



Royal Institute of Technology

# STAT. METH. IN CS – COLLAPSED GIBBS SAMPLER

## Lecture 12

# METROPOLIS HASTINGS (MH)

We want to compute  $p^*(x)$  (typically  $p(x|D)$ )

How?

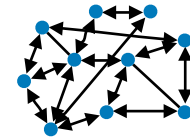
Implicitly construct Markov Chain  $M$  with stationary distribution  $p^*(x)$

Traverse it and sample every  $k$ :th visit

Use good or random starting point

Discard the first  $l$ :th samples

The remaining samples  $x_1, \dots, x_S$  is an approximation of  $p^*(x)$



$$p^*(x) \approx [\sum_i I(x=x_i)] / S$$

# GIBBS SAMPLING

- ★ Pick initial state  $x_1 = (x_{1,1}, \dots, x_{1,K})$
- ★ For  $s=1$  to  $S$ 
  - Sample  $k \sim_u [K]$
  - Sample  $x_{s+1,k} \sim p(x_{s+1,k} | x_{s,-k})$
  - Let  $x_{s+1} = (x_{s,1}, \dots, x_{1,k-1}, x_{s+1,k}, \dots, x_{s,K})$
  - If  $k|s$  record  $x_{s+1}$  (thinning)

Notation

$$\mathcal{D} = (x_1, \dots, x_N), \quad H = (z_1, \dots, z_N), \quad N_k = \sum_n I(z_i = k)$$

$$\boldsymbol{\pi} = (\pi_1, \dots, \pi_k), \quad \boldsymbol{\mu} = (\mu_1, \dots, \mu_k), \quad \boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_k), \quad \text{and} \quad \lambda_k = 1/\sigma_k^2$$

Hyperparameters  $\boldsymbol{\theta}_0 = (\mu_0, \lambda_0, \lambda_0, \beta_0, \alpha)$

Model

$$\boldsymbol{\pi} \sim \text{Dir}(\boldsymbol{\alpha}), \quad \mu_k \sim N(\mu_0, \lambda_0), \quad \lambda_k \sim \text{Ga}(\alpha_0, \beta_0), \quad z_i \sim \text{Cat}(\boldsymbol{\pi}), \quad \text{and}$$

$$p(x_n | Z_n = k) = N(\mu_k, \lambda_k)$$

# GIBBS SAMPLER FOR GMM

# A STATE

$$(H, \pi, \mu, \lambda)$$

Hyperparameters  $\theta_0 = (\mu_0, \lambda_0, \lambda_0, \beta_0, \alpha)$

Model  
 $\pi \sim \text{Dir}(\alpha)$ ,  $\mu_k \sim N(\mu_0, \lambda_0)$ ,  $\lambda_k \sim \text{Ga}(\alpha_0, \beta_0)$ ,  $z_i \sim \text{Cat}(\pi)$ , and  
 $p(x_n | Z_n = k) = N(\mu_k, \lambda_k)$

Likelihood

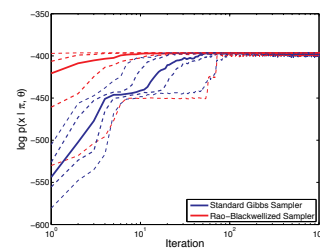
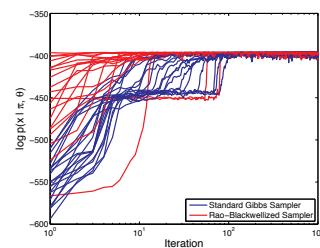
$$\begin{aligned} p(D, H, \pi, \mu, \lambda) &= p(D, H | \pi, \mu, \lambda) p(\pi) p(\mu, \lambda) \\ &= \prod_{n,k} [\pi_k N(x_n | \mu_k, \lambda_k)]^{I(z_n=k)} \text{Dir}(\pi | \alpha) \\ &\quad \prod_k N(\mu_k | \mu_0, \lambda_0) \text{Ga}(\lambda_k | \alpha_0, \beta_0) \end{aligned}$$

