

Markov Chain Monte Carlo

Shinichi Nakajima

Technische Universität Berlin

14. Jan. 2016

Bayesian Learning (easy cases)

Gaussian

$$\text{Posterior: } p(\mathbf{w}|\mathcal{D}) \propto p(\mathbf{w}, \mathcal{D}) \propto \underbrace{\exp\left(-\frac{\boldsymbol{\mu}^\top \mathbf{C}^{-1} \boldsymbol{\mu}}{2}\right)}_{\text{prior}} \underbrace{\prod_{i=1}^n \exp\left(-\frac{(\boldsymbol{\mu}-\mathbf{x}^{(i)})^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}-\mathbf{x}^{(i)})}{2}\right)}_{\text{likelihood}}$$

Complete the square! $\propto \exp\left(-\frac{(\boldsymbol{\mu}-\mathbf{m}_\mu)^\top \mathbf{S}_\mu^{-1} (\boldsymbol{\mu}-\mathbf{m}_\mu)}{2}\right)$

$$\therefore p(\mathbf{w}|\mathcal{D}) = \mathcal{N}(\boldsymbol{\mu}; \mathbf{m}_\mu, \mathbf{S}_\mu)$$

Multinomial

$$\text{Posterior: } p(\mathbf{w}|\mathcal{D}) \propto p(\mathbf{w}, \mathcal{D}) \propto \underbrace{\prod_{k=1}^K \theta_k^{\phi_k-1}}_{\text{prior}} \underbrace{\prod_{k=1}^K \theta_k^{x_k}}_{\text{likelihood}}$$

Add exponents! $\propto \prod_{k=1}^K \theta_k^{x_k + \phi_k - 1}$

$$\therefore p(\mathbf{w}|\mathcal{D}) = \text{Dir}(\boldsymbol{\theta}; \{x_k + \phi_k\}_{k=1}^K)$$

These two patterns cover most of the cases, including approximate learning!

Bayesian learning is computationally hard, but Bayesian do easy work. All what you need are **square completion**, **addition**, and **Wikipedia (to consult on moments)**!

If Wikipedia does not know, ...

moments (expectation of some function $f(\mathbf{w})$) are approximated by

$$\int f(\mathbf{w}) \cdot p(\mathbf{w}|\mathcal{D})d\mathbf{w} \approx \frac{1}{J} \sum_{j=1}^J f(\mathbf{w}^{(j)}), \quad \text{where } \mathbf{w}^{(j)} \sim p(\mathbf{w}|\mathcal{D}).$$

MCMC Sampling
(Metropolis-Hastings,
Gibbs Sampling)

or

$$\int f(\mathbf{w}) \cdot p(\mathbf{w}|\mathcal{D})d\mathbf{w} \approx \int f(\mathbf{w}) \cdot q(\mathbf{w})d\mathbf{w} \quad \text{where } q(\mathbf{w}) \approx p(\mathbf{w}|\mathcal{D}).$$

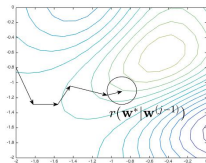
Deterministic Approximation
(Laplace Approximation,
Variational Bayes,
Expectation Propagation)

$q(\mathbf{w})$ must be in a known form.

Sampling from *unnormalized* distribution

We need samples such that

$$\mathbf{w} \sim \underbrace{p(\mathbf{w}|\mathcal{D})}_{\text{unknown}} \propto \underbrace{p(\mathbf{w}, \mathcal{D})}_{\text{known!}}$$



Metropolis-Hastings (MH) (a general method):

In the j -th iteration,

- 1 Draw a sample $\mathbf{w}^* \sim r(\mathbf{w}|\mathbf{w}^{(j-1)})$, e.g., $r(\mathbf{w}|\mathbf{w}^{(j-1)}) \propto \exp\left(-\frac{\|\mathbf{w}-\mathbf{w}^{(j-1)}\|^2}{2\gamma^2}\right)$.
- 2 Set the j -th sample to

$$\mathbf{w}^{(j)} = \begin{cases} \mathbf{w}^* & \text{with probability } T, \\ \mathbf{w}^{(j-1)} & \text{with probability } 1 - T, \end{cases}$$

with acceptance probability (which satisfies **detailed balance property** [appendix](#))

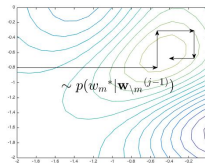
$$T = \min\left(1, \frac{p(\mathbf{w}^*, \mathcal{D}) r(\mathbf{w}^{(j-1)}|\mathbf{w}^*)}{p(\mathbf{w}^{(j-1)}, \mathcal{D}) r(\mathbf{w}^*|\mathbf{w}^{(j-1)})}\right).$$

Sampling from *unnormalized* distribution

We need samples such that

$$\mathbf{w} \sim \underbrace{p(\mathbf{w}|\mathcal{D})}_{\text{unknown}} \propto \underbrace{p(\mathbf{w}, \mathcal{D})}_{\text{known!}}.$$

If $p(w_m | \mathbf{w}_{\setminus m}, \mathcal{D})$ is in a known form (sampler available),

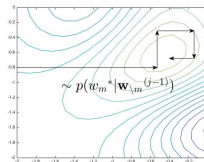


Sampling from *unnormalized* distribution

We need samples such that

$$\mathbf{w} \sim \underbrace{p(\mathbf{w}|\mathcal{D})}_{\text{unknown}} \propto \underbrace{p(\mathbf{w}, \mathcal{D})}_{\text{known!}}$$

If $p(w_m | \mathbf{w}_{\setminus m}, \mathcal{D})$ is in a known form (sampler available),



Gibbs Sampling (GS) (a special case of MH) is efficient:

In the j -th iteration,

- 1 Draw a sample *element* $w_m^* \sim p(w_m | \mathbf{w}_{\setminus m}^{(j-1)}, \mathcal{D})$, and set $\mathbf{w}_{\setminus m}^* = \mathbf{w}_{\setminus m}^{(j-1)}$.
- 2 Set $\mathbf{w}_m^{(j)} = \mathbf{w}^*$ (with probability 1).

