

Variational Inference (VI)

①

- In Bayesian learning, finding posterior $p(z|x)$ is difficult or finding moments w.r.t. $p(z|x)$ is difficult.
- So we need approximations
 - VI is such an approximation.
 - There are other techniques : (a) Sampling
(b) Expectation propagation
- Sampling (or MCMC) is quite accurate or in fact exact if we can sample many.
- Goal of VI: To find best possible approximation of the exact posterior distribution.
$$q(z) \approx p(z|x)$$
- How to do that? By some optimality criterion.

For VI, we need some assumption:

- We must need to have (or formulate) $p(x|z)$.
- We also need to have priors $p(z)$.
- So we have $p(x, z) = p(x|z) p(z)$ expression.

Example: We use GMM case.

prior: $p(z) = \prod_{k=1}^K \pi_k^{z_k}$ z is a 1-of- K representation vector.

conditional distribution: $p(x|z_{k=1}) = \mathcal{N}(x|\mu_k, \Sigma_k)$, $z_k \in \{0, 1\}$
 or $p(x|z) = \prod_{k=1}^K \mathcal{N}(x|\mu_k, \Sigma_k)^{z_k}$

\therefore ~~$p(x, z)$~~ $p(x, z) = p(x|z) \cdot p(z)$

$p(x|\mu, \Sigma, \pi) = \sum_z \prod_k \pi_k \mathcal{N}(x|\mu_k, \Sigma_k)^{z_k}$

$= \prod_{k=1}^K \left\{ \pi_k \mathcal{N}(x|\mu_k, \Sigma_k) \right\}^{z_k}$

$x = [x_1 \ x_2 \ \dots \ x_N]$ $z = [z_1 \ z_2 \ \dots \ z_N]$

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

$p(x, z|\pi, \mu, \Sigma) = \prod_{n=1}^N \prod_{k=1}^K \pi_k^{z_{n,k}} \mathcal{N}(x_n|\mu_k, \Sigma_k)^{z_{n,k}}$

- Now, if we need $p(z|x, \pi, \mu, \Sigma)$, that is a bit trouble.
- Just for the sake of an example, we want $q(z)$ that best approximates $p(z|x, \text{other parameters})$.
- Another example: $z = [z_1^T \ z_2^T \ \dots \ z_Q^T]^T$. Then if we want $p(z_i|x, \dots) = \int_{z_j, j \neq i} p(z|x, \pi, \mu, \Sigma) dz_j$, this is more difficult due to complex dependencies between $\{z_i\}$.

Optimization formulation

Notation: $E_q [h(z)] = \int q(z) h(z) dz$

$$\ln p(x) = E_q [\ln p(x)] \quad \rightarrow \text{because } p(x) \text{ does not depend on } q(z). \text{ The equality is valid for any } q.$$

$$= E_q \left[\ln \frac{p(x, z)}{p(z|x)} \right]$$

$$= E_q \left[\frac{\ln p(x, z)}{q(z)} \right] - E \left[\frac{\ln p(z|x)}{q} \right]$$

$$= E_q \left[\ln \frac{p(x, z)}{q(z)} \cdot \frac{q(z)}{p(z|x)} \right]$$

$$= \underbrace{E_q \left[\ln \frac{p(x, z)}{q(z)} \right]}_{L(q)} + \underbrace{E_q \left[\ln \frac{q(z)}{p(z|x)} \right]}_{KL(q \parallel p(z|x))}$$

$$= L(q) + KL(q \parallel p(z|x)).$$

As $KL(q \parallel p(z|x)) \geq 0$, $L(q)$ works as a LB.

- Note $\ln p(x)$ is fixed. So, if we define the optimization problem (cost) as

$$q^* = \arg \max_q L(q) = \arg \max_q E_q \left[\ln \frac{p(x, z)}{q(z)} \right],$$

then the $KL(q^* \parallel p(z|x))$ will be minimum.

By this logic q^* approaches to $p(z|x)$ or

$$q^* \approx p(z|x).$$

(4)

Factorized approximations

$$z = (z_1^* z_2^* \dots z_M^*)^T$$

$$P(z|x) \approx q(z) = \prod_{i=1}^M q_i(z_i).$$

Assumption: z_i 's are statistically independent.

