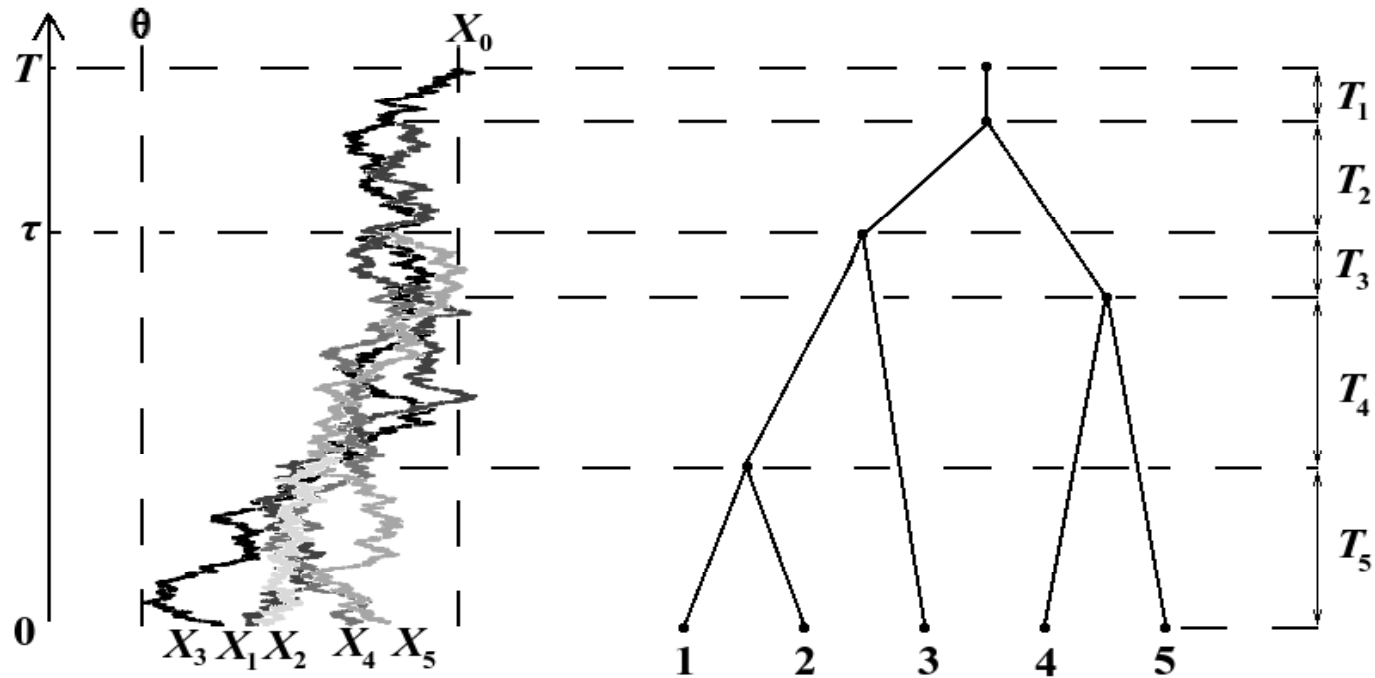
$$\alpha = 0 \Rightarrow X(t) = B(t)$$

- noise size $\sigma > 0$

$$\sigma = 0 \Rightarrow X(t) = \theta + e^{-\alpha t}(X_0 - \theta)$$



YULE-ORNSTEIN-UHLENBECK MODEL



T = time to the origin, T_i = inter-speciation times

τ = MRCA time for a random pair of species

θ = optimal trait value, X_0 = ancestral trait value

(X_1, \dots, X_n) = observed trait values identically distr. but dependent

YULE-ORNSTEIN-UHLENBECK MODEL

OU-process: the distribution of $X(t)$ is normal with

$$\begin{aligned} \mathbb{E}[X(t)] &= \theta + e^{-\alpha t}(X_0 - \theta) \\ \text{Var}[X(t)] &= \sigma_0^2(1 - e^{-2\alpha t}) \end{aligned} \quad \boxed{\sigma_0^2 = \frac{\sigma^2}{2\alpha}}$$

Stationary regime: $\mathbb{E}[X(t)] \rightarrow \theta$ and $\text{Var}[X(t)] \rightarrow \sigma_0^2$ as $t \rightarrow \infty$.

The Yule tree conditioned on n tips:

run forward a linear birth process from a single branch and stop at the n -th speciation disregarding the nascent branch.

Interspeciation times

T_i = time duration with exactly i branches in the tree

T_i are independent exponential with rates $i = 1, \dots, n$.

Backward description reminds a stretched Kingman's coalescent.

SAMPLE MEAN AND VARIANCE

Classical summary statistics: sample mean and sample variance

$$\bar{X}_n = \frac{X_1 + \dots + X_n}{n} \quad \text{and} \quad S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

Denoting X, Y a random pair of sampled trait values we get

$$\begin{aligned} \mathbb{E} [\bar{X}_n] &= \mathbb{E} [X] & \rho_n &= \text{Cov} [X, Y] / \text{Var} [X] \\ \text{Var} [\bar{X}_n] &= \text{Var} [X] - (n-1)n^{-1}(1-\rho_n) \text{Var} [X] \\ \mathbb{E} [S_n^2] &= (1-\rho_n) \text{Var} [X] \end{aligned}$$

and it becomes crucial to know the Laplace transforms of (T, τ) :

$$\begin{aligned} \mathbb{E} [X] &= \theta + (X_0 - \theta) \mathbb{E} [e^{-\alpha T}] \\ \text{Var} [X] &= \sigma_0^2(1 - \mathbb{E} [e^{-2\alpha T}]) + (X_0 - \theta)^2 \text{Var} [e^{-\alpha T}] \\ \text{Cov} [X, Y] &= \sigma_0^2(\mathbb{E} [e^{-2\alpha \tau}] - \mathbb{E} [e^{-2\alpha T}]) + (X_0 - \theta)^2 \text{Var} [e^{-\alpha T}] \end{aligned}$$

LAPLACE TRANSFORMS OF T AND τ

A Yule tree conditioned on n tips has $n - 1$ splittings.

