

STOCHASTIC VARIATIONAL INFERENCE

From Session-6-a, in multinomial mixture model:-

$$q(z_n) = \text{Multinomial}(z_n | v z_n)$$

$$q(\theta) = \text{Dirichlet}(\theta | v\theta)$$

$$q(\beta_k) = \text{Dirichlet}(\beta_k | v\beta_k) \quad k=1, \dots, K$$

$$(A) \quad v z_n \leftarrow \underset{v z_n}{\operatorname{argmax}} \underbrace{E_q \ln P(x, z)}_{\text{ELBO}} - E_q \ln q(z)$$

$$v z_{n,k} \propto \exp \left(E_{q(\theta)} \ln \theta_k + \sum_{w=1}^W m_{n,w} E_{q(\beta_w)} \ln \beta_k \right)$$

$$(B) \quad v\theta \leftarrow \underset{v\theta}{\operatorname{argmax}} E_q \ln P(x, z) - E_q \ln q(z)$$

$$v\theta_k = \alpha - 1 + \sum_{n=1}^N v z_{n,k}$$

$$(C) \quad v\beta_k \leftarrow \underset{v\beta_k}{\operatorname{argmax}} E_q \ln P(x, z) - E_q \ln q(z)$$

$$v\beta_{k,w} = \gamma - 1 + \sum_{n=1}^N m_{n,w} v z_{n,k}$$

* (A), (B) and (C) are solved by

$$\nabla_{\nu z_n} \text{ELBO} = 0, \quad \nabla_{\nu \theta} \text{ELBO} = 0, \quad \nabla_{\nu \beta_k} \text{ELBO} = 0$$

respectively!

- * This is the Batch variational inference
- * To update $\nu \theta$ and $\nu \beta_k$, all the observations have to be analyzed (notice resemblance to Batch EM and (non-collapsed) Gibbs Samplers)
- * Update to global parameters is slow!
- * Can local parameter (νz_n) updates be propagated faster to global parameters ($\nu \theta, \nu \beta_k, k=1, \dots, K$)

* STOCHASTIC OPTIMIZATION !

ALGORITHM :

Given data set of size N and f_t for $t=1, \dots, \infty$

$$\sum p_t = \infty, \quad \sum p_t^2 < \infty, \quad \lim_{t \rightarrow \infty} p_t = 0$$

Initialize global parameters \hat{q} randomly.

for $t: 0, \dots, \infty$

1) Sample an observation

2) Solve the local parameter

3) Construct global parameters \hat{g} from local parameter
of step (2), pretending observation in step (1)
was observed 'N' times.

4) Update global parameters

$$\hat{g}^{(t+1)} = \hat{g}^{(t)} + \eta_t (\hat{g} - \hat{g}^{(t)})$$

learning rate

$$\left\{ \begin{array}{l} \hat{g}^{(t+1)} = (1 - \eta_t) \hat{g}^{(t)} + \eta_t \hat{g} \\ \text{Same as above} \end{array} \right\}$$

running average

SVI for Multinomial Mixture Model :-

Initialize $v \hat{\beta}_k^{(0)}$, $k=1, \dots, K$ and $v \hat{\theta}^{(0)}$ randomly

for $t = 0, \dots, \infty$

1) $x_n \leftarrow$ uniformly sample a document

$$2) \nabla_{\theta} \ln \alpha \propto \exp \left(E_{q(\theta | \nabla \theta^{(t)})} \ln \theta_k + \sum_{w=1}^W m_{n,w} E_{q(\beta_k | \nabla \beta_k^{(t)})} \ln \beta_k \right)$$

$$3) \hat{\nabla \theta}_k = \alpha - 1 + N \nabla_{\theta} \ln \alpha, \forall k, w$$

$$\hat{\nabla \beta}_{k,w} = \gamma - 1 + N m_{n,w} \nabla_{\beta} \ln \alpha, \forall k, w$$

$$4) \begin{aligned} \nabla \theta^{(t+1)} &= (1 - f_t) \nabla \theta^{(t)} + f_t \hat{\nabla \theta} \\ \nabla \beta_k^{(t+1)} &= (1 - f_t) \nabla \beta_k^{(t)} + f_t \hat{\nabla \beta}_k, k = 1, \dots, K \end{aligned}$$

Question 1 :- Why does data-subsampling work?

Question 2 :- What is natural gradient?