The reason why the acceptance probability is one is

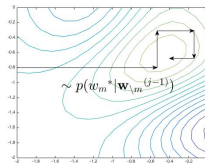
$$T = \min \left(1, \frac{p(\mathbf{w}^*, \mathcal{D}) r(\mathbf{w}^{(j-1)} | \mathbf{w}^*)}{p(\mathbf{w}^{(j-1)}, \mathcal{D}) r(\mathbf{w}^* | \mathbf{w}^{(j-1)})} \right)$$

Sampling from *unnormalized* distribution

We need samples such that

$$\mathbf{w} \sim \underbrace{p(\mathbf{w}|\mathcal{D})}_{\text{unknown}} \propto \underbrace{p(\mathbf{w}, \mathcal{D})}_{\text{known!}}$$

If $p(w_m|\mathbf{w}_{\setminus m}, \mathcal{D})$ is in a known form (sampler available),



Gibbs Sampling (GS) (a special case of MH) is efficient:

In the j -th iteration,

- 1 Draw a sample *element* $w_m^* \sim p(w_m|\mathbf{w}_{\setminus m}^{(j-1)}, \mathcal{D})$, and set $\mathbf{w}_{\setminus m}^* = \mathbf{w}_{\setminus m}^{(j-1)}$.
- 2 Set $\mathbf{w}_m^{(j)} = \mathbf{w}^*$ (with probability 1).

The reason why the acceptance probability is one is

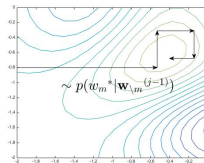
$$T = \min \left(1, \frac{p(\mathbf{w}^*, \mathcal{D}) r(\mathbf{w}^{(j-1)}|\mathbf{w}^*)}{p(\mathbf{w}^{(j-1)}, \mathcal{D}) r(\mathbf{w}^*|\mathbf{w}^{(j-1)})} \right) = \min \left(1, \frac{p(\mathbf{w}^*, \mathcal{D}) p(w_m^{(j-1)}|\mathbf{w}_{\setminus m}^*, \mathcal{D}) p(\mathbf{w}_{\setminus m}^{(j-1)}|\mathbf{w}_{\setminus m}^*, \mathcal{D})}{p(\mathbf{w}^{(j-1)}, \mathcal{D}) p(w_m^*|\mathbf{w}_{\setminus m}^{(j-1)}, \mathcal{D}) p(\mathbf{w}_{\setminus m}^*|\mathbf{w}_{\setminus m}^{(j-1)}, \mathcal{D})} \right).$$

Sampling from *unnormalized* distribution

We need samples such that

$$\mathbf{w} \sim \underbrace{p(\mathbf{w}|\mathcal{D})}_{\text{unknown}} \propto \underbrace{p(\mathbf{w}, \mathcal{D})}_{\text{known!}}$$

If $p(w_m | \mathbf{w}_{\setminus m}, \mathcal{D})$ is in a known form (sampler available),



Gibbs Sampling (GS) (a special case of MH) is efficient:

In the j -th iteration,

- 1 Draw a sample *element* $w_m^* \sim p(w_m | \mathbf{w}_{\setminus m}^{(j-1)}, \mathcal{D})$, and set $\mathbf{w}_{\setminus m}^* = \mathbf{w}_{\setminus m}^{(j-1)}$.
- 2 Set $\mathbf{w}_m^{(j)} = \mathbf{w}^*$ (with probability 1).

The reason why the acceptance probability is one is

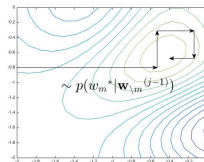
$$T = \min \left(1, \frac{p(\mathbf{w}^*, \mathcal{D}) r(\mathbf{w}^{(j-1)} | \mathbf{w}^*)}{p(\mathbf{w}^{(j-1)}, \mathcal{D}) r(\mathbf{w}^* | \mathbf{w}^{(j-1)})} \right) = \min \left(1, \frac{p(\mathbf{w}^*, \mathcal{D}) p(w_m^{(j-1)} | \mathbf{w}_{\setminus m}^*, \mathcal{D}) p(\mathbf{w}_{\setminus m}^{(j-1)} | \mathbf{w}_{\setminus m}^*, \mathcal{D})}{p(\mathbf{w}^{(j-1)}, \mathcal{D}) p(w_m^* | \mathbf{w}_{\setminus m}^{(j-1)}, \mathcal{D}) p(\mathbf{w}_{\setminus m}^* | \mathbf{w}_{\setminus m}^{(j-1)}, \mathcal{D})} \right).$$

Sampling from *unnormalized* distribution

We need samples such that

$$\mathbf{w} \sim \underbrace{p(\mathbf{w}|\mathcal{D})}_{\text{unknown}} \propto \underbrace{p(\mathbf{w}, \mathcal{D})}_{\text{known!}}$$

If $p(w_m|\mathbf{w}_{\setminus m}, \mathcal{D})$ is in a known form (sampler available),



Gibbs Sampling (GS) (a special case of MH) is efficient:

In the j -th iteration,

- 1 Draw a sample *element* $w_m^* \sim p(w_m|\mathbf{w}_{\setminus m}^{(j-1)}, \mathcal{D})$, and set $\mathbf{w}_{\setminus m}^* = \mathbf{w}_{\setminus m}^{(j-1)}$.
- 2 Set $\mathbf{w}_m^{(j)} = \mathbf{w}^*$ (with probability 1).

The reason why the acceptance probability is one is

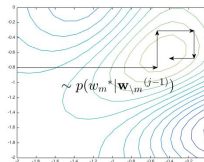
$$T = \min \left(1, \frac{p(\mathbf{w}^*, \mathcal{D}) r(\mathbf{w}^{(j-1)}|\mathbf{w}^*)}{p(\mathbf{w}^{(j-1)}, \mathcal{D}) r(\mathbf{w}^*|\mathbf{w}^{(j-1)})} \right) = \min \left(1, \frac{p(\mathbf{w}^*, \mathcal{D}) p(w_m^{(j-1)}|\mathbf{w}_{\setminus m}^*, \mathcal{D})}{p(\mathbf{w}^{(j-1)}, \mathcal{D}) p(\mathbf{w}_{\setminus m}^*|\mathbf{w}_m^{(j-1)}, \mathcal{D})} \right).$$

Sampling from *unnormalized* distribution

We need samples such that

$$\mathbf{w} \sim \underbrace{p(\mathbf{w}|\mathcal{D})}_{\text{unknown}} \propto \underbrace{p(\mathbf{w}, \mathcal{D})}_{\text{known!}}$$

If $p(w_m|\mathbf{w}_{\setminus m}, \mathcal{D})$ is in a known form (sampler available),



Gibbs Sampling (GS) (a special case of MH) is efficient:

In the j -th iteration,

- 1 Draw a sample *element* $w_m^* \sim p(w_m|\mathbf{w}_{\setminus m}^{(j-1)}, \mathcal{D})$, and set $\mathbf{w}_{\setminus m}^* = \mathbf{w}_{\setminus m}^{(j-1)}$.
- 2 Set $\mathbf{w}_m^{(j)} = \mathbf{w}^*$ (with probability 1).