Let κ_n be the splitting number (counted forward in time) corresponding to the coalescent of two randomly sampled tips.

Key Lemma. We claim that

$$\mathbb{P}(\kappa_n = k) = \frac{2(n+1)}{(n-1)(k+2)(k+1)}, \quad k = 1, \dots, n-1$$

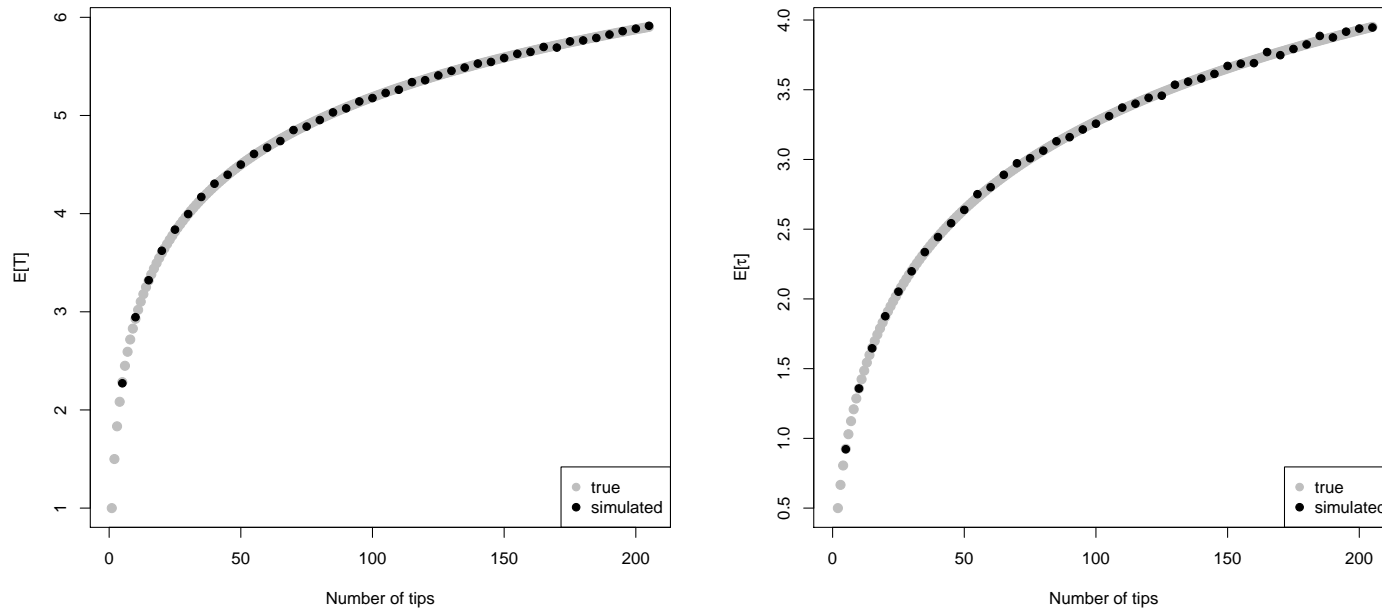
and that for any $x > -1$ and $y > -1$

$$\mathbb{E} \left[e^{-x(T-\tau)-y\tau} \right] = \frac{2(n+1)b_{n,y}}{(n-1)} \sum_{k=1}^{n-1} \frac{b_{k,x}}{(k+2)(k+1)b_{k,y}}$$

$$b_{n,x} = \frac{1}{1+x} \cdot \frac{2}{2+x} \cdots \frac{n}{n+x}$$

MOMENTS OF T AND τ

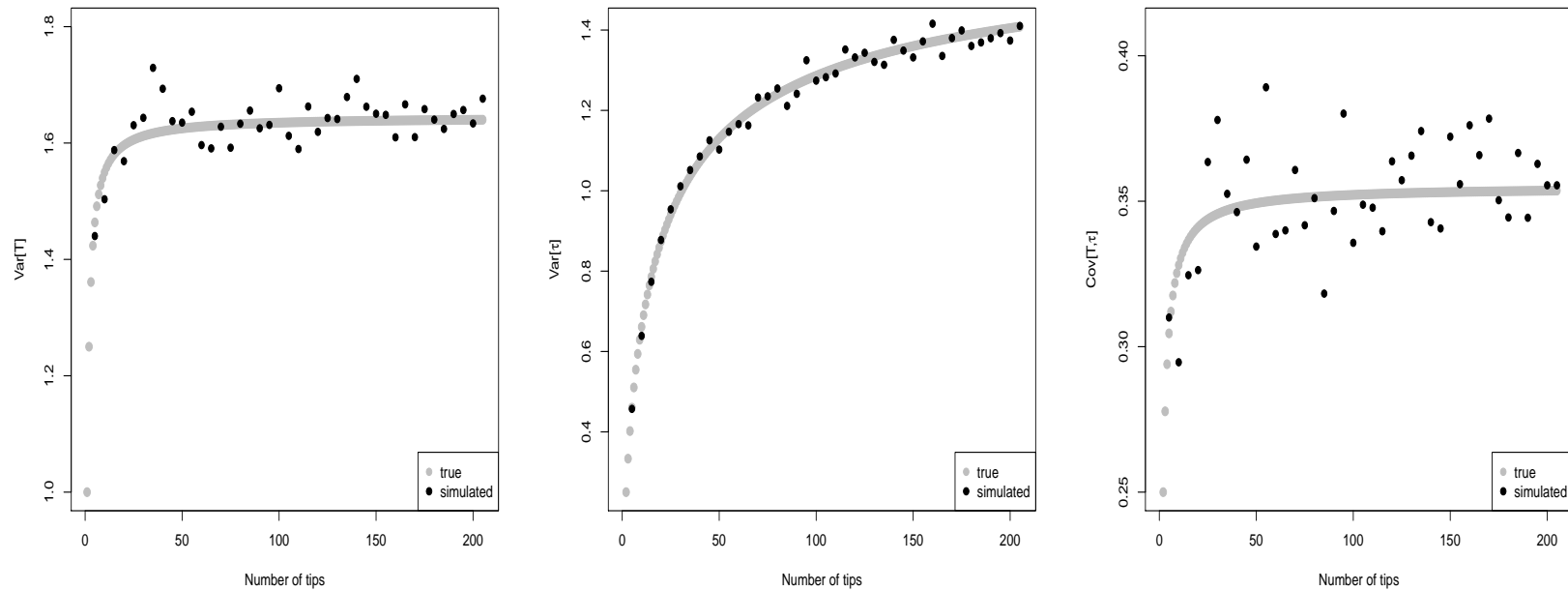
Using the Key Lemma we derive recursively all the joint moments $E[(T - \tau)^k \tau^m]$ for $k \geq 0$ and $m \geq 0$.



