

# Japanese Joint Statistical Meeting 2018

---

2018/09/09 (Sun.) – 2018/09/13 (Thu.)

## Stochastic Variational Inference of Mixture Models in Phylogenetics

Tung Dang and Hirohisa Kishino

Laboratory of Biometrics and Bioinformatics,  
University of Tokyo

**Corresponding author:**

[dangthanhtung91@vn-bml.com](mailto:dangthanhtung91@vn-bml.com)



2018年度

# 統計関連学会連合大会




Japanese Joint Statistical Meeting

2018年9月9日(日)～13日(木)

## ■E会場 (5333教室) 10:00 - 12:00

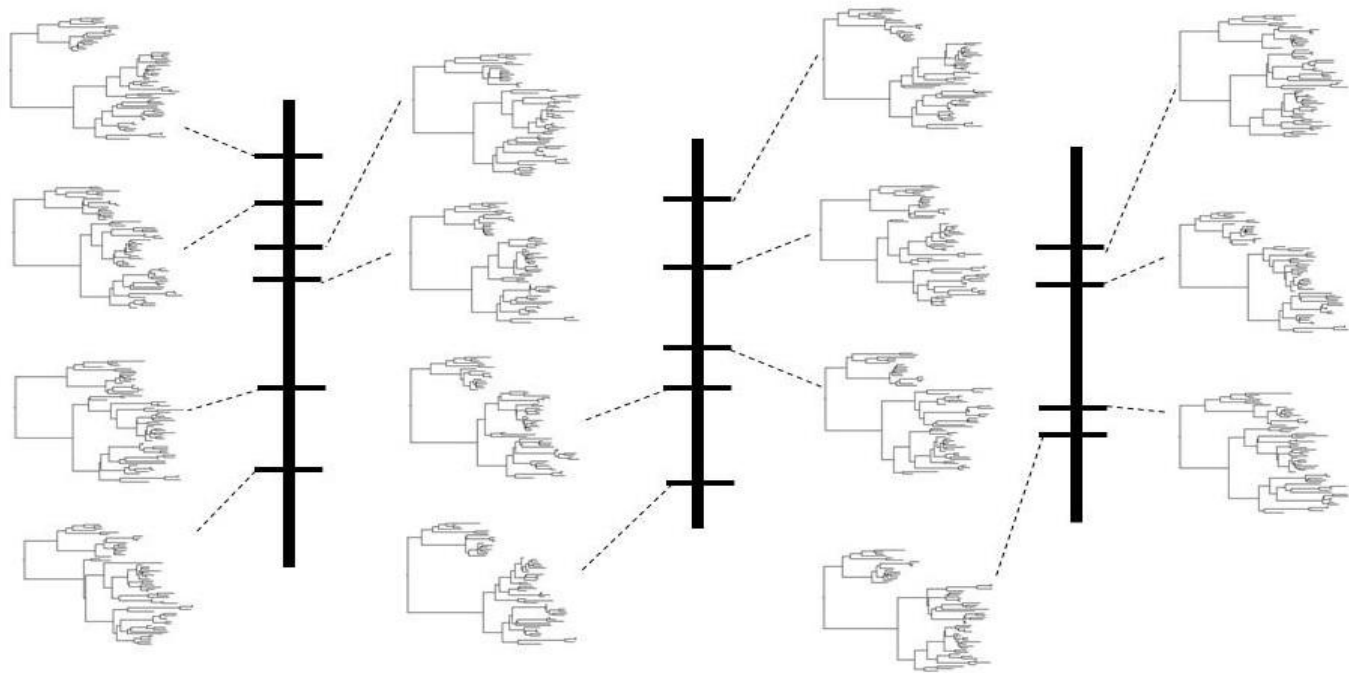
### 【一般講演】 English Session (4): Medical Statistics and Biostatistics

【Chair】 Ryuji Uozumi (Kyoto University)

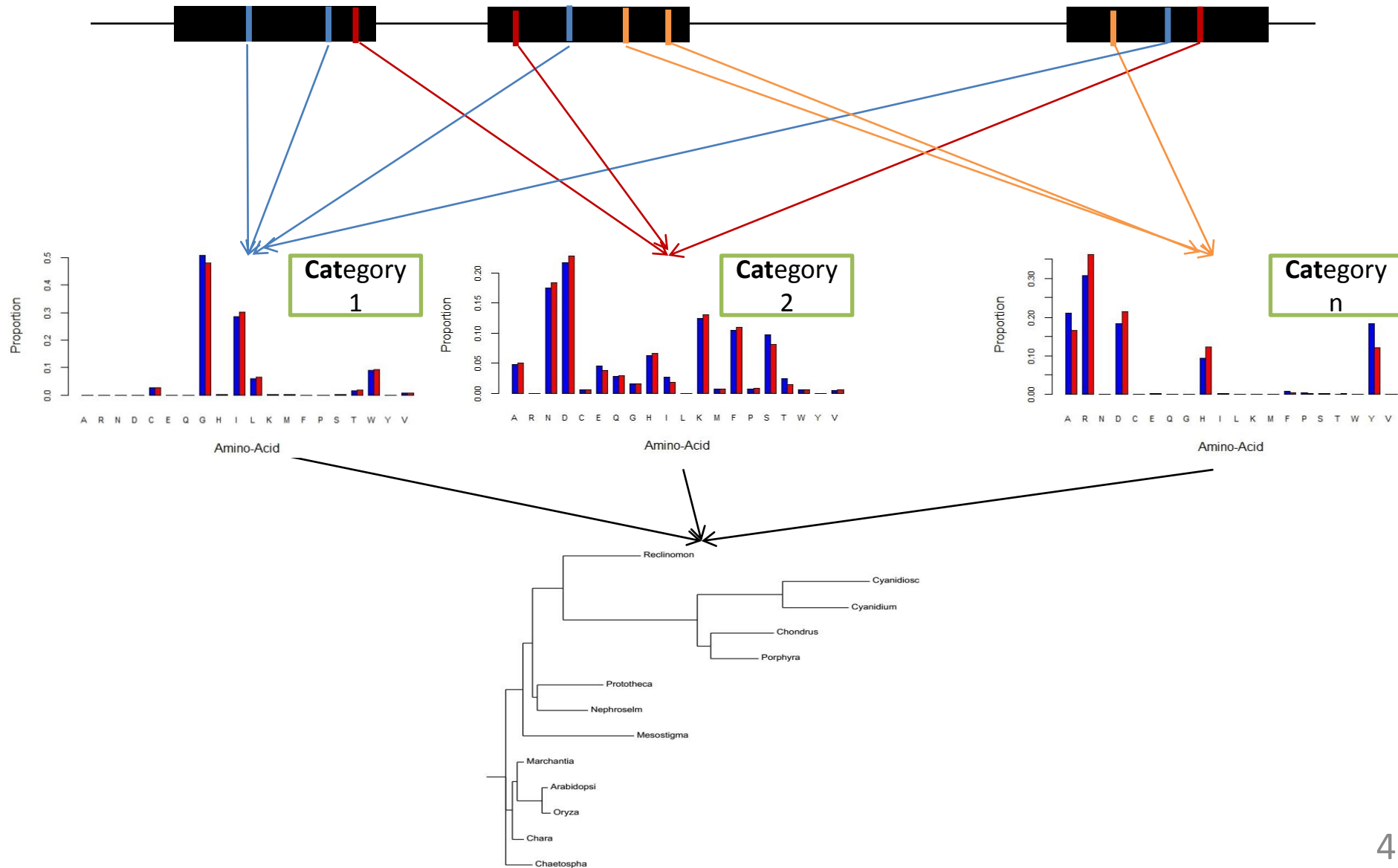
	講演タイトル	発表者 (所属)	共著者 (所属)	報告集
1	Use of external information for assessing efficacy equivalence in biosimilar clinical trials	Ryuji Uozumi (Kyoto University)	Shinjo Yada (A2 Healthcare)	 報告集
2	Bartlett-type corrections for the confidence interval of a treatment effect of the multivariate random effects meta-analysis model via analytical approach	Masahiro Kojima (Kyowa Hakko Kirin)		 報告集
3	The impact of unobserved heterogeneity and competing risks with shared frailty in radiation risk assessment	Kyoji Furukawa (Kurume University)		 報告集
4	Modelling life history under varying temperature conditions	Hideyasu Shimadzu (Loughborough University)		
5	Sufficient dimension reduction via random-partitions for the large-p-small-n problem	Hung Hung (National Taiwan University)		
6	Stochastic Variational Inference of Mixture Models in Phylogenetics	Tung Dang (University of Tokyo)	Hirohisa Kishino (University of Tokyo)	