The reason why the acceptance probability is one is

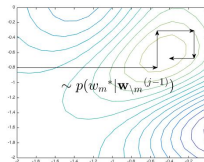
$$T = \min \left(1, \frac{p(\mathbf{w}^*, \mathcal{D}) r(\mathbf{w}^{(j-1)}|\mathbf{w}^*)}{p(\mathbf{w}^{(j-1)}, \mathcal{D}) r(\mathbf{w}^*|\mathbf{w}^{(j-1)})} \right) = \min \left(1, \frac{p(\mathbf{w}^*, \mathcal{D}) p(w_m^{(j-1)}|\mathbf{w}_{\setminus m}^*, \mathcal{D})}{p(\mathbf{w}^{(j-1)}, \mathcal{D}) p(w_m^*|\mathbf{w}_{\setminus m}^{(j-1)}, \mathcal{D})} \right).$$

Sampling from *unnormalized* distribution

We need samples such that

$$\mathbf{w} \sim \underbrace{p(\mathbf{w}|\mathcal{D})}_{\text{unknown}} \propto \underbrace{p(\mathbf{w}, \mathcal{D})}_{\text{known!}}$$

If $p(w_m|\mathbf{w}_{\setminus m}, \mathcal{D})$ is in a known form (sampler available),



Gibbs Sampling (GS) (a special case of MH) is efficient:

In the j -th iteration,

- 1 Draw a sample *element* $w_m^* \sim p(w_m|\mathbf{w}_{\setminus m}^{(j-1)}, \mathcal{D})$, and set $\mathbf{w}_{\setminus m}^* = \mathbf{w}_{\setminus m}^{(j-1)}$.
- 2 Set $\mathbf{w}_m^{(j)} = \mathbf{w}^*$ (with probability 1).

The reason why the acceptance probability is one is

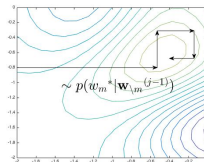
$$T = \min \left(1, \frac{p(\mathbf{w}^*, \mathcal{D}) r(\mathbf{w}^{(j-1)}|\mathbf{w}^*)}{p(\mathbf{w}^{(j-1)}, \mathcal{D}) r(\mathbf{w}^*|\mathbf{w}^{(j-1)})} \right) = \min \left(1, \frac{p(w_m^*|\mathbf{w}_{\setminus m}^*, \mathcal{D}) p(\mathbf{w}_{\setminus m}^*, \mathcal{D}) p(w_m^{(j-1)}|\mathbf{w}_{\setminus m}^*, \mathcal{D})}{p(w_m^{(j-1)}|\mathbf{w}_{\setminus m}^{(j-1)}, \mathcal{D}) p(\mathbf{w}_{\setminus m}^{(j-1)}, \mathcal{D}) p(w_m^*|\mathbf{w}_{\setminus m}^{(j-1)}, \mathcal{D})} \right)$$

Sampling from *unnormalized* distribution

We need samples such that

$$\mathbf{w} \sim \underbrace{p(\mathbf{w}|\mathcal{D})}_{\text{unknown}} \propto \underbrace{p(\mathbf{w}, \mathcal{D})}_{\text{known!}}$$

If $p(w_m | \mathbf{w}_{\setminus m}, \mathcal{D})$ is in a known form (sampler available),



Gibbs Sampling (GS) (a special case of MH) is efficient:

In the j -th iteration,

- 1 Draw a sample *element* $w_m^* \sim p(w_m | \mathbf{w}_{\setminus m}^{(j-1)}, \mathcal{D})$, and set $\mathbf{w}_{\setminus m}^* = \mathbf{w}_{\setminus m}^{(j-1)}$.
- 2 Set $\mathbf{w}_m^{(j)} = \mathbf{w}^*$ (with probability 1).

The reason why the acceptance probability is one is

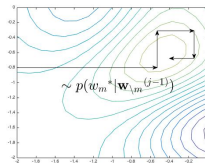
$$T = \min \left(1, \frac{p(\mathbf{w}^*, \mathcal{D}) r(\mathbf{w}^{(j-1)} | \mathbf{w}^*)}{p(\mathbf{w}^{(j-1)}, \mathcal{D}) r(\mathbf{w}^* | \mathbf{w}^{(j-1)})} \right) = \min \left(1, \frac{p(w_m^* | \mathbf{w}_{\setminus m}^*, \mathcal{D}) p(\mathbf{w}_{\setminus m}^*, \mathcal{D}) p(w_m^{(j-1)} | \mathbf{w}_{\setminus m}^*, \mathcal{D})}{p(w_m^{(j-1)} | \mathbf{w}_{\setminus m}^{(j-1)}, \mathcal{D}) p(\mathbf{w}_{\setminus m}^{(j-1)}, \mathcal{D}) p(w_m^* | \mathbf{w}_{\setminus m}^{(j-1)}, \mathcal{D})} \right)$$

Sampling from *unnormalized* distribution

We need samples such that

$$\mathbf{w} \sim \underbrace{p(\mathbf{w}|\mathcal{D})}_{\text{unknown}} \propto \underbrace{p(\mathbf{w}, \mathcal{D})}_{\text{known!}}$$

If $p(w_m|\mathbf{w}_{\setminus m}, \mathcal{D})$ is in a known form (sampler available),



Gibbs Sampling (GS) (a special case of MH) is efficient:

In the j -th iteration,

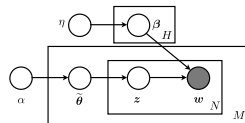
- 1 Draw a sample *element* $w_m^* \sim p(w_m|\mathbf{w}_{\setminus m}^{(j-1)}, \mathcal{D})$, and set $\mathbf{w}_{\setminus m}^* = \mathbf{w}_{\setminus m}^{(j-1)}$.
- 2 Set $\mathbf{w}_m^{(j)} = \mathbf{w}^*$ (with probability 1).