# LIKELIHOOD FOR GMM

## COLLAPSING

Integrating out some components of the state is called collapsing

It always improves convergence

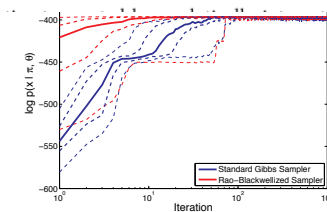
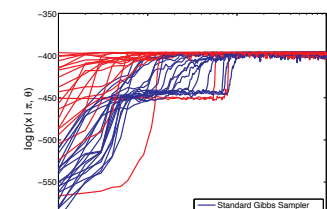


## COLLAPSING

**Theorem 24.2.1** (Rao-Blackwell). Let  $\mathbf{z}$  and  $\theta$  be dependent random variables, and  $f(\mathbf{z}, \theta)$  be some scalar function. Then

$$\text{var}_{\mathbf{z}, \theta} [f(\mathbf{z}, \theta)] \geq \text{var}_{\mathbf{z}} [E_{\theta} [f(\mathbf{z}, \theta) | \mathbf{z}]] \quad (24.20)$$

It always improves convergence



## COLLAPSED GIBBS SAMPLER FOR GMM

Integrate out  $\boldsymbol{\pi}$ ,  $\boldsymbol{\mu}$ , and  $\boldsymbol{\lambda}$  from

$p(D, H, \boldsymbol{\mu}, \boldsymbol{\lambda}, \boldsymbol{\pi}, \boldsymbol{\Theta}_0)$

$\boldsymbol{\Theta}_0$  is all the hyperparameters

Only full conditionals on  $z_n$  remain

Same joint as before except from conjugate prior on  $\mu_k$  and  $\lambda_k$

## NEW FULL CONDITIONAL

=3

$$\begin{aligned} p(z_n | D, H_{-n}, \boldsymbol{\Theta}_0) &= \frac{p(z_n, D | H_{-n}, \boldsymbol{\Theta}_0)}{p(D | H_{-n}, \boldsymbol{\Theta}_0)} \\ &\propto p(z_n, D | H_{-n}, \boldsymbol{\Theta}_0) \\ &\propto p(z_n | H_{-n}, \boldsymbol{\Theta}_0) p(D | z_n, H_{-n}, \boldsymbol{\Theta}_0) \\ &\propto p(z_n | H_{-n}, \boldsymbol{\Theta}_0) p(x_n | D_{-n}, z_n, H_{-n}, \boldsymbol{\Theta}_0) p(D_{-n} | z_n, H_{-n}, \boldsymbol{\Theta}_0) \\ &\propto p(z_n | H_{-n}, \boldsymbol{\Theta}_0) p(x_n | D_{-n}, z_n, H_{-n}, \boldsymbol{\Theta}_0) \end{aligned}$$

“easy” with a conjugate prior on  $\mu_k$  and  $\lambda_k$  but varies as we vary H

## FULL CONDITIONAL

Recall,  $\boldsymbol{\pi} \sim \text{Dir}(\boldsymbol{\alpha})$ ,  $\mu_k \sim N(\mu_0, \lambda_0)$ ,  $\lambda_k \sim \text{Ga}(\alpha_0, \beta_0)$ ,  $z_i \sim \text{Cat}(\boldsymbol{\pi})$ , and  $p(x_n | Z_n = k) = N(\mu_k, \lambda_k)$

i.e., marginal

$$p(z_1, \dots, z_N | \boldsymbol{\alpha}) = \frac{\Gamma(\alpha)}{\Gamma(N + \alpha)} \prod_{k=1}^K \frac{\Gamma(N_k + \alpha/K)}{\Gamma(\alpha/K)}$$