# Problem of Phylogenomics

- The pattern of molecular evolution varies among gene sites and genes in a genome.
- By taking into account the complex heterogeneity of evolutionary processes among sites in a genome, Bayesian infinite mixture models of genomic evolution enable robust phylogenetic inference.



# CAT model (Lartillot et al 2004)



# The likelihood of CAT model

$$p(\Xi_{ij} \mid \pi_{ka}, r_i, l_j) = \prod_{k=1}^{\infty} \prod_{a=1}^{20} \left( \pi_{ka} p(n_{ij} \mid r_i, l_j) \right)^{I[Z_i=k]}$$

$$p(\Xi_{ij} \mid \pi_{ka}, r_i, l_j) = \left( \prod_{k=1}^{\infty} \prod_{a=1}^{20} \pi_{ka}^{w_{ka}} \right)^{I[Z_i=k]} \left( \prod_i r_i^{\sum_j n_{ij}} \right) \left( \prod_j l_j^{\sum_i n_{ij}} \right) \left( \prod_{ij} \frac{e^{-r_i l_j}}{n_{ij}!} \right)$$

$l_j$  : branch lengths

$r_i$  : rate of substitution

$\pi_{ka}$  : equilibrium frequency profile

A stick-breaking construction

$V_k$  : unit length of the stick

$I[Z_i=k]$  : allocation variable

$\Xi_{ij}$  : substitution mapping data

$n_{ij}$  : number of substitutions

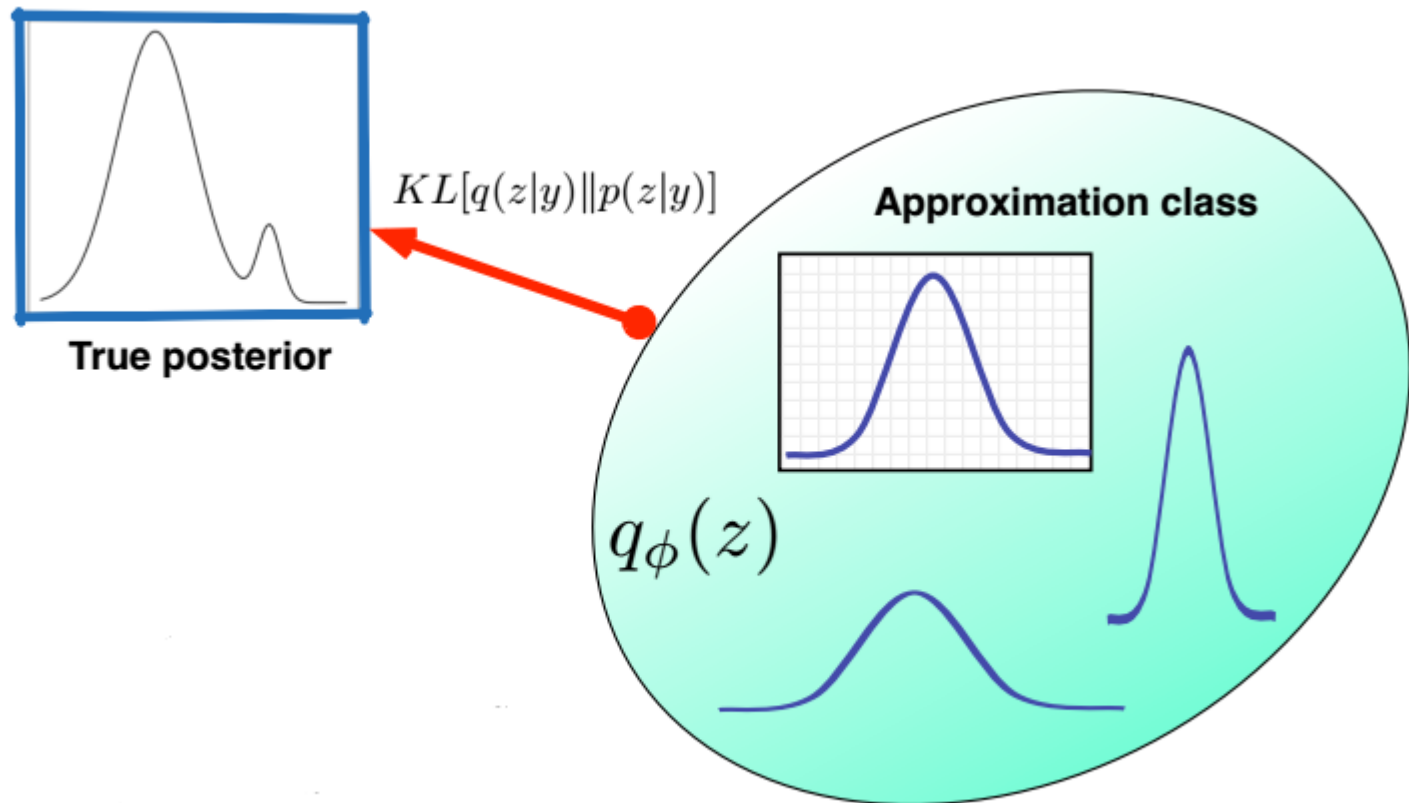
$w_{ka}$  : number of type of substitutions

# MCMC becomes infeasible for large datasets

- The computational burden of MCMC is prohibitive for large data sets.
- Even well-designed sampling schemes needs a large sample to achieve convergence.
- Diagnosis of convergence is difficult for high dimensional parameter space.

# Variational Method

approximate complicated densities by a simpler class of densities, called **variational distribution**.