Task: find $\{q_i^*\}_{i=1}^M$. Then for $q^* = \prod_{i=1}^M q_i^*$.

Statement: Optimal density for each group (i.e. q_i^*) can be found iteratively, by choosing

$$\begin{aligned} \ln q_i^*(z_i) &= E_{q_{j \neq i}} [\ln p(x, z)] + \text{const} \\ &= E_{q_{j \neq i}} [\ln p(x, z_1, z_2, \dots, z_{i-1}, z_i, z_{i+1}, \dots, z_M)] + c \end{aligned}$$

→ Here 'c' is a normalization constant, and $E_{q_{j \neq i}}[\cdot]$ means that the q_j 's are kept fixed for $j \neq i$, and these q_j 's are used to calculate expectation over all those other group of variables z_j , except z_i .

→ For $i=1, 2, \dots, M$, this is repeated for ~~each~~ each q_i , keeping all other $q_{j \neq i}$ fixed.

→ Optimality: $\{q_i^*\} = \arg \max_{\{q_i\}} E_{q_i} \left[\ln \frac{p(x, z)}{\prod_{i=1}^M q_i(z_i)} \right]$

→ $q_i^*(z_i) = \frac{\exp [E_{q_{j \neq i}} [\ln p(x, z)]]}{\int_{z_i} \exp [E_{q_{j \neq i}} [\ln p(x, z)]] dz_i}$.

⑤

Theorem: (Theorem 9.2 of Arne's compendium).

Assuming $z = [z_1, z_2]$, and given $p(z, x) = p(z_1, z_2, x)$, we want to have $q(z) = q_1(z_1) q_2(z_2) \approx p(z|x)$. Then, for any fixed q_2 , ~~we~~ a density function q_1 , obtained

as

$$\begin{aligned} \ln q_1(z_1) &= E_{q_2} [\ln p(z, x)] + c \\ &= E_{q_2} [\ln p(z_1, z_2, x)] + c \quad \dots \textcircled{1} \end{aligned}$$

maximizes $L(q) = E_q \left[\ln \frac{p(z, x)}{q(z)} \right]. \quad \dots \textcircled{2}$

proof: • After taking expectation over z_2 in $\textcircled{1}$, the remaining expression is a function of z_1 . With proper normalization, this function can be interpreted as the logarithm of a density function, called $\tilde{p}(z_1)$ here.

$$\begin{aligned} \text{So, } E_{q_2} [\ln p(z_1, z_2, x)] &= \int q_2(z_2) \ln p(z_1, z_2, x) dz_2 \\ &= f(z_1, x) + \text{const} \\ &\stackrel{\uparrow}{=} \ln \tilde{p}(z_1) + \text{const} \\ &\quad x \text{ is fixed and } f \geq 0. \end{aligned}$$

$$\begin{aligned} \text{Now, } L(q) &= \int_{z_1} \int_{z_2} q_1(z_1) q_2(z_2) \ln p(z, x) dz_1 dz_2 - \int_{z_1} \int_{z_2} q_1(z_1) q_2(z_2) [\ln q_1(z_1) + \ln q_2(z_2)] dz_1 dz_2 \\ &= \int_{z_1} q_1(z_1) \left[\int_{z_2} q_2(z_2) \ln p(z, x) dz_2 \right] dz_1 - \int_{z_1} q_1(z_1) \ln q_1(z_1) \underbrace{\int_{z_2} q_2(z_2) dz_2}_1 dz_1 \\ &\quad - \int_{z_1} q_1(z_1) dz_1 \underbrace{\int_{z_2} q_2(z_2) \ln q_2(z_2) dz_2}_{\text{const}} \\ &= \int_{z_1} q_1(z_1) \ln \tilde{p}(z_1) + \text{const} \underbrace{\int_{z_1} q_1(z_1) dz_1}_1 - \int_{z_1} q_1(z_1) \ln q_1(z_1) dz_1 - \text{const} \\ &= \int_{z_1} q_1(z_1) \ln \frac{\tilde{p}(z_1)}{q_1(z_1)} + \text{const} = - \int_{z_1} q_1(z_1) \ln \frac{q_1(z_1)}{\tilde{p}(z_1)} + \text{const} \\ &= - \text{KL}(q_1 \parallel \tilde{p}) + \text{const}. \end{aligned}$$

⑥

• Note that $q_1(z_1) = \tilde{p}(z_1)$ leads to $KL=0$, and hence minimizes $L(q)$.