The reason why the acceptance probability is one is

$$T = \min \left(1, \frac{p(\mathbf{w}^*, \mathcal{D}) r(\mathbf{w}^{(j-1)}|\mathbf{w}^*)}{p(\mathbf{w}^{(j-1)}, \mathcal{D}) r(\mathbf{w}^*|\mathbf{w}^{(j-1)})} \right) = \min \left(1, \frac{p(w_m^*|\mathbf{w}_{\setminus m}^*, \mathcal{D}) p(\mathbf{w}_{\setminus m}^*, \mathcal{D}) p(w_m^{(j-1)}|\mathbf{w}_{\setminus m}^*, \mathcal{D})}{p(w_m^{(j-1)}|\mathbf{w}_{\setminus m}^*, \mathcal{D}) p(\mathbf{w}_{\setminus m}^*, \mathcal{D}) p(w_m^*|\mathbf{w}_{\setminus m}^*, \mathcal{D})} \right) = 1.$$

Latent Dirichlet Allocation

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services," Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.



Likelihood and Prior

$$p(\mathbf{w}^{(n,m)} | \boldsymbol{\theta}, \mathbf{B}) = \prod_{l=1}^L ((\mathbf{B}\boldsymbol{\theta}^\top)_{l,m})^{w_l^{(n,m)}},$$

$$p(\boldsymbol{\theta}) \propto \prod_{m=1}^M \prod_{h=1}^H (\boldsymbol{\theta}_{m,h})^{\alpha-1},$$

$$p(\mathbf{B}) \propto \prod_{h=1}^H \prod_{l=1}^L (\mathbf{B}_{l,h})^{\eta-1}.$$

Word distribution is a mixture of multinomial *topic* distribution.

$\boldsymbol{\theta} \in [0, 1]^{M \times H}$: Document parameter

$\mathbf{B} \in [0, 1]^{L \times H}$: Topic parameter

M : # of documents

L : vocabulary size

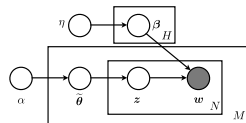
H : # of topics ($\leq \min(M, L)$)

$N^{(m)}$: # of words in m -th document

α, β : Hyperparameters

Latent Dirichlet Allocation

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services," Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.



Likelihood and Prior

$$p(\mathbf{w}^{(n,m)} | \boldsymbol{\theta}, \mathbf{B}) = \prod_{l=1}^L \left(\sum_{h=1}^H B_{l,h} \theta_{m,h} \right)^{w_l^{(n,m)}},$$

$$p(\boldsymbol{\theta}) \propto \prod_{m=1}^M \prod_{h=1}^H (\theta_{m,h})^{\alpha-1},$$

$$p(\mathbf{B}) \propto \prod_{h=1}^H \prod_{l=1}^L (B_{l,h})^{\eta-1}.$$

$\boldsymbol{\theta} \in [0, 1]^{M \times H}$: Document parameter

$\mathbf{B} \in [0, 1]^{L \times H}$: Topic parameter

M : # of documents

L : vocabulary size

H : # of topics ($\leq \min(M, L)$)

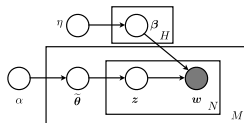
$N^{(m)}$: # of words in m -th document

α, β : Hyperparameters

Sum in the probability (mixture) is intractable.

Latent Dirichlet Allocation

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services," Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.



Complete Likelihood and Prior

$$p(\mathbf{w}^{(n,m)}, \mathbf{z}^{(n,m)} | \boldsymbol{\theta}, \mathbf{B}) = \prod_{l=1}^L \left(\prod_{h=1}^H (B_{l,h} \theta_{m,h}) z_h^{(n,m)} \right)^{w_l^{(n,m)}},$$

$$p(\boldsymbol{\theta}) \propto \prod_{m=1}^M \prod_{h=1}^H (\theta_{m,h})^{\alpha-1},$$

$$p(\mathbf{B}) \propto \prod_{h=1}^H \prod_{l=1}^L (B_{l,h})^{\eta-1}.$$

Latent variable $\mathbf{z}^{(n,m)}$ changes the sum to the product,

$\boldsymbol{\theta} \in [0, 1]^{M \times H}$: Document parameter

$\mathbf{B} \in [0, 1]^{L \times H}$: Topic parameter

$\mathbf{z}^{(n,m)} \in [0, 1]^H$: Topic assignment for each word

M : # of documents

L : vocabulary size

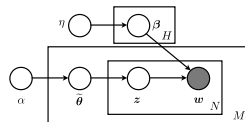
H : # of topics ($\leq \min(M, L)$)

$N^{(m)}$: # of words in m -th document

α, β : Hyperparameters

Latent Dirichlet Allocation

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services," Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.



Complete Likelihood and Prior

$$p(\mathbf{w}^{(n,m)}, \mathbf{z}^{(n,m)} | \boldsymbol{\theta}, \mathbf{B}) = \prod_{h=1}^H (\boldsymbol{\theta}_{m,h})^{z_h^{(n,m)}} \prod_{l=1}^L (\mathbf{B}_{l,h})^{w_l^{(n,m)} z_h^{(n,m)}},$$

$$p(\boldsymbol{\theta}) \propto \prod_{m=1}^M \prod_{h=1}^H (\boldsymbol{\theta}_{m,h})^{\alpha-1},$$

$$p(\mathbf{B}) \propto \prod_{h=1}^H \prod_{l=1}^L (\mathbf{B}_{l,h})^{\eta-1}.$$

Latent variable $\mathbf{z}^{(n,m)}$ changes the sum to the product, and makes likelihood separable.

$\boldsymbol{\theta} \in [0, 1]^{M \times H}$: Document parameter

$\mathbf{B} \in [0, 1]^{L \times H}$: Topic parameter

$\mathbf{z}^{(n,m)} \in [0, 1]^H$: Topic assignment for each word

M : # of documents

L : vocabulary size

H : # of topics ($\leq \min(M, L)$)

$N^{(m)}$: # of words in m -th document

α, β : Hyperparameters

Joint is NOT in known form, but ...