# Evidence Lower Bound (ELBO)

$$\text{KL}[q(\boldsymbol{\theta}; \boldsymbol{\Theta}) \parallel p(\boldsymbol{\theta}|\mathbf{X})]$$

$$= \mathbb{E}_q[\log q(\boldsymbol{\theta}; \boldsymbol{\Theta})] - \mathbb{E}_q[\log p(\mathbf{X}, \boldsymbol{\theta})] + \log p(\mathbf{X})$$

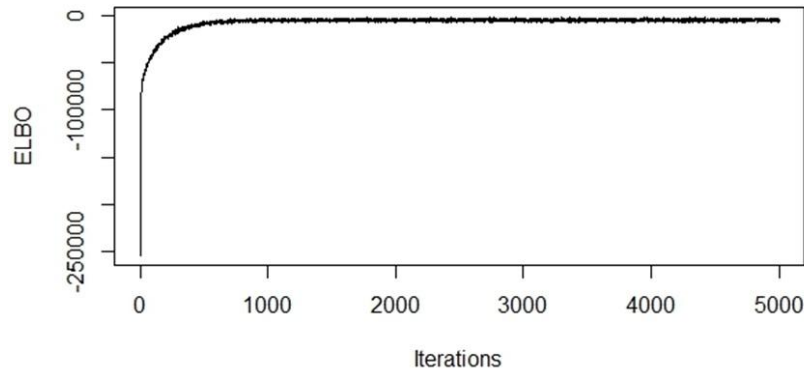
-ELBO: computable

Computational headache  
But constant !!



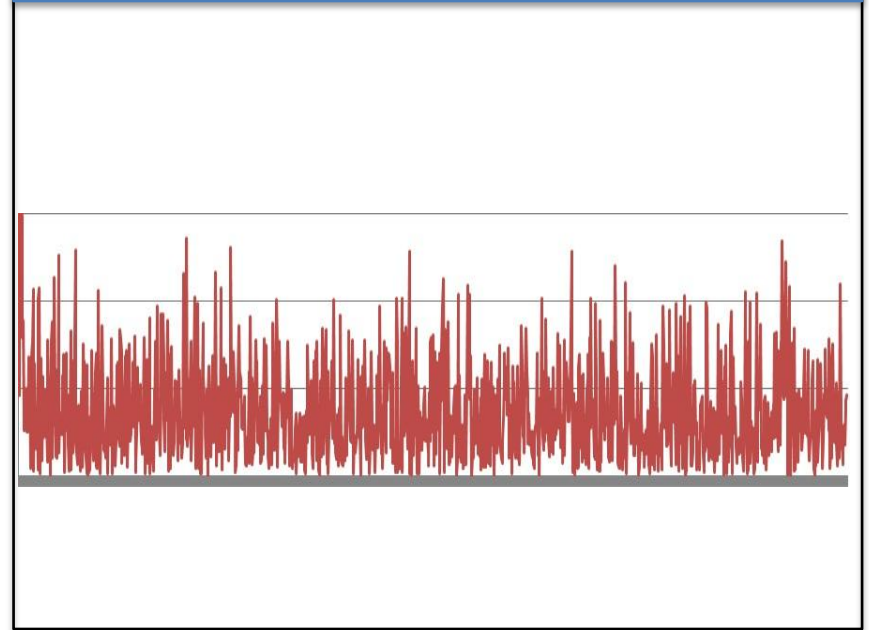
# Variational Inference and MCMC

**VI**



Maximizing ELBO

**MCMC**



Monte Carlo simulation with  
proposal and acceptance steps

# Example: *Variational Inference of Bayesian Ridge Regression*

- The likelihood is

$$p(t|x, w, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | w^T x_n, \beta^{-1})$$

- The regression coefficient  $p(w|\alpha) = \mathcal{N}(w|0, \alpha^{-1}I)$

- The extent of shrinkage  $p(\alpha) = \text{Gam}(\alpha|a_0, b_0)$

- The variational mean-field representation

$$q(w, \alpha) = q(w) q(\alpha)$$

- The practical variational distributions

$$\begin{aligned} q(w) &= \mathcal{N}(w|m_N, S_N) \\ q(\alpha) &= \text{Gam}(\alpha|a_N, b_N) \end{aligned}$$

# Example: *Variational Inference of Bayesian Ridge Regression*

- The ELBO for Bayesian ridge regression model.

$$\mathcal{L}[q(w, \alpha | m_N, S_N, a_N, b_N)] = \mathbb{E}_q[\log p(t|w)] + \mathbb{E}_q[\log p(w|\alpha)] + \mathbb{E}_q[\log p(\alpha)] \\ - \mathbb{E}_q[\log q(w|m_N, S_N)] - \mathbb{E}_q[\log q(\alpha|a_N, b_N)]$$

- Update  $m_N, S_N$  : Given  $a_N = a_N^0, b_N = b_N^0, m_N = m_N^0, S_N = S_N^0$

$$m_N = \beta S_N^0 X^T t$$

$$S_N = (\mathbb{E}_q[\alpha] + \beta X^T X)^{-1} I = \left( \frac{a_N^0}{b_N^0} + \beta X^T X \right)^{-1} I$$

- Update  $a_N, b_N$  : Given  $a_N = a_N^0, b_N = b_N^0, m_N = m_N^0, S_N = S_N^0$

$$a_N = a_0 + \frac{M}{2}$$

$$b_N = b_0 + \frac{1}{2} \mathbb{E}_q[w^T w] = b_0 + \frac{1}{2} (m_N^{0T} m_N^0 + S_N^0)$$

# Example: *Variational Inference of Dirichlet mixture model*

- The PDF of a DMM can be represented

$$f(\mathbf{X}; \mathbf{\Pi}, \mathbf{U}) = \prod_{n=1}^N \sum_{i=1}^I \pi_i \text{Dir}(\mathbf{x}_n; \mathbf{u}_i), \quad \pi_i > 0, \quad \sum_{i=1}^I \pi_i = 1$$

- The mixture weights  $\mathbf{\Pi} = [\pi_1, \dots, \pi_I]^T$
- The parameter matrix  $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_I]$
- The joint distribution of the observation  $\mathbf{X}$  and all the latent variables

$$\begin{aligned} f(\mathbf{X}, \mathbf{Z}) &= f(\mathbf{X}, \mathbf{U}, \mathbf{\Pi}, \mathbf{Z}) \\ &= f(\mathbf{X}|\mathbf{Z}, \mathbf{U})f(\mathbf{Z}|\mathbf{\Pi})f(\mathbf{\Pi})f(\mathbf{U}) \\ &= \prod_{n=1}^N \prod_{i=1}^I \left[ \pi_i \frac{\Gamma(\sum_{k=1}^{K+1} u_{ki})}{\prod_{k=1}^{K+1} \Gamma(u_{ki})} \prod_{k=1}^{K+1} x_{kn}^{u_{ki}-1} \right]^{z_{ni}} \times \frac{\Gamma(\sum_{i=1}^I c_{i_0})}{\prod_{i=1}^I \Gamma(c_{i_0})} \prod_{i=1}^I \pi_i^{c_{i_0}-1} \\ &\quad \times \prod_{i=1}^I \prod_{k=1}^{K+1} \frac{\alpha_{ki_0}^{\mu_{ki_0}}}{\Gamma(\mu_{ki_0})} u_{ki}^{\mu_{ki_0}-1} e^{-\alpha_{ki_0} u_{ki}}. \end{aligned}$$