• So the optimal choice is $\ln q_1(z_1) = \ln \tilde{p}(z_1)$
 $= E_{q_2}(\ln p(z, x)) + \text{const.}$
 (Theorem proved).

• This proof covers general statement:

$$\ln q_i^*(z_i) = E_{q_{j \neq i}} [\ln p(x, z)] + \text{const.}$$

where we use $q_i = q_1$ and $q_{j \neq i} = q_2$.

⑦

Comment on the use of VI for Gaussian mixture (section 10.2 of Bishop Book).

It is shown that the use of VI with factorized distribution leads to two optimization steps analogous to E and M steps of the maximum likelihood EM algorithm.

EM - Special case of VI (from Arne's compendium)

• Let us explore relation between EM and VI.

• ~~First~~ First, we recall the EM algo. (General EM)

Given a joint distribution $p(x, z | \theta)$ over observed variables x and latent variables z , governed by the parameters θ , the goal is to maximize likelihood $p(x | \theta)$ w.r.t. θ . That is $\theta^* = \arg \max_{\theta} p(x | \theta)$.

1. Choose an initial θ^{old} .

2. E Step: Evaluate $p(z | x, \theta^{\text{old}})$.

3. M Step: Define $Q(\theta, \theta^{\text{old}}) = \sum_z p(z | x, \theta^{\text{old}}) \ln p(x, z | \theta)$

$$\theta^{\text{new}} = \arg \max_{\theta} Q(\theta, \theta^{\text{old}}).$$

4. Check convergence. If not converged $\theta^{\text{old}} \leftarrow \theta^{\text{new}}$.

• Now we investigate VI & EM connection

Assumptions: (a) $P(x|z, \theta)$ is known

(b) Prior $P(z|\theta)$ and $P(\theta)$ are known.

$\therefore \ln P(x, z, \theta) = \ln P(x|z, \theta) P(z|\theta) P(\theta)$ is known.

→ Let us assume the goal is to estimate a point estimate $\hat{\theta}$.
(in analogy with the EM)

↓ VI factorization

$$\rightarrow \phi(\theta, z|x) \approx q(\theta, z) = q_1(\theta|\hat{\theta}) q_2(z|\hat{\theta}) = q_1(\theta|\hat{\theta}, x) q_2(z|\hat{\theta}, x)$$

$\hat{\theta}$ is a hyper-parameter
Why $q_1(\theta|\hat{\theta})$? We force the posterior density q_1 to be very sharply peaked at $\hat{\theta}$. So we model

$$q_1(\theta|\hat{\theta}) = \prod_k \frac{1}{e_k} \delta\left(\frac{\theta_k - \hat{\theta}_k}{e_k}\right) \rightarrow \text{Dirac Delta}$$

$$\boxed{E_{q_1}[g(\theta)] = \int g(\theta) q_1(\theta|\hat{\theta}) d\theta \rightarrow g(\hat{\theta})}$$

→ Strategy: In each step of VI iteration, we apply the general factorization solution $\ln q_i^*(z_i) = E_{q_{j \neq i}}[\ln P(x, z, \theta)]$
where we assume $z = [z, \theta]$.

(a) Let us have $\hat{\theta} = \theta^{\text{old}}$.

(b) $q_1(\theta|\hat{\theta}) = q_1(\theta|\theta^{\text{old}})$.

(c) $\ln q_2(z|\hat{\theta}) = E_{q_1}[\ln P(x, z, \theta)] + \text{const} = \ln P(x, z, \theta^{\text{old}}) + \text{const}$.

Using proper normalization, $q_2(z|\hat{\theta}) \propto P(x, z, \theta^{\text{old}}) \propto P(z|x, \theta^{\text{old}})$
Combine with E step of EM.

(d) To find a new point estimate $\hat{\theta} = \theta'$, we have $q_1(\theta|\hat{\theta}) = q_1(\theta|\theta')$.