Posterior: $p(\{\mathbf{z}^{(n,m)}\}, \boldsymbol{\theta}, \mathbf{B} | \{\mathbf{w}^{(n,m)}\}) \propto p(\{\mathbf{z}^{(n,m)}\}, \boldsymbol{\theta}, \mathbf{B}, \{\mathbf{w}^{(n,m)}\})$, where

$$\begin{aligned} p(\{\mathbf{z}^{(n,m)}\}, \boldsymbol{\theta}, \mathbf{B}, \{\mathbf{w}^{(n,m)}\}) &= p(\{\mathbf{w}^{(n,m)}\}, \{\mathbf{z}^{(n,m)}\} | \boldsymbol{\theta}, \mathbf{B}) p(\boldsymbol{\theta}) p(\mathbf{B}) \\ &\propto \prod_{m=1}^M \prod_{n=1}^{N^{(m)}} \prod_{h=1}^H (\boldsymbol{\theta}_{m,h})^{z_h^{(n,m)} + (\alpha-1)/N^{(m)}} \prod_{l=1}^L (\mathbf{B}_{l,h})^{w_l^{(n,m)} z_h^{(n,m)} + (\eta-1)/(MN^{(m)})}. \end{aligned}$$

Joint distribution (on $\{\mathbf{z}^{(n,m)}\}, \boldsymbol{\theta}, \mathbf{B}$) is not in a known form, but *conditionals*

$$p(\{\mathbf{z}^{(n,m)}\} | \boldsymbol{\theta}, \mathbf{B}, \{\mathbf{w}^{(n,m)}\}) \propto \prod_{m=1}^M \prod_{n=1}^{N^{(m)}} \prod_{h=1}^H \left(\boldsymbol{\theta}_{m,h} \prod_{l=1}^L (\mathbf{B}_{l,h})^{w_l^{(n,m)}} \right)^{z_h^{(n,m)}} \quad \text{Multinomial}$$

$$p(\boldsymbol{\theta} | \{\mathbf{z}^{(n,m)}\}, \{\mathbf{w}^{(n,m)}\}) \propto \prod_{m=1}^M \prod_{n=1}^{N^{(m)}} \prod_{h=1}^H (\boldsymbol{\theta}_{m,h})^{z_h^{(n,m)} + (\alpha-1)/N^{(m)}} \quad \text{Dirichlet}$$

$$p(\mathbf{B} | \{\mathbf{z}^{(n,m)}\}, \{\mathbf{w}^{(n,m)}\}) \propto \prod_{m=1}^M \prod_{n=1}^{N^{(m)}} \prod_{h=1}^H \prod_{l=1}^L (\mathbf{B}_{l,h})^{w_l^{(n,m)} z_h^{(n,m)} + (\eta-1)/(MN^{(m)})} \quad \text{Dirichlet}$$

are in known forms. → **Gibbs sampling!**

Gibbs Sampling

(Naive) Gibbs Sampling:

In j -th step,

- 1 For each (n, m) **independently**, draw $\mathbf{z}^{(j,n,m)} \sim p(\mathbf{z}^{(n,m)} | \boldsymbol{\theta}^{(j-1)}, \mathbf{B}^{(j-1)}, \{\mathbf{w}^{(n,m)}\})$,
- 2 For each m **independently**, draw $\tilde{\boldsymbol{\theta}}_m^{(j)} \sim p(\tilde{\boldsymbol{\theta}}_m | \{\mathbf{z}^{(j,n,m)}\}, \{\mathbf{w}^{(n,m)}\})$,
- 3 For each h **independently**, draw $\beta_h^{(j)} \sim p(\beta_h | \{\mathbf{z}^{(j,n,m)}\}, \{\mathbf{w}^{(n,m)}\})$.

$$p(\{\mathbf{z}^{(n,m)}\} | \boldsymbol{\theta}, \mathbf{B}, \{\mathbf{w}^{(n,m)}\}) \propto \prod_{m=1}^M \prod_{n=1}^{N^{(m)}} \prod_{h=1}^H \left(\boldsymbol{\theta}_{m,h} \prod_{l=1}^L (\mathbf{B}_{l,h})^{w_l^{(n,m)}} \right)^{z_h^{(n,m)}} \quad \text{Multinomial}$$

$$p(\boldsymbol{\theta} | \{\mathbf{z}^{(n,m)}\}, \{\mathbf{w}^{(n,m)}\}) \propto \prod_{m=1}^M \prod_{n=1}^{N^{(m)}} \prod_{h=1}^H (\boldsymbol{\theta}_{m,h})^{z_h^{(n,m)} + (\alpha - 1)/N^{(m)}} \quad \text{Dirichlet}$$

$$p(\mathbf{B} | \{\mathbf{z}^{(n,m)}\}, \{\mathbf{w}^{(n,m)}\}) \propto \prod_{m=1}^M \prod_{n=1}^{N^{(m)}} \prod_{h=1}^H \prod_{l=1}^L (\mathbf{B}_{l,h})^{w_l^{(n,m)} z_h^{(n,m)} + (\eta - 1)/(MN^{(m)})} \quad \text{Dirichlet}$$

Conditionally independent! (\rightarrow easily parallelized) **But we can do better.**

Gibbs Sampling

Since joint is in (independent) Dirichlet forms of $\boldsymbol{\theta}$ and \mathbf{B} , given $\{\mathbf{z}^{(n,m)}\}$, we can marginalize:

$$p(\{\mathbf{z}^{(n,m)}\}, \{\mathbf{w}^{(n,m)}\}) = \int p(\{\mathbf{z}^{(n,m)}\}, \boldsymbol{\theta}, \mathbf{B}, \{\mathbf{w}^{(n,m)}\}) d\boldsymbol{\theta} d\mathbf{B}$$

$$\propto \prod_{h=1}^H \left\{ \left(\prod_{m=1}^M \Gamma\left(\alpha + \sum_{n=1}^{N(m)} z_h^{(n,m)}\right) \right) \left(\frac{\prod_{l=1}^L \Gamma\left(\eta + \sum_{m=1}^M \sum_{n=1}^{N(m)} w_l^{(n,m)} z_h^{(n,m)}\right)}{\Gamma\left(L\eta + \sum_{m=1}^M \sum_{n=1}^{N(m)} z_h^{(n,m)}\right)} \right) \right\}.$$