- The ELBO for Dirichlet Mixture model

$$\mathcal{L} = \mathbb{E}_{\mathbf{Z}}[\ln f(\mathbf{X}, \mathbf{Z})] - \mathbb{E}_{\mathbf{Z}}[\ln f(\mathbf{Z})]$$

# Example: *Variational Inference of Dirichlet mixture model*

- The optimal solution to the posterior distribution of  $u_{ki}$  is

$$\begin{aligned} \ln f^*(u_{ki}; \mu_{ki}, \alpha_{ki}) &= \sum_{n=1}^N \mathbb{E}[Z_{ni}] \times \mathbb{E}_{\mathcal{Z} \setminus u_{ki}} \left[ \ln \frac{\Gamma(\sum_{k=1}^{K+1} u_{ki})}{\prod_{k=1}^{K+1} \Gamma(u_{ki})} \right] \\ &\quad + u_{ki} \sum_{n=1}^N \mathbb{E}[Z_{ni}] \ln x_{kn} + (\mu_{ki_0} - 1) \ln u_{ki} - \alpha_{ki_0} u_{ki} + \text{const.} \end{aligned}$$

- The optimal solution to the posterior distribution

$$\ln f^*(\pi_i; c_i) = \mathbb{E}_{\mathcal{Z} \setminus \pi_i} [\ln f(\mathbf{X}, \mathcal{Z})] = \ln \pi_i \times \sum_{n=1}^N \mathbb{E}[Z_{ni}] + \ln \pi_i (c_{i_0} - 1) + \text{const.}$$

- For the variable  $z_{ni}$ , the optimal approximation to the posterior distribution is

$$\begin{aligned} \ln f^*(z_{ni}) &= \mathbb{E}_{\mathcal{Z} \setminus z_{ni}} [\ln f(\mathbf{X}, \mathcal{Z})] \\ &= z_{ni} \times \left[ \mathbb{E}[\ln \pi_i] + \sum_{k=1}^{K+1} (u_{ki} - 1) \ln x_{kn} \right] \\ &\quad + z_{ni} \times \mathbb{E} \left[ \ln \frac{\Gamma(\sum_{k=1}^{K+1} u_{ki})}{\prod_{k=1}^{K+1} \Gamma(u_{ki})} \right] + \text{const.} \end{aligned}$$

# Example: Variational Inference of Dirichlet mixture model

## Definition [\[ edit \]](#)

The Taylor series of a [real](#) or [complex-valued](#) function  $f(x)$  that is infinitely differentiable at a [real](#) or [complex number](#)  $a$  is the power series

$$f(a) + \frac{f'(a)}{1!}(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \frac{f'''(a)}{3!}(x-a)^3 + \dots,$$

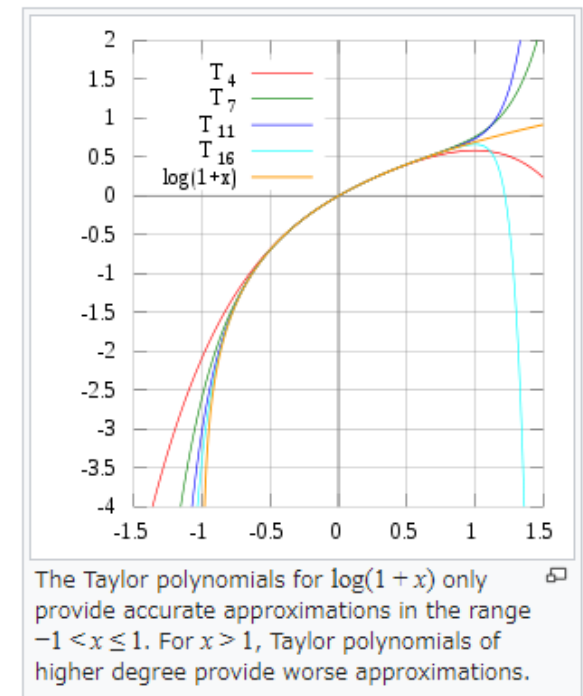
which can be written in the more compact [sigma notation](#) as

$$\sum_{n=0}^{\infty} \frac{f^{(n)}(a)}{n!}(x-a)^n,$$

- A lower-bound approximation can be obtained as

$$\begin{aligned} \ln \frac{\Gamma(\sum_{k=1}^{K+1} u_{ki})}{\prod_{k=1}^{K+1} \Gamma(u_{ki})} &\geq \ln \frac{\Gamma(\sum_{k=1}^{K+1} \bar{u}_{ki})}{\prod_{k=1}^{K+1} \Gamma(\bar{u}_{ki})} \\ &+ \sum_{k=1}^{K+1} \left[ \psi \left( \sum_{m=1}^k \bar{u}_{mi} + \sum_{l=k+1}^{K+1} u_{li} \right) - \psi(\bar{u}_{ki}) \right] \bar{u}_{ki} (\ln u_{ki} - \ln \bar{u}_{ki}), \end{aligned}$$

$\psi(\cdot)$  is the digamma function defined as  $\psi(x) = \partial \ln \Gamma(x) / \partial x$ .



# Example: Variational Inference of Dirichlet mixture model

- The optimal solution to the posterior distribution of  $u_{ki}$  is

$$\begin{aligned}
 \ln f^*(u_{ki}; \mu_{ki}, \alpha_{ki}) &\approx \mathbb{E}_{Z \setminus u_{ki}} [\ln \tilde{f}(\mathbf{X}, Z)] \\
 &= \sum_{n=1}^N \mathbb{E}[Z_{ni}] \times \mathbb{E}_{Z \setminus u_{ki}} \left[ \ln \frac{\Gamma(\sum_{k=1}^{K+1} u_{ki})}{\prod_{k=1}^{K+1} \Gamma(u_{ki})} \right] \\
 &\quad + u_{ki} \sum_{n=1}^N \mathbb{E}[Z_{ni}] \ln x_{kn} + (\mu_{ki_0} - 1) \ln u_{ki} - \alpha_{ki_0} u_{ki} + \text{const.}
 \end{aligned}$$

$$\begin{aligned}
 &= \sum_{n=1}^N \mathbb{E}[Z_{ni}] \times \left[ \psi \left( \sum_{k=1}^{K+1} \bar{u}_{ki} \right) - \psi(\bar{u}_{ki}) \right] \bar{u}_{ki} \ln u_{ki} - \ln \bar{u}_{ki} \\
 &\quad + u_{ki} \sum_{n=1}^N \mathbb{E}[Z_{ni}] \ln x_{kn} + (\mu_{ki_0} - 1) \ln u_{ki} - \alpha_{ki_0} u_{ki} + \text{const.,}
 \end{aligned}$$

