

STOCHASTIC VARIATIONAL INFERENCE

1. Review variational inference (mean-field)
2. natural gradient
3. stochastic variational inference

1. VI. Evidence lower bound.

$$\log p(x) \geq \mathbb{E}_q[\log p(x, \theta)] - \mathbb{E}_q[\log q(\theta)]$$

mean-field assumption. $\log p(x) \geq \mathbb{E}_q[\log p(x, \theta)] - \sum_i \mathbb{E}_{q_i}[\log q_i(\theta_i)]$

co-ordinate ascent.

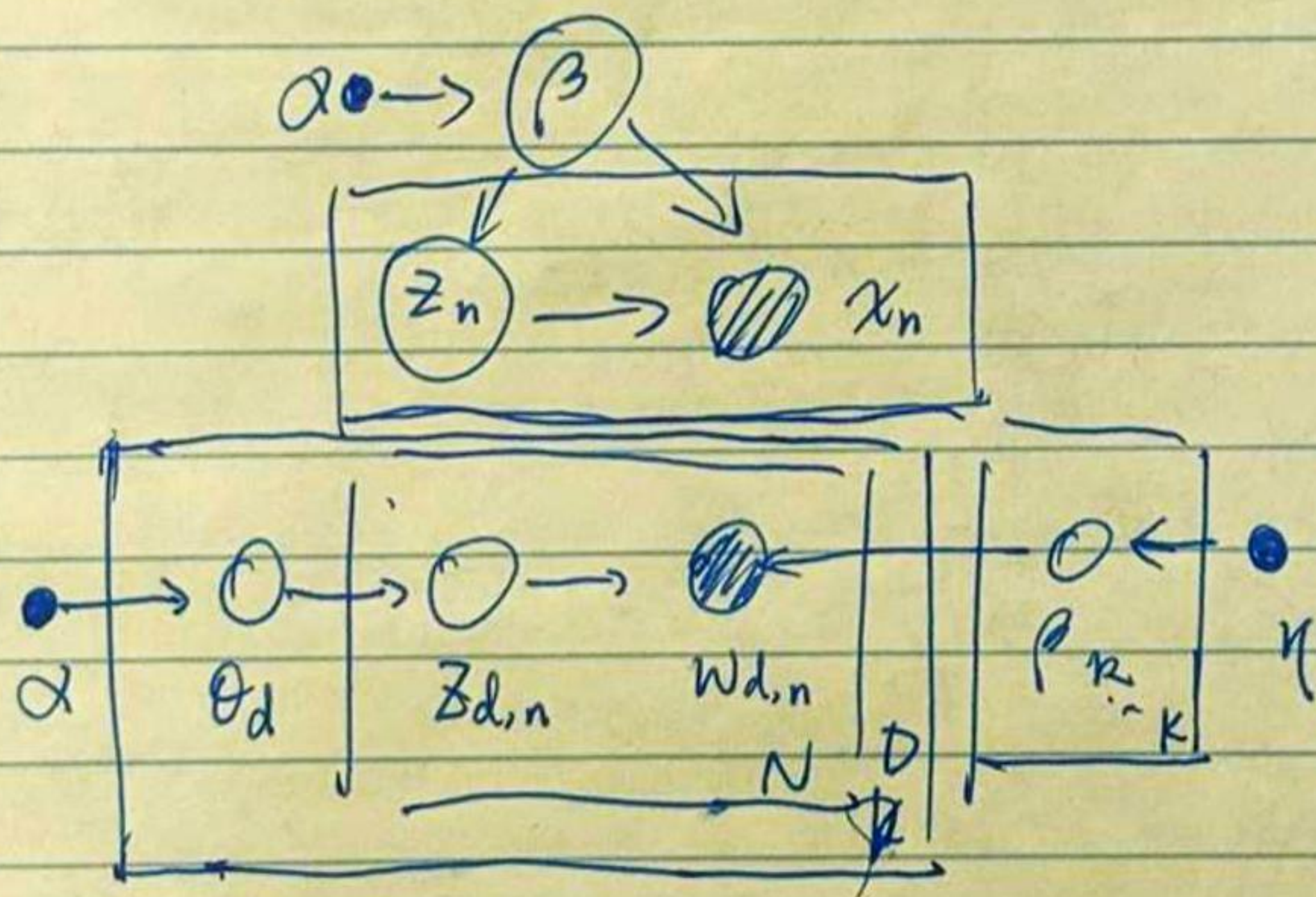
fixed point update $q_i^*(\theta_i) \propto \exp\{\mathbb{E}_{q_{-i}}[\log p(\theta_i | \theta_{-i}, x)]\} H(q_i)$

θ - global hidden variables β . complete statistics: conditions

local hidden variables z

fixed param α .

$$p(x, z, \beta | \alpha) = p(\beta | \alpha) \prod_{n=1}^N p(x_n, z_n | \beta)$$



$$p(\beta | x, z, \alpha) = h(\beta) \exp\{\eta_g(x, z, \alpha)^T t(\beta) - a_g(\eta_g(x, z, \alpha))\}$$

$$p(z_{nj} | x_n, z_{n,-j}, \beta) = h(z_{nj}) \exp\{\eta_l(x_n, z_{n,-j}, \beta)^T t(z_{nj}) - a_l(\eta_l)\}$$

$$p(x_n, z_n | \beta) = \langle \beta^T, a_l(\beta) \rangle \quad p(\beta) = \langle \alpha^T, a_g(\alpha) \rangle$$