Simulated and true values of $E[T]$ and $E[\tau]$.

Each point comes from 10^4 simulated Yule trees.

MOMENTS OF T AND τ



Simulated and true values of $\text{Var}[T]$, $\text{Var}[\tau]$, and $\text{Cov}[T, \tau]$.

Each point comes from 10^4 simulated Yule trees.

EXACT FORMULAE FOR THE MOMENTS OF \bar{X}_n AND S_n^2

Using the Key Lemma we compute

$$\begin{aligned} \mathbf{E} [e^{-\alpha T}] &= b_{n,\alpha} = \frac{1}{1+\alpha} \cdot \frac{2}{2+\alpha} \cdots \frac{n}{n+\alpha} \\ \mathbf{E} [e^{-\alpha \tau}] &= \frac{2 - (n+1)(y+1)b_{n,\alpha}}{(n-1)(\alpha-1)} \end{aligned}$$

which imply

$$\begin{aligned} \mathbf{E} [\bar{X}_n] &= b_{n,\alpha} X_0 + (1 - b_{n,\alpha}) \theta \\ \text{Var} [\bar{X}_n] &= \sigma_0^2 \cdot \frac{2\alpha + 1 - (4\alpha n + 2\alpha + 1)b_{n,2\alpha}}{(2\alpha - 1)n} + (X_0 - \theta)^2 (b_{n,2\alpha} - b_{n,\alpha}^2) \\ \mathbf{E} [S_n^2] &= \sigma_0^2 \left(1 + \frac{(2\alpha + 1)(n+1)b_{n,2\alpha} - 2}{(2\alpha - 1)(n-1)} \right) \end{aligned}$$

PHYLOGENETIC POINT ESTIMATES

For example, with $\alpha = 1$ we get $\mathbb{E} [\bar{X}_n] = \frac{1}{n+1} X_0 + \frac{n}{n+1} \theta$.

\bar{X}_n and S_n^2 are asymptotically unbiased estimates of the optimum θ and the stationary OU-variance σ_0^2 respectively.

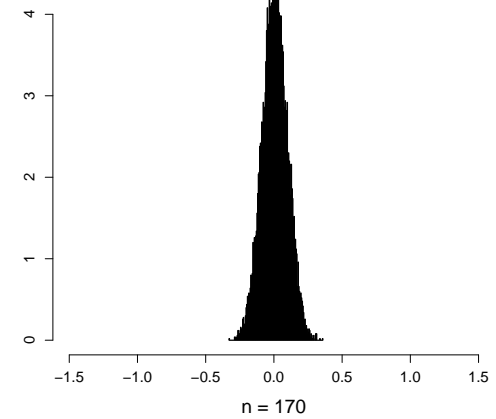
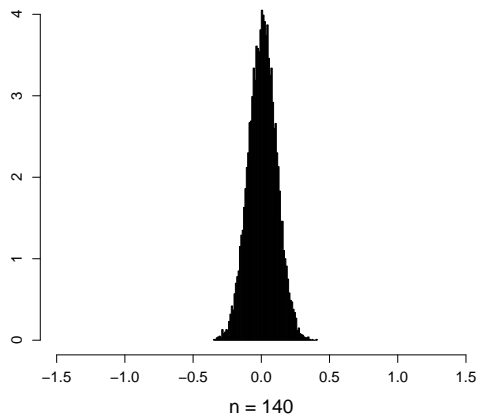
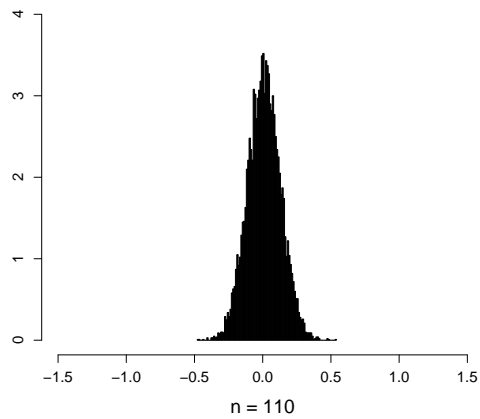
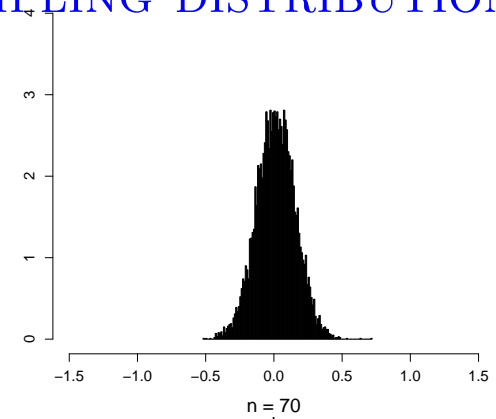
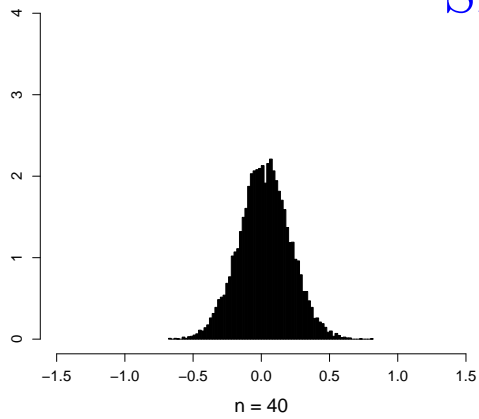
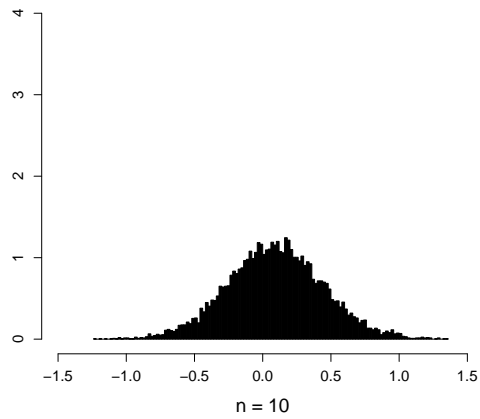
Consistency of the estimate \bar{X}_n