- For the variable  $z_{ni}$ , the optimal approximation to the posterior distribution is

$$\begin{aligned}
 \ln f^*(z_{ni}) &= \mathbb{E}_{Z \setminus z_{ni}} [\ln f(\mathbf{X}, Z)] \\
 &= z_{ni} \times \left[ \mathbb{E}[\ln \pi_i] + \sum_{k=1}^{K+1} (u_{ki} - 1) \ln x_{kn} \right] \\
 &\quad + z_{ni} \times \mathbb{E} \left[ \ln \frac{\Gamma(\sum_{k=1}^{K+1} u_{ki})}{\prod_{k=1}^{K+1} \Gamma(u_{ki})} \right] + \text{const.}
 \end{aligned}$$

$$\begin{aligned}
 &\approx \mathbb{E}_{Z \setminus z_{ni}} [\ln f(\mathbf{X}, Z)] \\
 &= z_{ni} \times \left[ \mathbb{E}[\ln \pi_i] + \sum_{k=1}^{K+1} (u_{ki} - 1) \ln x_{kn} \right] \\
 &\quad + z_{ni} \times \underbrace{\sum_{k=1}^{K+1} \left[ \psi \left( \sum_{k=1}^{K+1} \bar{u}_{ki} \right) - \psi(\bar{u}_{ki}) \right] \bar{u}_{ki} \mathbb{E}_{u_{ki}} [\ln u_{ki}] - \ln \bar{u}_{ki}}_{P_i} + \text{const.}
 \end{aligned}$$

# Example: *Variational Inference of Dirichlet mixture model*

**Algorithm 1.** Variational DMM.

**Input:** observation  $\mathbf{X}$ , number of mixture components  $I$   
Initialize  $\alpha_{ki_0}, \mu_{ki_0}, c_{i_0}$ , for  $i=1, \dots, I, k=1, \dots, K+1$ <sup>3</sup>;

**repeat**

  for each  $k, i$

$$\alpha_{ki}^* = \alpha_{ki_0} - \sum_{n=1}^N \mathbb{E}[Z_{ni}] \ln x_{kn}$$

$$\mu_{ki}^* = \mu_{ki_0} + \sum_{n=1}^N \mathbb{E}[Z_{ni}] \bar{u}_{ki} [\psi(\sum_{k=1}^{K+1} \bar{u}_{ki}) - \psi(\bar{u}_{ki})]$$

$$c_i^* = c_{i_0} + \sum_{n=1}^N \mathbb{E}[Z_{ni}]$$

**until** stop criteria are reached.

**Output:** the optimal hyperparameters  $\alpha_{ki}^*, \mu_{ki}^*, c_i^*$ .

(The quantities  $\bar{u}_{ki}$  and  $\mathbb{E}[Z_{ni}]$  are calculated

$$\bar{u}_{ki} = \frac{\mu_{ki}}{\alpha_{ki}}, \quad \mathbb{E}[Z_{ni}] = \frac{\rho_{ni}}{\sum_{i=1}^I \rho_{ni}},$$

$$\ln \rho_{ni} = \psi(c_i) - \psi(\mathbf{c}^T \mathbf{1}_I) + P_i + (\mathbf{u}_i - 1)^T \ln \mathbf{x}_n$$



# Variational distribution for CAT model

- The mean-field variational representation

$$\begin{aligned}
 & q(\Xi, z, V, \pi, l, r | \Theta) \\
 = & \prod_j q(l_j | \gamma_j, \gamma'_j) \times \prod_i q(r_i | \zeta_i, \zeta'_i) \\
 & \times \prod_{k=1}^{K_{\max}} \prod_{a=1}^{20} q(\pi_a^k | \lambda_a^k) \times \prod_{k=1}^{K_{\max}} q(V_k | \vartheta_k, \vartheta'_k) \\
 & \times \prod_i \prod_{k=1}^{K_{\max}} q(z_i^k | \phi_i^k) \times \prod_{ij} q(n_{ij} | \omega_{ij}) \\
 & \times \prod_{k=1}^{K_{\max}} \prod_{a=1}^{20} q(w_a^k | \iota_a^k)
 \end{aligned}$$

$$q(l_j | \gamma_j, \gamma'_j) = \text{Gamma}(l_j | \gamma_j, \gamma'_j)$$

$$q(r_i | \zeta_i, \zeta'_i) = \text{Gamma}(r_i | \zeta_i, \zeta'_i)$$

$$q(\pi_a^k | \lambda_a^k) = \text{Dirichlet}(\pi_a^k | \lambda_a^k)$$

$$q(V_k | \vartheta_k, \vartheta'_k) = \text{Beta}(V_k | \vartheta_k, \vartheta'_k)$$

$$q(z_i^k | \phi_i^k) = \text{Multinomial}(z_i^k | \phi_i^k)$$

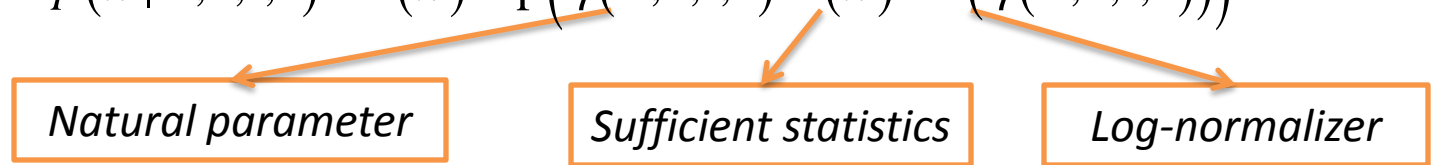
$$q(n_{ij} | \omega_{ij}) = \text{Poisson}(n_{ij} | \omega_{ij})$$

$$q(w_a^k | \iota_a^k) = \text{Multinomial}(w_a^k | \iota_a^k)$$

# Stochastic Variational Inference

- These distributions are in the exponential family

$$p(\pi | \Xi, z, l, r) = h(\pi) \exp\left(\eta(\Xi, z, l, r)^T t(\pi) - a(\eta(\Xi, z, l, r))\right)$$



- The variational parameters  $q(\pi | \lambda) = h(\pi) \exp(\lambda^T t(\pi) - a(\lambda))$

- The natural gradient of the ELBO [Amari, 1998]

$$\nabla_{\lambda}^{\text{nat}} ELBO(\lambda) = \left( \text{prior} + NE_q \left[ t(\Xi, z, l, r) \right] \right) - \lambda$$

- It is good for stochastic optimization.
  - Its expectation is the exact gradient (*unbiased*).
  - It only depends on optimized parameters of one data point (*cheap*).