Answer 1 :-

Consider :- $\mathbf{z} \sim P(z|\theta)$

$$E_{P(z|\theta)}[z] = \int z P(z|\theta) dz$$

$$= f(\theta)$$

find θ^* such that $f(\theta^*) = 0$

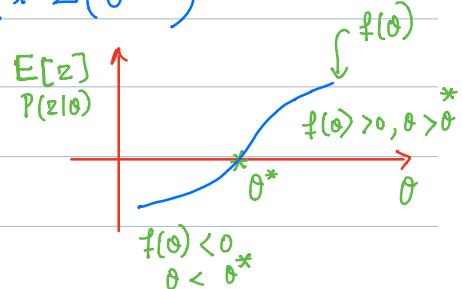
$$\text{i.e. } E_{P(z|\theta^*)}[z] = \int z P(z|\theta^*) dz = 0$$

then iterate in this fashion :-

$$\theta^{(t)} = \theta^{(t-1)} - f_t * z(\theta^{t-1})$$

where $z(\theta) \sim P(z|\theta^{(t-1)})$

$$\text{and } \sum f_t = \infty, \sum f_t^2 < \infty, \lim_{t \rightarrow \infty} f_t = 0$$



- * This root finding procedure (1951, ROBBINS-MONRO) forms the basis of stochastic optimization

Specifically, let's say,

$$\nabla_{\theta} g(\theta) = f(\theta) = E_{P(z|\theta)}[z] = \int z p(z|\theta) dz$$

- * then ' \bar{z} ' has the interpretation of sample of gradient and the above iterative procedure gives the optima of $g(\theta)$

- * In many optimization procedures, $\nabla_{\theta} g(\theta)$ appears as a sum over data-points.

i.e $\nabla_{\theta} g(\theta) = \sum_n \underbrace{\left[\nabla_{\theta} g(\theta) \right]_n}_{\text{contribution of } n\text{-th datapoint to gradient}}$

Under uniform sampling of data point (Step-1 of SVI), we have

$$\begin{aligned} \mathbb{E}_{i \sim \text{uniform}(1, N)} [\nabla_{\theta} g(\theta)]_i &= \sum_{n=1}^N \left(\frac{1}{N} \right) [\nabla_{\theta} g(\theta)]_n \\ &= \frac{1}{N} \sum_{n=1}^N [\nabla_{\theta} g(\theta)]_n \end{aligned}$$

likelihood thickening = $\frac{1}{N} \nabla_{\theta} g(\theta)$

Hence $\mathbf{N}[\nabla_{\theta} g(\theta)]_i$ where $i \sim \text{uniform}(1, N)$

is an unbiased estimate of gradient.

\Rightarrow We can follow $N[\nabla_{\theta} g(\theta)]_i$ instead of batch gradient $\underline{\underline{N[\nabla_{\theta} g(\theta)]_i}}$

Answer 2 :- (What is natural gradient?)

$$F(\lambda) \stackrel{-1}{\nabla_{\lambda}} \text{ELBO}$$

gradient of ELBO w.r.t
variational parameter λ

Two Views:-

Fisher information Matrix.

1) Optimization & 2) information theoretic.

View 1) Optimization

Let $g(\theta)$ be some objective function,

Goal :- optimize $g(\theta)$ using iterative procedure

* Instead of following direction of regular gradient $\nabla_{\theta} g(\theta)$, follow instead e -

$$-\left[H(\theta)\right]^{-\frac{1}{2}} \nabla_{\theta} g(\theta)$$

where $H(\theta)$ is the "Hessian" of objective function

defined as $H(\theta)[i,j] = i,j$ entry of Hessian matrix

$$= \frac{\partial^2}{\partial \theta_i \partial \theta_j} g(\theta)$$

(Well-known Newton method)

* Fisher-information matrix plays the same role in optimizing parameters of some probability distribution (e.g $v\theta$)

$$\text{i-e } F(\lambda)[i,j] = -E_{q(\beta|\lambda)} \frac{\partial^2}{\partial \lambda_i \partial \lambda_j} \ln q(\beta|\lambda)$$

true
for exponential
families.

$$\Rightarrow E_{q(\beta|\lambda)} \frac{\partial}{\partial \lambda_i} \ln q(\beta|\lambda) \cdot \frac{\partial}{\partial \lambda_j} \ln q(\beta|\lambda)$$

Hence use $\nabla_{\lambda}^{-1} \text{ELBO}$

where λ is some variational parameter.