So,

$$\begin{aligned} p(z_i = k | \mathbf{z}_{-i}, \boldsymbol{\alpha}) &= \frac{p(\mathbf{z}_{1:N} | \boldsymbol{\alpha})}{p(\mathbf{z}_{-i} | \boldsymbol{\alpha})} = \frac{\frac{1}{\Gamma(N + \alpha)}}{\frac{1}{\Gamma(N + \alpha - 1)}} \times \frac{\Gamma(N_k + \alpha/K)}{\Gamma(N_{k,-i} + \alpha/K)} \\ &= \frac{\Gamma(N + \alpha - 1)}{\Gamma(N + \alpha)} \frac{\Gamma(N_{k,-i} + 1 + \alpha/K)}{\Gamma(N_{k,-i} + \alpha/K)} = \frac{N_{k,-i} + \alpha/K}{N + \alpha - 1} \end{aligned}$$

where  $N_{k,-i} \triangleq \sum_{n \neq i} \mathbb{I}(z_n = k) = N_k - 1$ ,  $\boldsymbol{\alpha} = \sum_k \alpha_k$

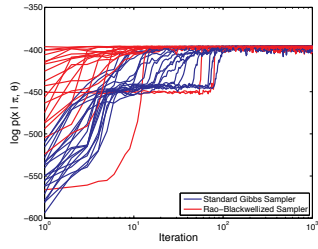
## THE COLLAPSED ALGORITHM

**Algorithm 24.1:** Collapsed Gibbs sampler for a mixture model

```

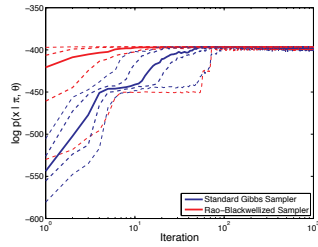
1 for each  $i = 1 : N$  in random order do
2   Remove  $\mathbf{x}_i$ 's sufficient statistics from old cluster  $z_i$ ;
3   for each  $k = 1 : K$  do
4     Compute  $p_k(\mathbf{x}_i) \triangleq p(\mathbf{x}_i | \{\mathbf{x}_j : z_j = k, j \neq i\})$ ;
5   Compute  $p(z_i = k | \mathbf{z}_{-i}, \mathcal{D}) \propto (N_{k,-i} + \alpha/K) p_k(\mathbf{x}_i)$ ;
6   Sample  $z_i \sim p(z_i | \cdot)$ ;
7   Add  $\mathbf{x}_i$ 's sufficient statistics to new cluster  $z_i$ 

```

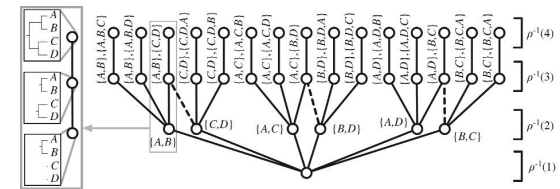
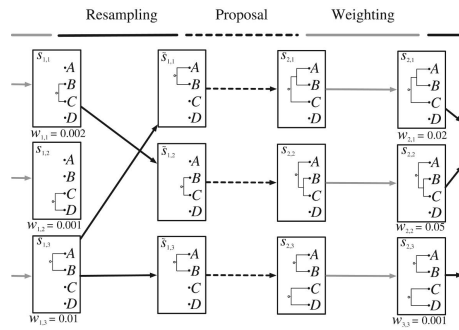


# CONVERGENCE COLLAPSED VS STANDARD

- Blue is standard
- Red is collapsed
- x iterations
- y loglikelihood



# THE END



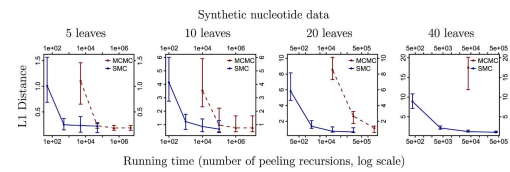


FIGURE 3. Comparison of the convergence time of PoseSMC and MCMC. We generated coalescent trees of different sizes and data sets of 1000 nucleotides. We computed the L1 distance of the minimum Bayes risk reconstruction to the true generating tree as a function of the running time (in units of the number of peeling recursions, on a log scale). The missing MCMC data points are due to MrBayes stalling on these executions.