# Variational approximation for model of nucleotide substitution

- Evidence Lower Bound (ELBO)**

$$\begin{aligned}
 \mathcal{L} = & \sum_{k=1}^{K_{max}} \mathbb{E}_q [\log (p (V_k | 1, \kappa))] - \sum_{k=1}^{K_{max}} \mathbb{E}_q \left[ \log \left( q \left( V_k | \vartheta_k, \vartheta'_k \right) \right) \right] \\
 & + \sum_i \sum_{k=1}^{K_{max}} \mathbb{E}_q [\log p(Z_i^k | V_1, V_2, \dots, V_{K_{max}})] \\
 & - \sum_i \sum_{k=1}^{K_{max}} \mathbb{E}_q [\log (q (Z_i^k | \phi_i^k))] \\
 & + \sum_{k=1}^{K_{max}} \sum_{a=1}^{20} \mathbb{E}_q [\log p(\pi_a^k | v_a)] - \sum_{k=1}^{K_{max}} \sum_{a=1}^{20} \mathbb{E}_q [\log (q (\pi_a^k | \lambda_a^k))] \\
 & + \sum_j \mathbb{E}_q [\log (p (l_j | 1, \beta))] - \sum_j \mathbb{E}_q \left[ \log \left( q \left( l_j | \gamma_j, \gamma'_j \right) \right) \right] \\
 & + \sum_i \mathbb{E}_q [\log (p (r_i | \alpha, \alpha))] - \sum_i \mathbb{E}_q \left[ \log \left( q \left( r_i | \zeta_i, \zeta'_i \right) \right) \right] \\
 & + \int \int \int \sum_{\pi} \sum_r \sum_l \sum_z \sum_n q(z) q(n) q(\pi) q(r) q(l) \log (p(n | D, z, \pi, r, l)) d(\pi) d(r) d(l) \\
 & - \sum_n q(n) \log (q(n)) \\
 & + \int \int \int \sum_{\pi} \sum_r \sum_l \sum_z \sum_w q(z) q(w) q(\pi) q(r) q(l) \log (p(w | D, z, \pi, r, l)) d(\pi) d(r) d(l) \\
 & - \sum_w q(w) \log (q(w))
 \end{aligned}$$

Calculations of these integrals are **analytically intractable**

# Variational approximation for model of nucleotide substitution

- We consider a first-order Taylor expansion to preserve a bound, intractable expectations are avoided

The product of  $r_i l_j$  is denoted  $t_{ij}$ , firstly, we consider a first-order Taylor expansion of  $\log(p(a, b|\pi, t_{ij}))$  for  $\pi_a$  and  $t_{ij}$  at  $\pi'_a = \frac{v_a}{\sum_{a'=1}^{20} v_{a'}}$  and  $t'_{ij} = \frac{1}{\beta} \frac{\alpha}{\alpha} = \frac{1}{\beta}$

$$\begin{aligned} \log[e^{-t_{ij}} + (1 - e^{-t_{ij}}) \pi_a] &\approx \log[e^{-t'_{ij}} + (1 - e^{-t'_{ij}}) \pi'_a] \\ &+ \frac{\partial \log[e^{-t_{ij}} + (1 - e^{-t_{ij}}) \pi'_a]}{\partial t_{ij}} (t_{ij} - t'_{ij}) \\ &+ \frac{\partial \log[e^{-t'_{ij}} + (1 - e^{-t'_{ij}}) \pi_a]}{\partial \pi_a} (\pi_a - \pi'_a) \end{aligned} \quad (40)$$

# Stochastic Variational Inference for CAT model

- Iterates  $t$  times to update local variational parameters  $\Theta_l = \{\mathcal{G}, \mathcal{G}', \phi\}$  of local variables  $\Phi_l = \{V, z\}$  based on mapping data

$$\Theta_l^* = E_{\Theta_g} \left\{ \eta[\Phi_g, \Xi] \right\}$$

- By using the natural gradient, iterates  $t$  times to update global variational parameters  $\Theta_g = \{\gamma, \gamma', \zeta, \zeta', \lambda, \omega, \iota\}$  of global variables  $\Phi_g = \{\Xi, \pi, l, r\}$  based on mapping data with step size

$$\widehat{\nabla_{\Theta_g} \mathcal{L}} = \text{prior} + N \{ E_{\Theta_l} [t(\Phi_n, \Xi_n), 1] \} - \Theta_g$$

$$\Theta_g^{(t)} = \Theta_g^{(t-1)} + \rho_t \widehat{\nabla_{\Theta_g} \mathcal{L}}$$

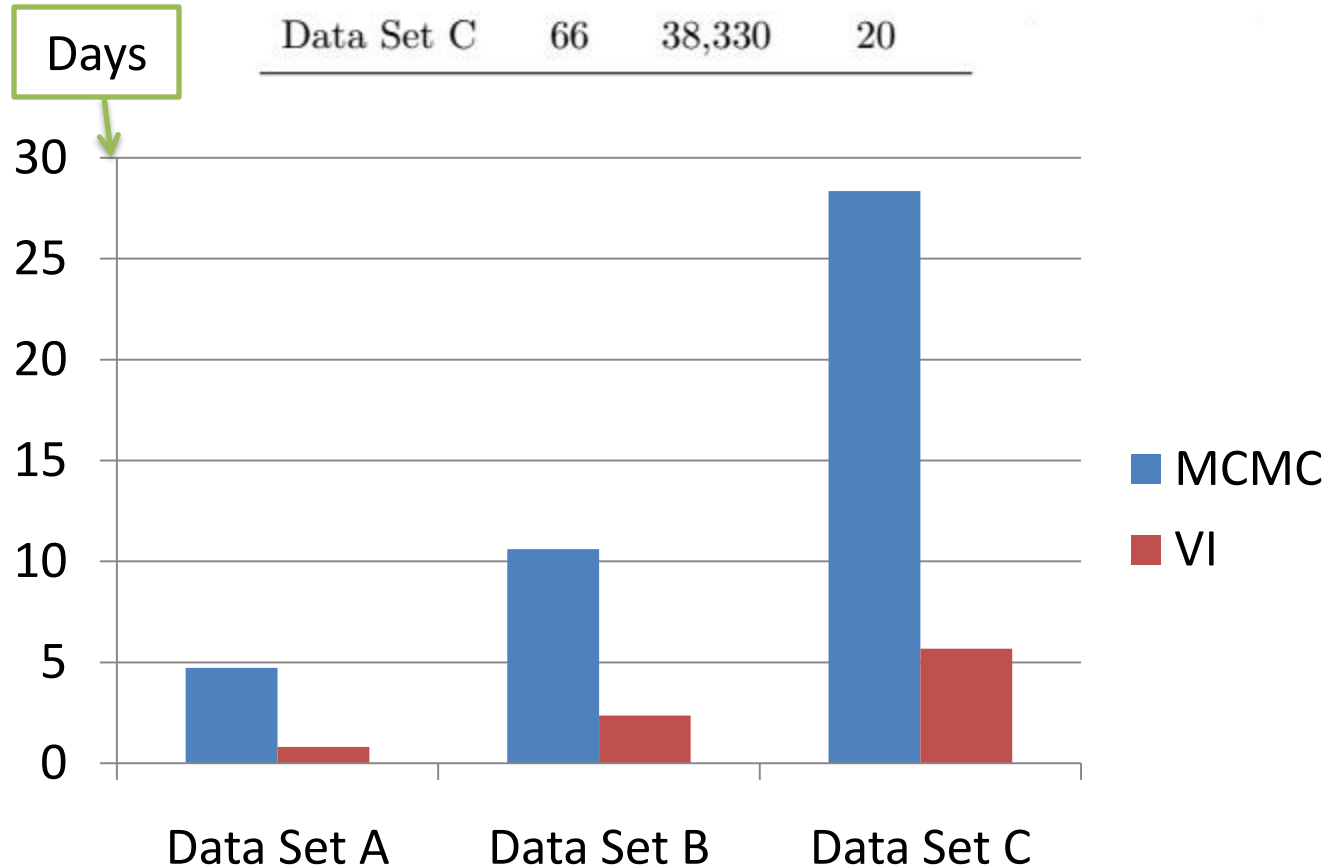
# Checking the performance

- the mitochondrial data set which consisted of 33 proteins, a total of 6,622 amino acid positions, from 13 species (data set A).
- the plastid data set which consisted of 50 plastidencoded proteins, a total of 10,137 amino acid positions, from 28 species (data set B).
- mitochondrial protein sequences, a large alignment from EST and genome data, which consists of 197 genes, a total of 38,330 amino-acid positions from 66 species (data set C).