Use the relation by lower bound

$$\theta^{\text{new}} = \arg \max_{\theta'} L(q) = \arg \max_{\theta'} E_{q_1} \left[\ln \frac{P(x, z, \theta)}{q_1(\theta|\theta') \cdot q_2(z|\theta')} \right]$$

$$\begin{aligned}
 L(\alpha) &= E_q \left[\ln \frac{P(x, z, \theta)}{q_1(\theta|\theta') q_2(z)} \right] \\
 &= \int_{\theta, z} q_1(\theta|\theta') q_2(z) \left[\ln P(x, z, \theta) - \ln q_1(\theta|\theta') - \ln q_2(z) \right] d\theta dz \\
 &= \int_z q_2(z) \left[\int_{\theta} q_1(\theta|\theta') \ln P(x, z, \theta) d\theta \right] dz \\
 &\quad - \int_{\theta} q_1(\theta|\theta') \ln q_1(\theta|\theta') d\theta \int_z q_2(z) dz \\
 &\quad - \int_{\theta} q_1(\theta|\theta') d\theta \int_z q_2(z) \ln q_2(z) dz \\
 &= \int_z q_2(z) \ln P(x, z, \theta') dz - \int_{\theta} q_1(\theta|\theta') \ln q_1(\theta|\theta') d\theta \\
 &\quad - \int_z q_2(z) \ln q_2(z) dz
 \end{aligned}$$

Note: (i) The first term depends on θ' and θ^{old} (because $q_2(z)$ is a function of θ^{old}). We express the first term as $c(\theta', \theta^{\text{old}})$.

(ii) The second term is constant. Because the change in θ' only leads to change in peak for Dirac distribution $q_1(\theta|\theta')$.

(iii) The third term is independent of θ' .

So, maximization of $L(\alpha) \equiv$ maximization of $c(\theta', \theta^{\text{old}})$.

From step (e) $q_2(z|x) = P(x, z, \theta^{\text{old}}) \cdot e^{\text{const}}$

$$= e_1 P(x, z, \theta^{\text{old}})$$

$$= e_1 P(z|x, \theta^{\text{old}}) \cdot P(x, \theta^{\text{old}})$$

Now, x is fixed and θ^{old} is fixed, then $P(x, \theta^{\text{old}}) = \text{constant}$

$$\therefore q_2(z|x) = e_1 P(z|x, \theta^{\text{old}}) \cdot x \text{ constant}$$

$$= e_2 P(z|x, \theta^{\text{old}})$$

$$\int_z q_2(z|x) dz = 1 \quad \text{and} \quad \int_z P(z|x, \theta^{\text{old}}) dz = 1$$

$$\Rightarrow e_2 = 1 \quad \text{or,} \quad q_2(z|x) = P(z|x, \theta^{\text{old}})$$

$$\begin{aligned}\theta^{\text{new}} &= \arg \max_{\theta'} L(\theta') \\ &= \arg \max_{\theta'} c(\theta', \theta^{\text{old}})\end{aligned}$$

$$\begin{aligned}c(\theta', \theta^{\text{old}}) &= \int_{\mathbf{z}} q_{\mathbf{z}}(\mathbf{z} | x, \theta^{\text{old}}) \ln P(x, \mathbf{z}, \theta') d\mathbf{z} \\ &= \int_{\mathbf{z}} P(\mathbf{z} | x, \theta^{\text{old}}) \ln P(x, \mathbf{z}, \theta') \cdot P(\theta') d\mathbf{z} \\ &= \int_{\mathbf{z}} P(\mathbf{z} | x, \theta^{\text{old}}) \ln P(x, \mathbf{z}, \theta') d\mathbf{z} + \int_{\mathbf{z}} P(\mathbf{z} | x, \theta^{\text{old}}) \ln P(\theta') d\mathbf{z}\end{aligned}$$

Compare with
M-step of \rightarrow
EM.

Points: (i) If $P(\theta')$ is a non-informative prior ~~there~~, such as constant, then $\ln P(\theta')$ is constant. We have the same solution as general EM (M-step) for ML.

(ii) If $P(\theta')$ is non-constant, then we have MAP.

② If not converged, $\theta^{\text{old}} \leftarrow \theta^{\text{new}}$.

Conclusion: (i) VI and EM both lead to same computational algorithm.
(ii) Bayesian VI considers θ as random. But EM assumes θ some unknown non-random quantity.
(iii) So, VI is general than EM and EM can be viewed as a special case of VI.

Some comments of Gibbs sampling and Expectation propagation.