$$\text{Var} [\bar{X}_n] \sim \sigma^2 \cdot \begin{cases} C_{\alpha, \delta} \cdot n^{-2\alpha}, & \text{if } 0 < \alpha < 0.5, \\ 2n^{-1} \ln n, & \text{if } \alpha = 0.5, \\ \frac{2\alpha+1}{2\alpha(2\alpha-1)n}, & \text{if } \alpha > 0.5, \end{cases}$$

where $\delta = \frac{|X_0 - \theta|}{\sigma}$ is the normalized deviation of X_0 from θ , and

$$C_{\alpha, \delta} = \frac{2}{1 - 2\alpha} \Gamma(2\alpha + 1) + \delta^2 \Gamma(2\alpha + 1) - \delta^2 \Gamma^2(\alpha + 1).$$

SAMPLING DISTRIBUTIONS



Histograms of \bar{X}_n for $\alpha = 1$, $X_0 = 1$, $\theta = 0$, $\sigma = 1$.

Each histogram is based on 10^4 independently simulated samples.

PHYLOGENETIC CONFIDENCE INTERVAL

For sufficiently strong adaptation when $\alpha > 0.5$ as $n \rightarrow \infty$ we have

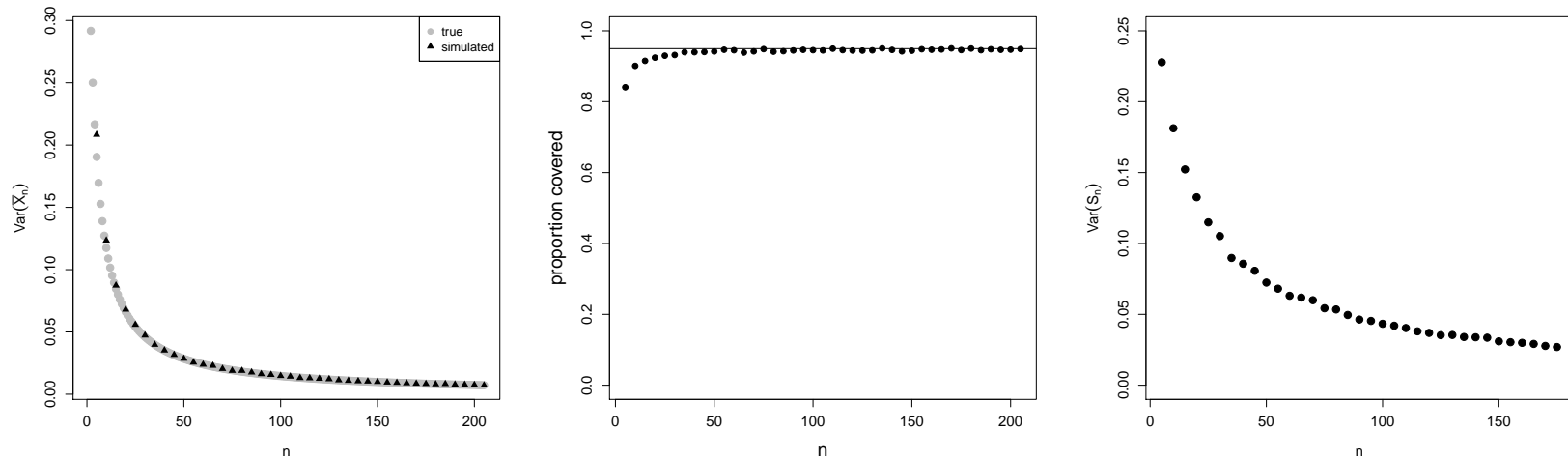
$$\begin{aligned} \mathbb{E} [\bar{X}_n] &= \theta + O(\delta n^{-\alpha}) \\ \text{Var} [\bar{X}_n] &\sim \frac{2\alpha + 1}{2\alpha - 1} \cdot \frac{\sigma_0^2}{n} \\ \mathbb{E} [S_n^2] &\rightarrow \sigma_0^2 \qquad \sigma_0^2 = \sigma^2 / (2\alpha) \end{aligned}$$

Let $\alpha > 0.5$ and $\delta = \frac{|X_0 - \theta|}{\sigma}$ stay bounded as $n \rightarrow \infty$. We propose a convenient approximate formula for a 95% confidence interval for θ :

$$\boxed{\bar{X}_n \pm 1.96 \cdot S_n \cdot \frac{K_\alpha}{\sqrt{n}}} \qquad K_\alpha = \sqrt{\frac{2\alpha + 1}{2\alpha - 1}}$$

Compared to $(\bar{X}_n \pm 1.96 \cdot S_n / \sqrt{n})$ it has an extra factor $K_\alpha > 1$ reflecting a positive correlation among the sample observations.