Now $\{\mathbf{z}^{(n,m)}\}$ are mutually dependent, and

$$p(\mathbf{z}^{(n,m)} | \{\mathbf{z}^{(n',m')}\}_{(n',m') \neq (n,m)}, \{\mathbf{w}^{(n,m)}\}) = \frac{p(\{\mathbf{z}^{(n,m)}\}, \{\mathbf{w}^{(n,m)}\})}{p(\{\mathbf{z}^{(n',m')}\}_{n' \neq n, m' \neq m}, \{\mathbf{w}^{(n,m)}\})}$$

$$\propto \prod_{h=1}^H \left\{ \frac{\left(\alpha + \sum_{n' \neq n} z_h^{(n',m')}\right) \left(\eta + \sum_{(n',m') \neq (n,m)} w_{l(w^{(n,m)})}^{(n',m')} z_h^{(n',m')}\right) z_h^{(n,m)}}{\left(L\eta + \sum_{(n',m') \neq (n,m)} z_h^{(n',m')}\right) z_h^{(n,m)}} \right\}.$$

For derivation, see [Wikipedia](#) for Dirichlet moments, and use $\Gamma(\tau + 1) = \tau\Gamma(\tau)$.

Gibbs Sampling

Collapsed Gibbs Sampling:

In j -th step,

- 1 For a pair (n, m) chosen, draw a sample $\mathbf{z}^{(j,n,m)} \sim p(\mathbf{z}^{(n,m)} | \{\mathbf{z}^{(j-1,n',m')}\}_{(n',m') \neq (n,m)}, \{\mathbf{w}^{(n,m)}\})$, and set $\mathbf{z}^{(j,n',m')} = \mathbf{z}^{(j-1,n',m')}$ for $(n', m') \neq (n, m)$.

Usually, (n, m) is chosen sequentially, and pick a sample every after all $\{\mathbf{z}^{(n,m)}\}$ are updated.

$$\begin{aligned}
 p(\mathbf{z}^{(n,m)} | \{\mathbf{z}^{(n',m')}\}_{(n',m') \neq (n,m)}, \{\mathbf{w}^{(n,m)}\}) &= \frac{p(\{\mathbf{z}^{(n,m)}\}, \{\mathbf{w}^{(n,m)}\})}{p(\{\mathbf{z}^{(n',m')}\}_{n' \neq n, m' \neq m}, \{\mathbf{w}^{(n,m)}\})} \\
 &\propto \prod_{h=1}^H \left\{ \frac{\left(\alpha + \sum_{n' \neq n} z_h^{(n',m')} \right) \left(\eta + \sum_{(n',m') \neq (n,m)} w_{i(w^{(n,m)})}^{(n',m')} z_h^{(n',m')} \right) z_h^{(n,m)}}{\left(L\eta + \sum_{(n',m') \neq (n,m)} z_h^{(n',m')} \right)} \right\}
 \end{aligned}$$

Gibbs Sampling

Collapsed Gibbs Sampling:

In j -th step,

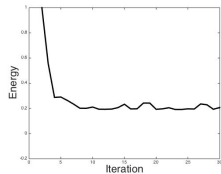
- 1 For a pair (n, m) chosen, draw a sample $\mathbf{z}^{(j,n,m)} \sim p(\mathbf{z}^{(n,m)} | \{\mathbf{z}^{(j-1,n',m')}\}_{(n',m') \neq (n,m)}, \{\mathbf{w}^{(n,m)}\})$, and set $\mathbf{z}^{(j,n',m')} = \mathbf{z}^{(j-1,n',m')}$ for $(n', m') \neq (n, m)$.

Usually, (n, m) is chosen sequentially, and pick a sample every after all $\{\mathbf{z}^{(n,m)}\}$ are updated.

* The energy in the right figure is a collapsed likelihood

$$p(\{\mathbf{z}^{(n,m)}\}, \{\mathbf{w}^{(n,m)}\}) = \int p(\{\mathbf{z}^{(n,m)}\}, \boldsymbol{\theta}, \mathbf{B}, \{\mathbf{w}^{(n,m)}\}) d\boldsymbol{\theta} d\mathbf{B},$$

which is proportional to the posterior $p(\{\mathbf{z}^{(n,m)}\} | \{\mathbf{w}^{(n,m)}\})$.



Gibbs Sampling

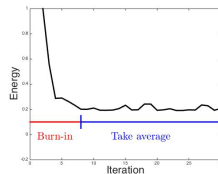
Collapsed Gibbs Sampling:

In j -th step,

- 1 For a pair (n, m) chosen, draw a sample $\mathbf{z}^{(j,n,m)} \sim p(\mathbf{z}^{(n,m)} | \{\mathbf{z}^{(j-1,n',m')}\}_{(n',m') \neq (n,m)}, \{\mathbf{w}^{(n,m)}\})$, and set $\mathbf{z}^{(j,n',m')} = \mathbf{z}^{(j-1,n',m')}$ for $(n', m') \neq (n, m)$.

Usually, (n, m) is chosen sequentially, and pick a sample every after all $\{\mathbf{z}^{(n,m)}\}$ are updated.

With the estimator $\bar{\mathbf{z}}^{(n,m)} = \frac{1}{J - J_{\text{burn-in}}} \sum_{j=J_{\text{burn-in}}+1}^J \mathbf{z}^{(j,n,m)}$,



Gibbs Sampling

Collapsed Gibbs Sampling:

In j -th step,

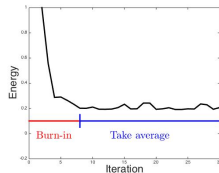
- 1 For a pair (n, m) chosen, draw a sample $\mathbf{z}^{(j,n,m)} \sim p(\mathbf{z}^{(n,m)} | \{\mathbf{z}^{(j-1,n',m')}\}_{(n',m') \neq (n,m)}, \{\mathbf{w}^{(n,m)}\})$, and set $\mathbf{z}^{(j,n',m')} = \mathbf{z}^{(j-1,n',m')}$ for $(n', m') \neq (n, m)$.

Usually, (n, m) is chosen sequentially, and pick a sample every after all $\{\mathbf{z}^{(n,m)}\}$ are updated.