# CPU time

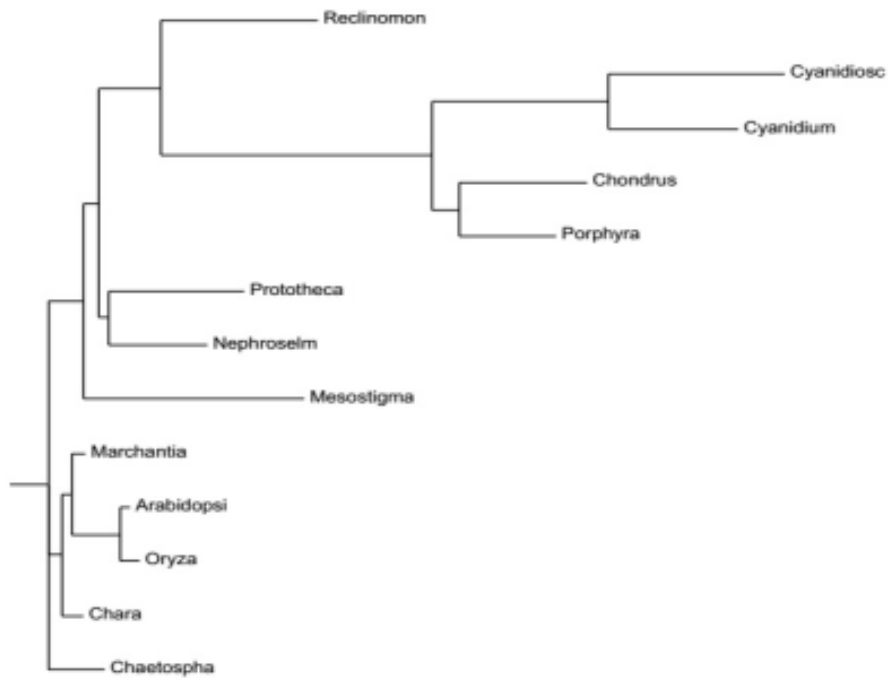
Run times of variational inference and MCMC algorithms on real data

Data set	Taxa	Sites	States
Data Set A	13	6,622	20
Data Set B	28	10,137	20
Data Set C	66	38,330	20