$$\eta_g(x, z, \alpha) = (\alpha_1 + \sum_{n=1}^N t(z_n, x_n), \alpha_2 + N)$$

$$\mathcal{L}(q) = \mathbb{E}_q[\log P(x, z, \beta)] - \mathbb{E}_q[\log q(z, \beta)]$$

$$q(z, \beta) = q(\beta | \lambda) \prod_{n=1}^N \prod_{j=1}^J q(z_{nj} | \phi_{nj})$$

$$q(\beta | \lambda) = h(\beta) \exp \{ \lambda^T t(\beta) - \alpha g(\lambda) \}$$

$$q(z_{nj} | \phi_{nj}) = h(z_{nj}) \exp \{ \phi_{nj}^T t(z_{nj}) - \alpha l(\phi_{nj}) \}$$

$$\mathcal{L}(\lambda) = \mathbb{E}_q[\log(P(x, z) \cdot P(\beta | x, z))] - \left(\mathbb{E}_\lambda[\log q(\beta)] + \sum_{n=1}^N \sum_{j=1}^J \mathbb{E}_{\phi_{nj}}[\log q(z_{nj})] \right)$$

const

$$= \mathbb{E}_q[\log P(x, z)] + \mathbb{E}_q[\log P(\beta | x, z)] - \mathbb{E}_\lambda[\log q(\beta)] + \text{const}$$

const

$$= \mathbb{E}_q[\log P(\beta | x, z)] - \mathbb{E}_\lambda[\log q(\beta)]$$

$$\nabla \mathcal{L} = 0 \rightarrow \lambda^* = \mathbb{E}_q[\eta_g(x, z, \lambda)]$$

Similar

$$\phi_{nj} = \mathbb{E}_q[\eta_l(x_n, z_{n,j}, \beta)]$$

Gradient find a maximum $f(\lambda)$

$$\lambda^{(t+1)} = \lambda^{(t)} + \rho \nabla_\lambda f(\lambda^{(t)})$$

learning rate $\mathcal{N}(0, 10000)$

$\mathcal{N}(10, 0.1), \mathcal{N}(10, 10000)$

$\mathcal{N}(1, 0.1)$

symmetrized KL divergence

$$D_{KL}^{sym}(\lambda, \lambda') = \mathbb{E}_\lambda \left[\log \frac{q(\beta | \lambda)}{q(\beta | \lambda')} \right] + \mathbb{E}_{\lambda'} \left[\log \frac{q(\beta | \lambda')}{q(\beta | \lambda)} \right]$$

~~KL~~

$$dx^T G(x) dx = D_{KL}^{sym}(\lambda, \lambda + d\lambda)$$

$$\hat{\nabla}_\lambda f(\lambda) \stackrel{\Delta}{=} G(\lambda)^{-1} \nabla_\lambda f(\lambda)$$

Riemannian metric

$G(\lambda) \triangleq$ fisher information matrix of $q(w)$
 = log normalizer = $\nabla_{\lambda}^2 \log \zeta(\lambda)$

$$\hat{\nabla}_{\lambda} \zeta = \mathbb{E}_q[\eta_g(x, z, \alpha)] - \lambda$$

$\hat{\nabla}_{\lambda} \zeta_i = \mathbb{E}_q[\eta_g(x_i^{(N)}, z_i^{(N)}, \alpha)] - \lambda$
 intermediate global params $\hat{\lambda} = \mathbb{E}_q[\eta_g(x_i^{(N)}, z_i^{(N)}, \alpha)]$

$$\lambda^{(t)} = \lambda^{(t-1)} + \rho_t \hat{\nabla}_{\lambda} \zeta_i$$

$$= \lambda^{(t-1)} + \rho_t (\mathbb{E}_q[\eta_g(x_i^{(N)}, z_i^{(N)}, \alpha)] - \lambda^{(t-1)})$$

$$= (1 - \rho_t) \lambda^{(t-1)} + \rho_t \mathbb{E}_q[\eta_g(x_i^{(N)}, z_i^{(N)}, \alpha)]$$

$$\rho_t = (t + \tau)^{-k} \leftarrow \text{forgetting gate, } \tau \in (0.5, 1]$$

mini batch. $\lambda^{(t)} = (1 - \rho_t) \lambda^{(t-1)} + \frac{\rho_t}{S} \sum_S \hat{\lambda}_s$

LDA example.

$$p(z_{dn} = k | \theta_d, \beta_{1:k}, w_{dn}) \propto \alpha \exp \{ \log \theta_{dk} + \log \beta_{k, w_{dn}} \}$$

$$q(z_{dn}) \sim \text{Multinomial}(\phi_{dn})$$

$$p(\theta_d | z_d) = \text{Dirichlet}(\alpha + \sum_{n=1}^N z_{dn})$$

$$q(\theta_d) \sim \text{Dirichlet}(\gamma_d)$$

local update. $\phi_{dn}^k \propto \gamma_d^k$

$$\phi_{dn}^k \propto \exp \{ \mathbb{E}[\log \theta_{dk}] + \mathbb{E}[\log \beta_{k, w_{dn}}] \}, k \in \{1, \dots, k\}$$

$$\gamma_d = \alpha + \sum_{n=1}^N \phi_{dn}^k$$

global update.

$$p(\beta_k | z, w) = \text{Dirichlet}(\eta + \sum_d \sum_n z_{dn}^k w_{dn})$$

$$q(\beta_k) = \text{Dirichlet}(\lambda_k)$$

update: $\lambda_k = \eta + \sum_n \phi_{dn}^k w_{dn}$

set $\lambda^{(t)} = (1 - \rho_t) \lambda^{(t-1)} + \rho_t \hat{\lambda}_t$