PHYLOGENETIC CONFIDENCE INTERVAL



Simulations results for $\alpha = 1$, $\sigma = 1$, $X_0 = 1$, $\theta = 0$.

Left: variance of the sample mean, simulated and theoretical values.

Center: each dot gives the coverage proportion of 10^4 confidence intervals. The horizontal line: predicted coverage of 95%.

Right: simulated values of $\text{Var}[S_n^2]$ decrease toward zero with n indicating that S_n^2 is a consistent estimate of the stationary variance.

ASYMPTOTIC NORMALITY

A strict mathematical justification of this formula requires two results which will be addressed elsewhere.

On one hand, we have to prove that

$$\text{Var} [S_n^2] \rightarrow 0, \quad n \rightarrow \infty.$$

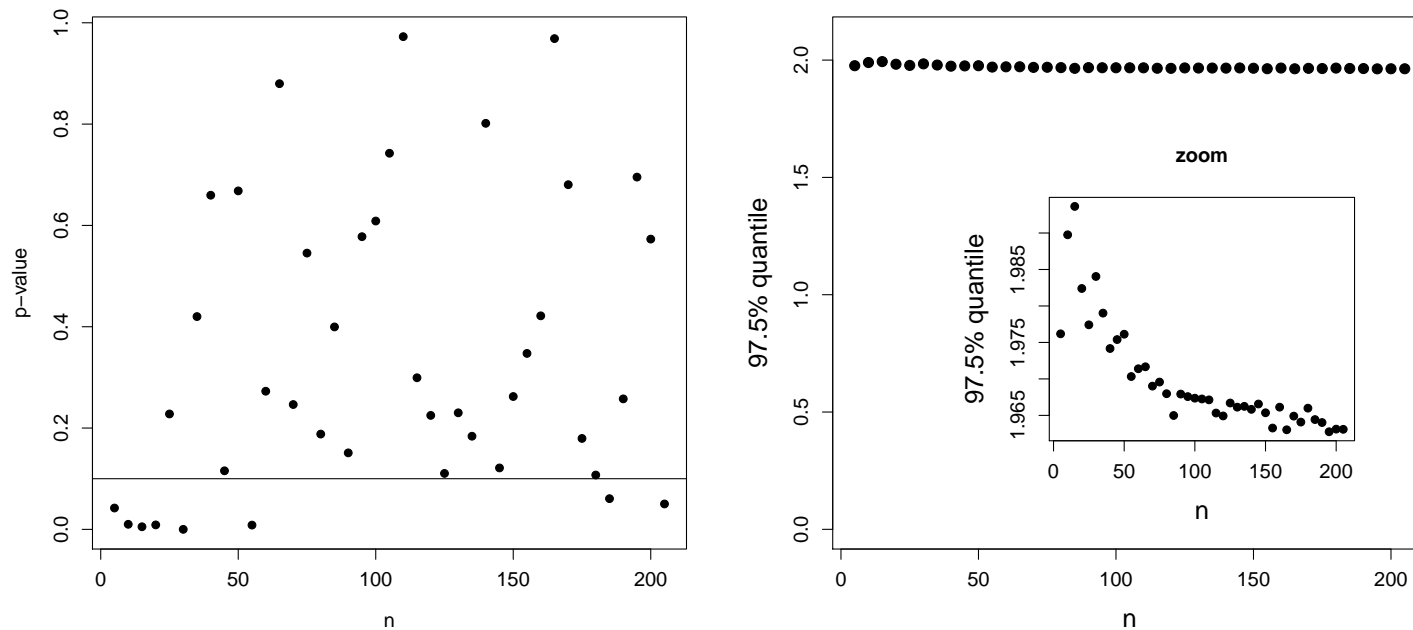
On the other hand, the following Central Limit Theorem has to be verified

Conjecture. As $n \rightarrow \infty$, the standardized sample mean

$$Z_n = (\bar{X}_n - \text{E} [\bar{X}_n]) (\text{Var} [\bar{X}_n])^{-1/2}$$

is asymptotically normally distributed with parameters (0,1).

ASYMPTOTIC NORMALITY



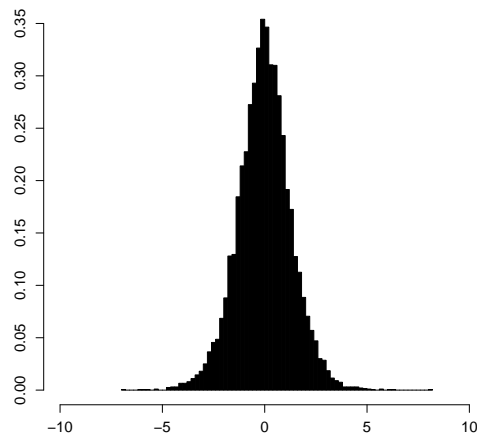
Simulation results for $\alpha = 1$, $X_0 = 1$, $\theta = 0$, $\sigma = 1$.

Left: p-value of Shapiro–Wilk test for normality of Z_n .