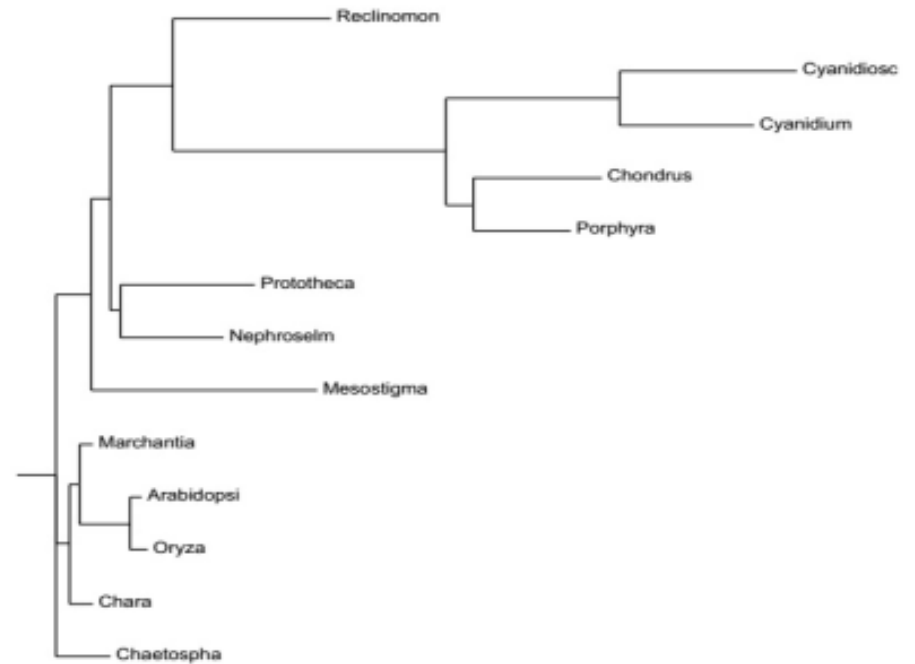


# Reconstructed tree

MCMC



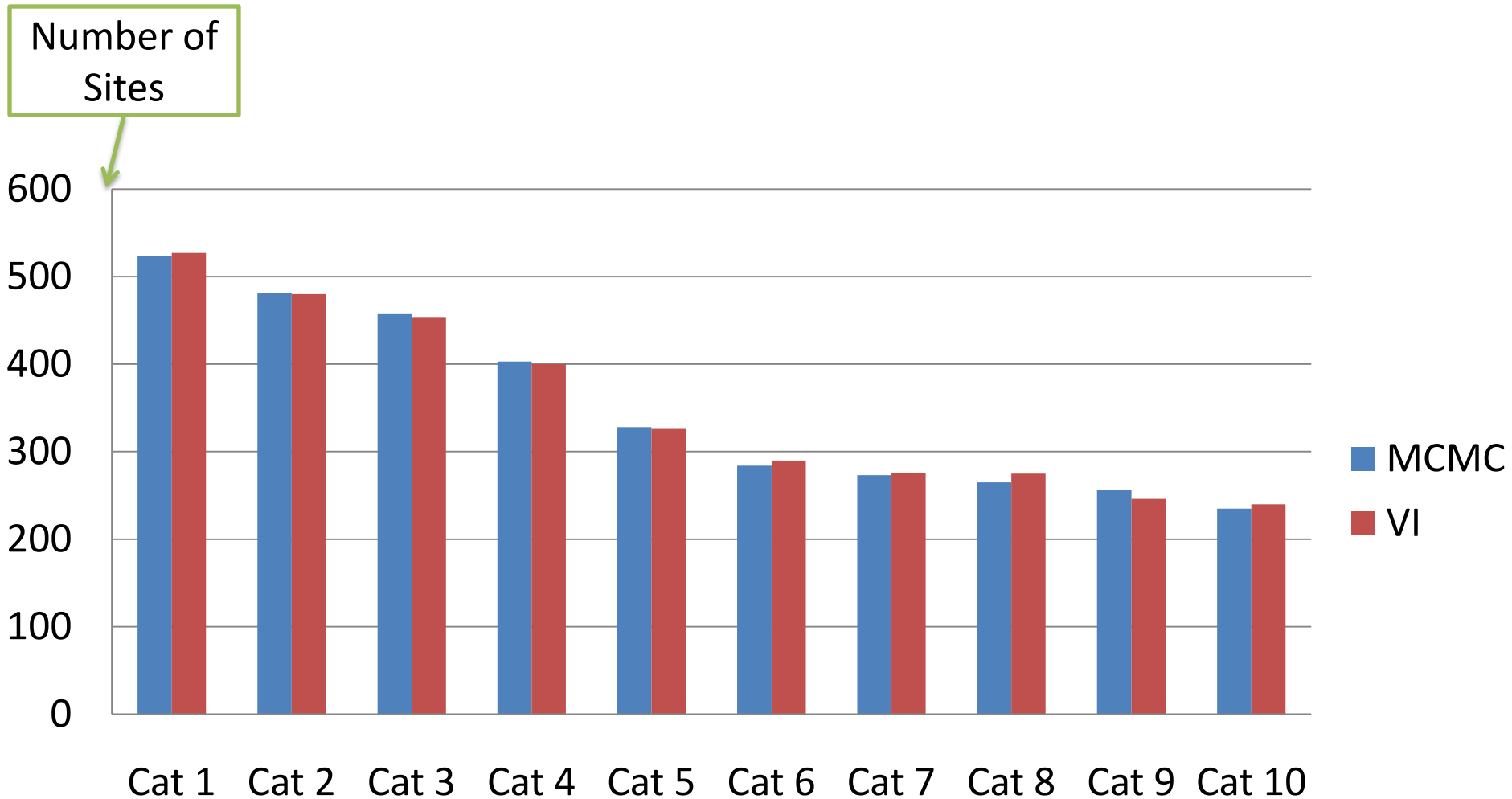
VI



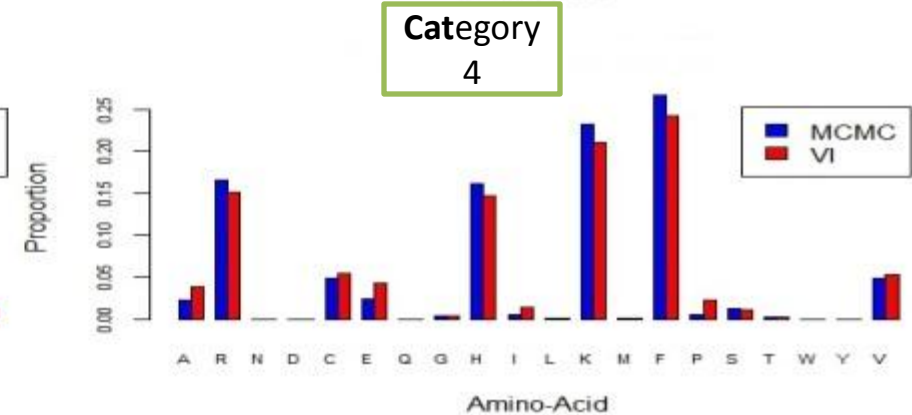
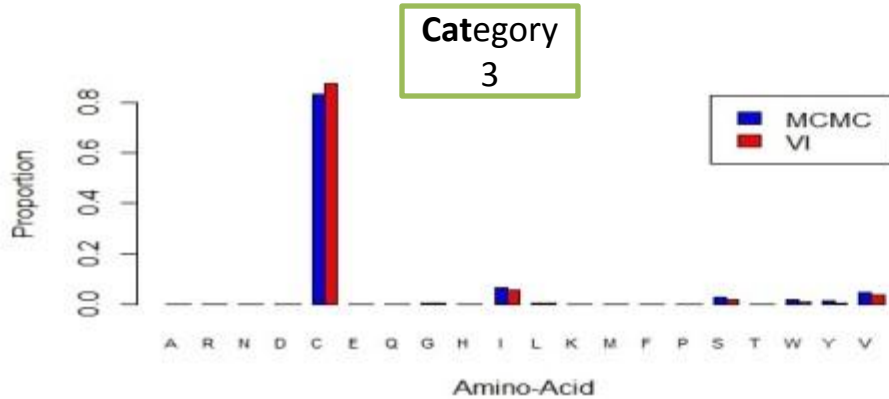
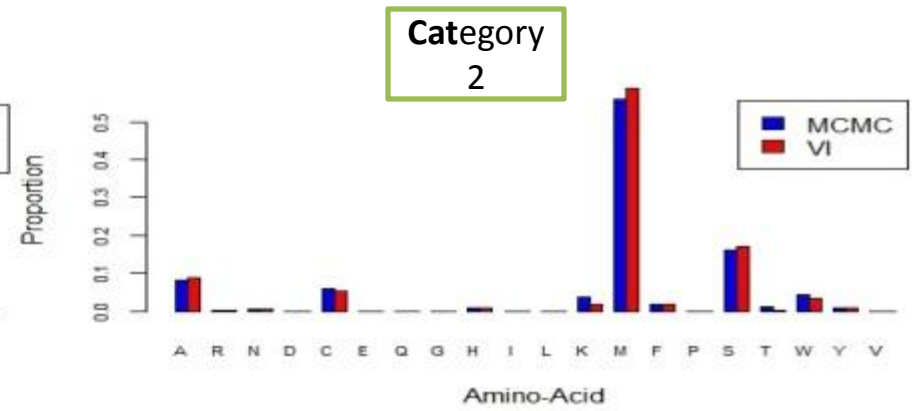
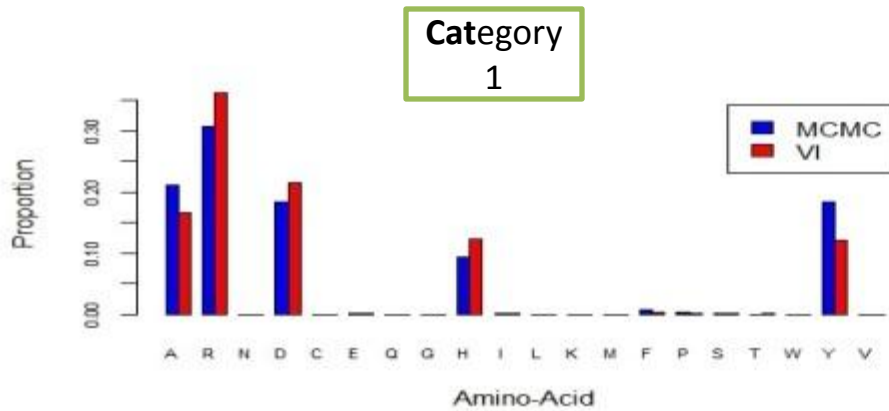
A mitochondrial data set (13 taxa and 6,622 amino acid positions (Rodríguez-Ezpeleta *et al.* 2006))



# Size of profile categories



# Look into the profile of each category



A mitochondrial data set (13 taxa and 6,622 amino acid positions (Rodríguez-Ezpeleta *et al.* 2006))

# Summary

- The pattern of molecular evolution varies among gene sites and genes in a genome.
- By taking into account the complex heterogeneity of evolutionary processes among sites in a genome, Bayesian infinite mixture models of genomic evolution enable robust phylogenetic inference.
- With large modern data sets, however, the computational burden of Markov chain Monte Carlo sampling techniques becomes prohibitive.
- Here, we have developed a variational Bayesian procedure to speed up the widely used PhyloBayes MPI program, which deals with the heterogeneity of amino acid propensity.