With the estimator $\bar{\mathbf{z}}^{(n,m)} = \frac{1}{J - J_{\text{burn-in}}} \sum_{j=J_{\text{burn-in}}+1}^J \mathbf{z}^{(j,n,m)}$,

$$p(\bar{\boldsymbol{\theta}}_m | \{\bar{\mathbf{z}}^{(n,m)}\}, \{\mathbf{w}^{(n,m)}\}) \propto \prod_{h=1}^H (\boldsymbol{\theta}_{m,h})^{\alpha-1 + \sum_{n=1}^{N(m)} \bar{z}_h^{(n,m)}}$$

$$p(\boldsymbol{\beta}_h | \{\bar{\mathbf{z}}^{(n,m)}\}, \{\mathbf{w}^{(n,m)}\}) \propto \prod_{l=1}^L (B_{l,h})^{\eta-1 + \sum_{m=1}^M \sum_{n=1}^{N(m)} w_l^{(n,m)} \bar{z}_h^{(n,m)}}$$



Gibbs Sampling

Collapsed Gibbs Sampling:

In j -th step,

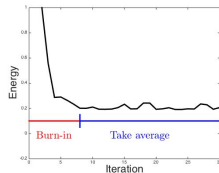
- 1 For a pair (n, m) chosen, draw a sample $\mathbf{z}^{(j,n,m)} \sim p\left(\mathbf{z}^{(n,m)} \mid \{\mathbf{z}^{(j-1,n',m')}\}_{(n',m') \neq (n,m)}, \{\mathbf{w}^{(n,m)}\}\right)$, and set $\mathbf{z}^{(j,n',m')} = \mathbf{z}^{(j-1,n',m')}$ for $(n', m') \neq (n, m)$.

Usually, (n, m) is chosen sequentially, and pick a sample every after all $\{\mathbf{z}^{(n,m)}\}$ are updated.

With the estimator $\bar{\mathbf{z}}^{(n,m)} = \frac{1}{J - J_{\text{burn-in}}} \sum_{j=J_{\text{burn-in}}+1}^J \mathbf{z}^{(j,n,m)}$,

$$p(\bar{\boldsymbol{\theta}}_m \mid \{\bar{\mathbf{z}}^{(n,m)}\}, \{\mathbf{w}^{(n,m)}\}) = \text{Dir}\left(\{\alpha + \sum_{n=1}^{N^{(m)}} \bar{\mathbf{z}}_h^{(n,m)}\}_{h=1}^H\right)$$

$$p(\bar{\boldsymbol{\beta}}_h \mid \{\bar{\mathbf{z}}^{(n,m)}\}, \{\mathbf{w}^{(n,m)}\}) = \text{Dir}\left(\{\eta + \sum_{m=1}^M \sum_{n=1}^{N^{(m)}} w_l^{(n,m)} \bar{\mathbf{z}}_h^{(n,m)}\}_{l=1}^L\right)$$



More efficient! But lost *independence*... ($\{\mathbf{z}^{(n,m)}\}$ are mutually dependent).

Gibbs Sampling

Collapsed Gibbs Sampling:

In j -th step,

- 1 For a pair (n, m) chosen, estimate

$$\widehat{\Theta}_{\setminus(n,m)} = \langle \Theta \rangle_{p(\Theta | \{z^{(j-1,n',m')}\}_{\setminus(n,m)}, \{w^{(n,m)}\})}, \quad \widehat{B}_{\setminus(n,m)} = \langle B \rangle_{p(B | \{z^{(j-1,n',m')}\}_{\setminus(n,m)}, \{w^{(n,m)}\})}.$$

- 2 Draw a sample $z^{(j,n,m)} \sim p(z^{(n,m)} | \widehat{\Theta}_{\setminus(n,m)}, \widehat{B}_{\setminus(n,m)}, \{w^{(n,m)}\})$,

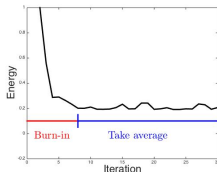
and set $z^{(j,n',m')} = z^{(j-1,n',m')}$ for $(n', m') \neq (n, m)$.

Usually, (n, m) is chosen sequentially, and pick a sample every after all $\{z^{(n,m)}\}$ are updated.

With the estimator $\bar{z}^{(n,m)} = \frac{1}{J - J_{\text{burn-in}}} \sum_{j=J_{\text{burn-in}}+1}^J z^{(j,n,m)}$,

$$p(\bar{\Theta}_m | \{\bar{z}^{(n,m)}\}, \{w^{(n,m)}\}) = \text{Dir}\left(\{\alpha + \sum_{n=1}^{N^{(m)}} \bar{z}_h^{(n,m)}\}_{h=1}^H\right)$$

$$p(\bar{\beta}_h | \{\bar{z}^{(n,m)}\}, \{w^{(n,m)}\}) = \text{Dir}\left(\{\eta + \sum_{m=1}^M \sum_{n=1}^{N^{(m)}} w_l^{(n,m)} \bar{z}_h^{(n,m)}\}_{l=1}^L\right)$$



Amounts to global parameter update after each element sampling.

Gibbs Sampling

Approximate Collapsed Gibbs Sampling

In j -th step,

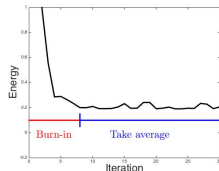
- 1 For chosen pairs (n, m) , **independently** draw $\mathbf{z}^{(j,n,m)} \sim p(\mathbf{z}^{(n,m)} | \widehat{\boldsymbol{\theta}}^{(j-1)}, \widehat{\mathbf{B}}^{(j-1)}, \{\mathbf{w}^{(n,m)}\})$,
- 2 Estimate $\widehat{\boldsymbol{\theta}}^{(j)} = \langle \boldsymbol{\theta} \rangle_{p(\boldsymbol{\theta} | \{\mathbf{z}^{(j,n,m)}\}, \{\mathbf{w}^{(n,m)}\})}$,
- 3 Estimate $\widehat{\mathbf{B}}^{(j)} = \langle \mathbf{B} \rangle_{p(\mathbf{B} | \{\mathbf{z}^{(j,n,m)}\}, \{\mathbf{w}^{(n,m)}\})}$.