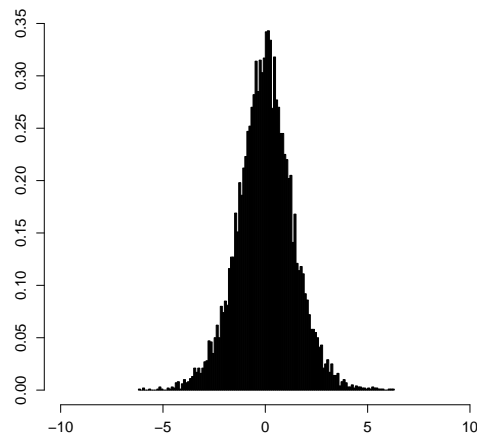
Right: A more thorough 0.975 level quantile check of normality for Z_n is given by a hybrid Monte–Carlo–numerical approach.

X_0 ESTIMATE FOR THE YMB-MODEL

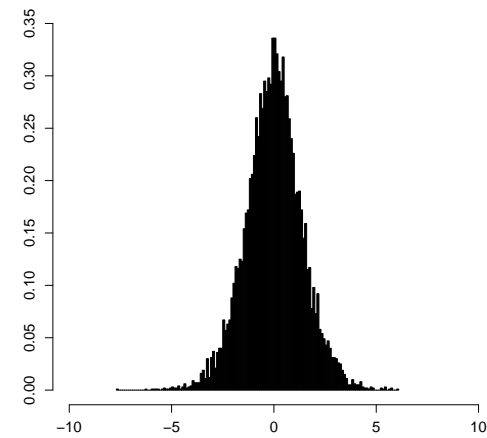
If $\alpha = 0$ and $\sigma = 1$, then $E[\bar{X}_n] = X_0$ with $\text{Var}[\bar{X}_n] \rightarrow 2$ as $n \rightarrow \infty$.



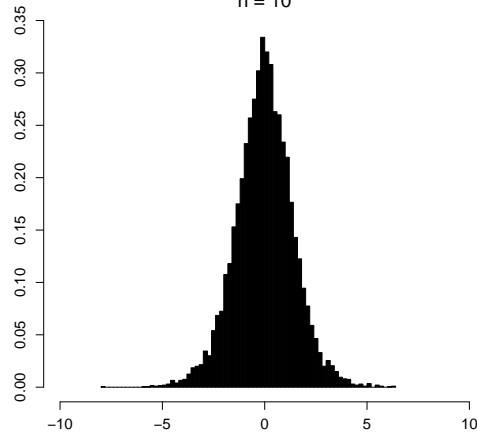
$n = 10$



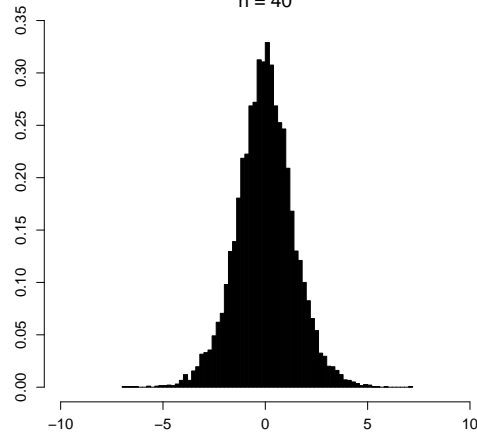
$n = 40$



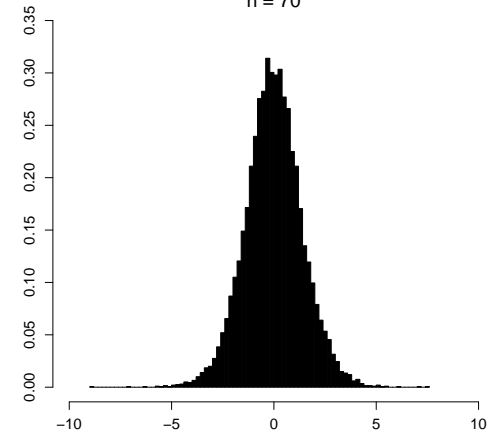
$n = 70$



$n = 110$



$n = 140$



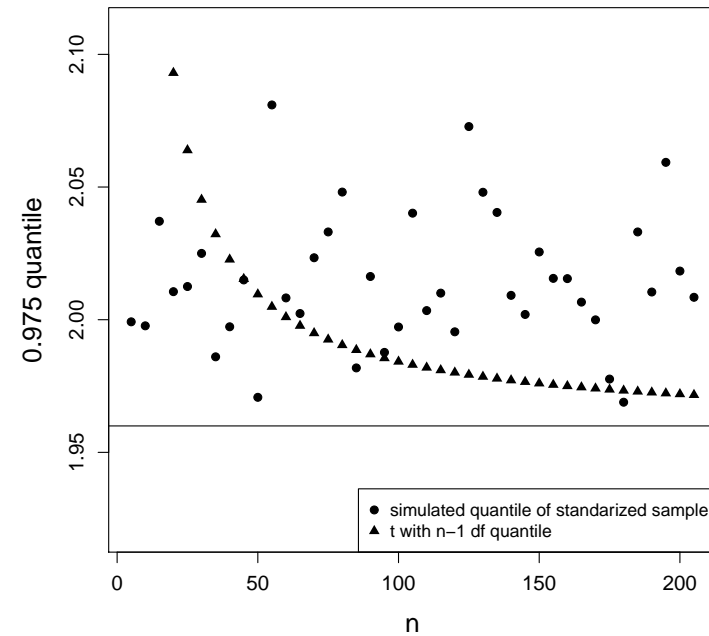
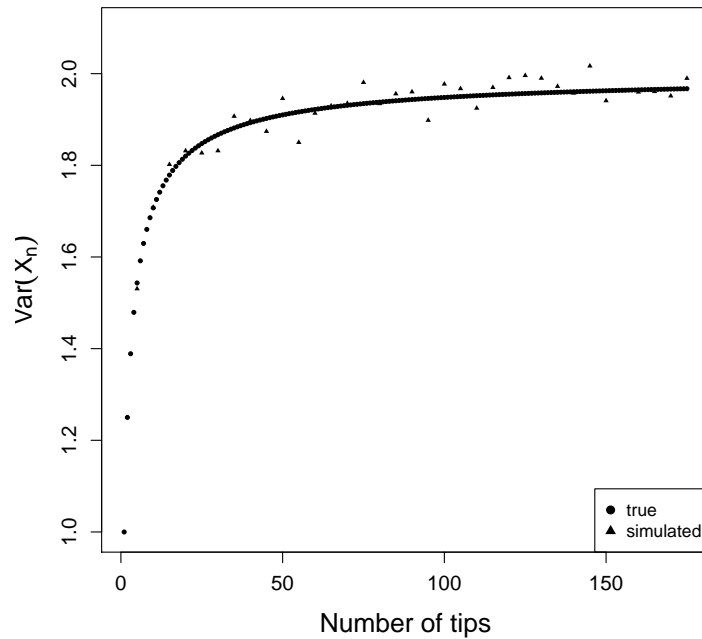
$n = 170$