Most IMPORTANT :-

$$\nabla_{\lambda} \text{ELBO} = F(\lambda)(\lambda - E_q \text{ complete conditional natural parameters})$$

Hence $F(\lambda)^{-1}$ is not needed to be computed

PROOF :-

X :- data, and Z, θ - two latent variables.

$$\begin{aligned} \Pr(z | hz) &= a(z) b(hz) \exp(hz^T t(z)) \\ \Pr(\theta | h\theta) &= a(\theta) b(h\theta) \exp(h\theta^T t(\theta)) \end{aligned} \quad \left. \right\} \text{Priors}$$

$$\Pr(\theta | h\theta, x, z) = a(\theta) b(\eta(h\theta, x, z)) \exp(\eta^T(h\theta, x, z) t(\theta)) \quad \text{--- A}$$

$$\Pr(z | hz, x, \theta) = a(z) b(\eta(hz, x, \theta)) \exp(\eta^T(hz, x, \theta) t(z)) \quad \text{--- B}$$

(A) and (B) are $\Pr(\text{latent variable} \mid \text{everything else besides latent-variable})$

Let $q_z(z|vz)$ and $q_\theta(\theta|v\theta)$ be variational distributions for z and θ in Mean-Field Approx.

$$q_z(z|vz) = a(z) b(vz) \exp(vz^T t(z))$$

$$q_\theta(\theta|v\theta) = a(\theta) b(v\theta) \exp(v\theta^T t(\theta))$$

To find vz and $v\theta$, we maximize ELBO

$$L(vz, v\theta) = E_{q_z} \ln P(x, z, \theta | vz, v\theta) - E_{q_z} \ln q_z(z, \theta)$$

$$= E_{q_z} \ln P(x, z | vz, v\theta) + E_{q_z} \ln P(\theta | x, z, v\theta)$$

$$- E_{q_z} \ln q_z(z|vz) - E_{q_z} \ln q_\theta(\theta|v\theta)$$

optimizing w.r.t $v\theta$, let's write ELBO in terms of $v\theta$ dependent terms.

$$L(v\theta) = E_{q_\theta} \ln P(\theta | x, z, v\theta) - E_{q_\theta} \ln q_\theta(\theta|v\theta)$$

using A

$$= E_q \left[\ln a(\theta) + \ln b(\eta(x, z, \theta)) + \right.$$

$$\eta^T(x, z, \theta) t(\theta) - \ln a(\theta) - \ln b(v\theta)$$

$$- v\theta^T t(\theta) \Big]$$

ignore as no $v\theta$ involved

$$= E_q \left[\ln b(\eta(x, z, \theta)) \right] + E_{q(2)} \left[\eta(x, z, \theta) \right]^T E_{q(1)} t(\theta)$$

$$- \ln b(v\theta) - (v\theta)^T E_q(t(\theta))$$

$$L(v\theta) = -E_{q(2)} \left[\eta(x, z, \theta) \right]^T \nabla_{v\theta} \ln b(v\theta) + (v\theta)^T \nabla_{v\theta} \ln b(v\theta)$$

$$- \ln b(v\theta)$$

$$\bar{\nabla}_{v\theta} L(v\theta) = - \left[\nabla_{v\theta}^2 \ln b(v\theta) \right] E_{q(2)} \left[\eta(x, z, \theta) \right] +$$

$$\left[\nabla_{v\theta}^2 \ln b(v\theta) \right] v\theta + \nabla_{v\theta} \ln b(v\theta) - \nabla_{v\theta} \ln b(v\theta)$$

$$\bar{\nabla}_{v\theta} L(v\theta) = - \left[\nabla_{v\theta}^2 \ln b(v\theta) \right] \left(E_{q(2)} \left[\eta(x, z, \theta) \right] - v\theta \right)$$

for exponential family, this is $F(v\theta)$

2nd View :- information theoretic

steepest descent direction $\leftarrow \underset{d\theta}{\operatorname{argmax}} f(\theta + d\theta)$

such that $\|d\theta\|_2 < \epsilon$

steepest descent direction in euclidean. $\leftarrow \frac{f'(\theta)}{\|f'(\theta)\|_2} \epsilon = \text{direction of gradient at } \theta$

natural gradient $\leftarrow \underset{d\theta}{\operatorname{argmax}} f(\theta + d\theta)$

such that $\text{KL}(p(\cdot | \theta) || p(\cdot | \theta + d\theta)) < \epsilon$