With the estimator $\widehat{\mathbf{z}}^{(n,m)} = \frac{1}{J - J_{\text{burn-in}}} \sum_{j=J_{\text{burn-in}}+1}^J \mathbf{z}^{(j,n,m)}$,

$$p(\widehat{\boldsymbol{\theta}}_m | \{\widehat{\mathbf{z}}^{(n,m)}\}, \{\mathbf{w}^{(n,m)}\}) = \text{Dir}\left(\{\alpha + \sum_{n=1}^{N^{(m)}} \widehat{\mathbf{z}}_h^{(n,m)}\}_{h=1}^H\right)$$

$$p(\widehat{\boldsymbol{\beta}}_h | \{\widehat{\mathbf{z}}^{(n,m)}\}, \{\mathbf{w}^{(n,m)}\}) = \text{Dir}\left(\{\eta + \sum_{m=1}^M \sum_{n=1}^{N^{(m)}} w_l^{(n,m)} \widehat{\mathbf{z}}_h^{(n,m)}\}_{l=1}^L\right)$$

This should work.



Variational Bayesian Learning

Variational Bayesian Approximation:

In j -th step,

- 1 For chosen pairs (n, m) , estimate $\hat{\mathbf{z}}^{(j,n,m)} = \langle \mathbf{z}^{(n,m)} \rangle_{p(\mathbf{z}^{(n,m)} | \hat{\boldsymbol{\theta}}^{(j-1)}, \hat{\mathbf{B}}^{(j-1)}, \{\mathbf{w}^{(n,m)}\})}$,
- 2 Estimate $\hat{\boldsymbol{\theta}}^{(j)} = \langle \boldsymbol{\theta} \rangle_{p(\boldsymbol{\theta} | \{\hat{\mathbf{z}}^{(j,n,m)}\}, \{\mathbf{w}^{(n,m)}\})}$,
- 3 Estimate $\hat{\mathbf{B}}^{(j)} = \langle \mathbf{B} \rangle_{p(\mathbf{B} | \{\hat{\mathbf{z}}^{(j,n,m)}\}, \{\mathbf{w}^{(n,m)}\})}$.

Variational Bayes is similar.

Variational Bayesian Learning

Variational Bayesian Approximation:

In j -th step,

- 1 For chosen pairs (n, m) , estimate $\hat{\mathbf{z}}^{(j,n,m)} = \langle \mathbf{z}^{(n,m)} \rangle_{p(\mathbf{z}^{(n,m)} | \hat{\boldsymbol{\theta}}^{(j-1)}, \hat{\mathbf{B}}^{(j-1)}, \{\mathbf{w}^{(n,m)}\})}$,
- 2 Estimate $\hat{\boldsymbol{\theta}}^{(j)} = \langle \boldsymbol{\theta} \rangle_{p(\boldsymbol{\theta} | \{\hat{\mathbf{z}}^{(j,n,m)}\}, \{\mathbf{w}^{(n,m)}\})}$,
- 3 Estimate $\hat{\mathbf{B}}^{(j)} = \langle \mathbf{B} \rangle_{p(\mathbf{B} | \{\hat{\mathbf{z}}^{(j,n,m)}\}, \{\mathbf{w}^{(n,m)}\})}$.

For each k (a part of unknowns) in turn,

Gibbs sampling draws a sample from conditional $w_k \sim p(w_k | \mathbf{w}_{\setminus k}, \mathcal{D})$.

- Slow (each iteration gives **one sample** from the distribution).
- Accurate (correlation between $\{w_k\}$ is **taken into account**).

Variational Bayes estimates the mean of conditional $\hat{w}_k = \langle w_k \rangle_{p(w_k | \hat{\mathbf{w}}_{\setminus k}, \mathcal{D})}$.

- Fast (each iteration estimates **whole distribution**).
- Inaccurate (correlation between $\{w_k\}$ is **neglected**).

Variational Bayes is similar.

Properties of MCMC methods

- Converges to the Bayesian posterior.
- (Generally) slower than deterministic methods.
- To get independent samples, we have to subsample from the MCMC sequence.
- Efficient methods are being developed (Hamiltonian Monte Carlo, distributed computation, stochastic gradient).

Appendix: Detailed Balance Property

Transition probability is (all probability below is conditional on \mathcal{D})

$$\begin{aligned} p(\mathbf{w}^{(j)}|\mathbf{w}^{(j-1)}) &= T(\mathbf{w}^{(j)}|\mathbf{w}^{(j-1)}) \cdot r(\mathbf{w}^{(j)}|\mathbf{w}^{(j-1)}) \\ &= \min\left(r(\mathbf{w}^{(j)}|\mathbf{w}^{(j-1)}), \frac{p(\mathbf{w}^{(j)})r(\mathbf{w}^{(j-1)}|\mathbf{w}^{(j)})}{p(\mathbf{w}^{(j-1)})}\right). \end{aligned}$$

Therefore,

$$\begin{aligned} p(\mathbf{w}^{(j)}|\mathbf{w}^{(j-1)})p(\mathbf{w}^{(j-1)}) &= \min(p(\mathbf{w}^{(j-1)})r(\mathbf{w}^{(j)}|\mathbf{w}^{(j-1)}), p(\mathbf{w}^{(j)})r(\mathbf{w}^{(j-1)}|\mathbf{w}^{(j)})) \\ &= \min(p(\mathbf{w}^{(j)})r(\mathbf{w}^{(j-1)}|\mathbf{w}^{(j)}), p(\mathbf{w}^{(j-1)})r(\mathbf{w}^{(j)}|\mathbf{w}^{(j-1)})) \\ &= T(\mathbf{w}^{(j-1)}|\mathbf{w}^{(j)})r(\mathbf{w}^{(j-1)}|\mathbf{w}^{(j)})p(\mathbf{w}^{(j)}) \\ &= p(\mathbf{w}^{(j-1)}|\mathbf{w}^{(j)})p(\mathbf{w}^{(j)}) \quad \text{Detailed Balance Property} \end{aligned}$$

$$\therefore \int p(\mathbf{w}^{(j)}|\mathbf{w}^{(j-1)})p(\mathbf{w}^{(j-1)})d\mathbf{w}^{(j-1)} = \int p(\mathbf{w}^{(j-1)}|\mathbf{w}^{(j)})p(\mathbf{w}^{(j)})d\mathbf{w}^{(j-1)} = p(\mathbf{w}^{(j)})$$

Detailed balance $\Rightarrow p(\mathbf{w}^{(j)})$ is stationary of Markov process. [back](#)