X_0 ESTIMATE FOR THE YMB-MODEL

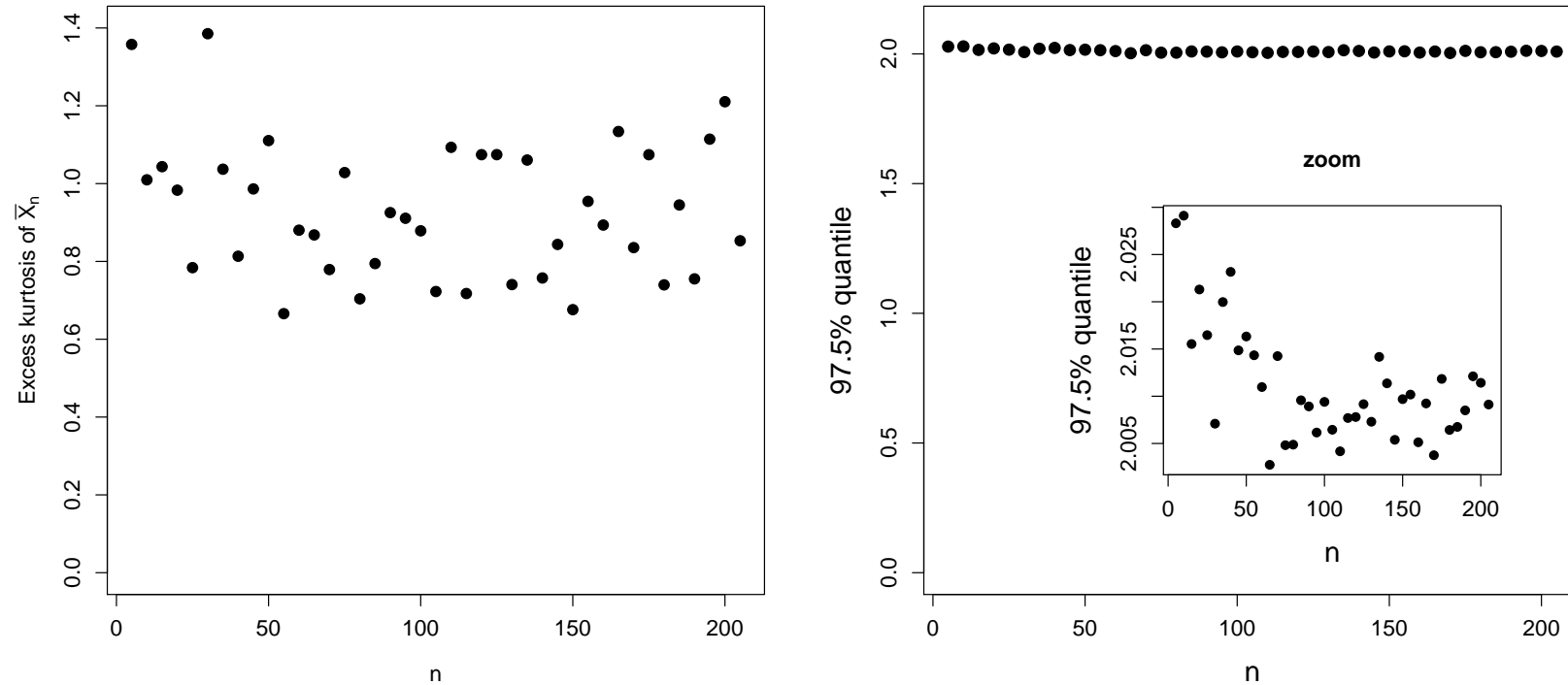
For YBM-model \bar{X}_n is an unbiased but not consistent estimate X_0 .

Assuming σ is known one can think of an approximate 95% CI for X_0

$$\bar{X}_n \pm q_n \frac{\sigma}{\sqrt{n}} \sqrt{2n - a_n}$$



X_0 ESTIMATE FOR THE YMB-MODEL



Asymptotic distribution of \bar{X}_n for the Brownian motion model with $X_0 = 0$ and $\sigma = 1$ is symmetric but not normal.

INTERSPECIES CORRELATION

From

$$\rho_n = 1 - \frac{\sigma_0^2(1 - \mathbb{E}[e^{-2\alpha\tau}])}{\sigma_0^2(1 - \mathbb{E}[e^{-2\alpha T}]) + (X_0 - \theta)^2 \text{Var}[e^{-\alpha T}]}$$

we obtain putting $\delta = \frac{|X_0 - \theta|}{\sigma}$

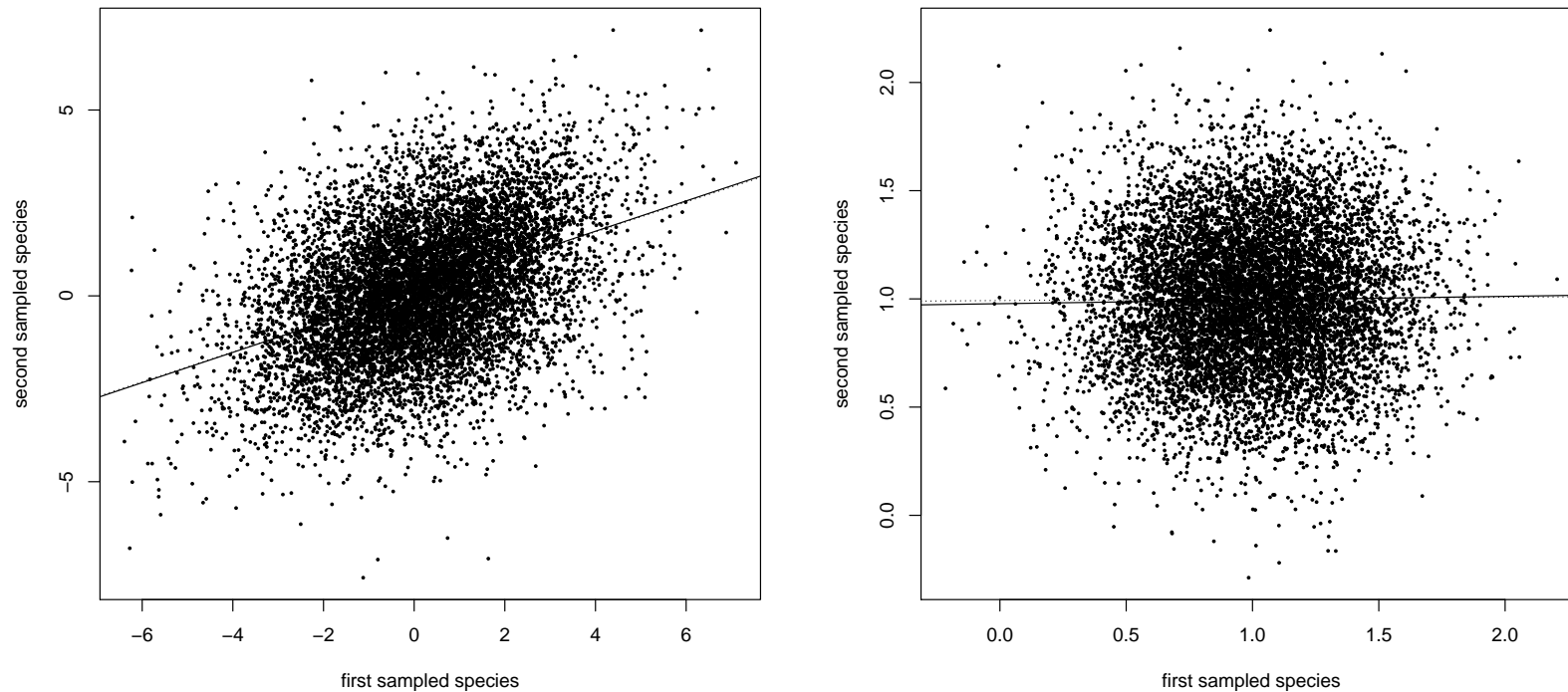
$$\rho_n = 1 - \frac{2\alpha(n-1) + (n+1)((2\alpha+1)b_{n,2\alpha} - 1)}{(n-1)(2\alpha-1)(1 + (2\alpha\delta^2 - 1)b_{n,2\alpha}) - 2\alpha\delta^2 b_{n,\alpha}^2}$$

As $n \rightarrow \infty$

$$\rho_n \sim \begin{cases} 2(\ln n)^{-1}, & \text{if } \alpha = 0 \\ c_{\alpha,\delta} n^{-2\alpha}, & \text{if } 0 < \alpha < 0.5 \\ 2n^{-1} \ln n, & \text{if } \alpha = 0.5 \\ \frac{2}{2\alpha-1} n^{-1}, & \text{if } \alpha > 0.5 \end{cases}$$

$$\rho_n = \frac{\mathbb{E}[T - \tau]}{\mathbb{E}[T]}$$

INTERSPECIES CORRELATION



Simulation results for $n = 30$, $\sigma = 1$, $X_0 = 0$, $\theta = 1$. Thick lines fitted to simulated data are indistinguishable from the predicted (dashed) lines $y = \rho_{30}x + (1 - b_{30,\alpha})(1 - \rho_{30})$. Left: $\alpha = 0.05$. Right: $\alpha = 5$.

RELATED PROJECTS

Ongoing project 1 (with K.Bartoszek)

YOU-model with jumps at speciation events.

Ongoing project 2 (with K.Bartoszek, G.Jones, B.Oxelman):

suppose a tetraploid comes as a hybrid of a pair in a set of n diploids.

The time to the hybridization event is asymptotically exponential.

Assumption 1: conditioned Yule tree for n diploid species with parameter λ . Assumption 2: Poissonian clock for hybridization events with rate β per a pair of diploids.



Acknowledgement. This work was supported by the Swedish Research Council grant 621-2010-5623.

REFERENCES

1. Sagitov S. and Bartoszek K.
Interspecies correlation for neutrally evolving traits. *J. Theor. Biol.* 309 (2012) 11-19
2. Bartoszek K. and Sagitov S.
A phylogenetic confidence interval for the optimal trait value.
Submitted to *J. Math. Biol.* <http://arxiv.org/abs/1207.6488>
3. Graham J., Sagitov S., and Oxelman B.
Statistical inference of allopolyploid species networks in the presence of incomplete lineage sorting.
Submitted to *Syst. Biol.* <http://arxiv.org/abs/1208.3606>
4. Bartoszek K., Jones G., Oxelman B., and Sagitov S.
Time to a single hybridization event in a group of species with unknown ancestral history. In progress.