

4. Variational method. EM algorithm as a variational method. Variational Bayes.

Def. KL-divergence:  $D_{KL}(Q \parallel P) = \sum_x Q(x) \log \frac{Q(x)}{P(x)}$ ;

$$D_{KL}(Q \parallel P) \geq 0, D_{KL}(Q \parallel P) = 0 \Leftrightarrow Q = P. D_{KL}(Q \parallel P) \neq D_{KL}(P \parallel Q).$$

~~Suppose we are to derive the posterior  $P(\theta | \{x_n\}_{n=1}^N)$~~

~~Given prior  $p(\theta)$ .~~

Suppose we have an assumed model  $\mathcal{H}$ , a set of data points  $\{x^{(n)}\}_{n=1}^N$ , a set of latent variables

$\{z^{(n)}\}_{n=1}^N$  and a set of parameters  $\theta$ .

(e.g.,  $z^{(n)}$  is a class label in k-means)

Our task is to get the best (MAP) parameters  $\theta_{ML}$ , given  $x$  and  $\mathcal{H}$ .

~~log~~  $\left( \prod_{i=1}^n \pi_{z_i} \mathcal{N}(x_i | \mu_{z_i}, \sigma_{z_i}^2) \right)$  e.g.)

1) If we know  $p(x, z | \theta, \mathcal{H})$ , we can do the following:

$$\log p(x | \theta, \mathcal{H}) = \int q(z) \log p(x | \theta, \mathcal{H}) dz =$$

$$= \int q(z) \log \frac{p(x, z | \theta, \mathcal{H})}{q(z)} dz + \int q(z) \log \frac{q(z)}{p(z | \theta, \mathcal{H})} dz =$$

$$\equiv \cancel{D_{KL}(q \parallel p(x, z | \theta, \mathcal{H}))} +$$

$$= \underbrace{\int q(z) \log \frac{p(x, z | \theta, \mathcal{H})}{q(z)} dz}_{\text{Variational lower bound}} + \underbrace{D_{KL}(q \parallel p(z | x, \theta, \mathcal{H}))}_{\geq 0}$$

||  
 $\mathcal{L}_V(q, \theta)$

2) As a result, instead of ~~solving~~

$$\mathcal{L}(\theta) \equiv \cancel{\log p(x | \theta, \mathcal{H})} \rightarrow \max_{\theta}$$

we solve  $L_V(q, \theta) \rightarrow \max_{q, \theta}$

by moving step-by-step on  $q$  and  $\theta$ :

$$(M) \quad \theta^{(t)} = \arg \max_{\theta} L_V(q^{(t-1)}, \theta^{(t-1)}) \stackrel{!}{=} \arg \max_{\theta} \mathbb{E} \log p(x, z | \theta, H)$$

$$(E) \quad q^{(t)} = \arg \max_q L_V(q, \theta^{(t)}) = p(z | x, \theta^{(t)}, H) \\ \left( \arg \min_q D_{KL}(q, p(z | x, \theta^{(t)}, H)) \right)$$

and we get EM algorithm (~~in general~~) (as a general idea).

3) If all  $z_i \in \{1, \dots, k\}$ , then:

$$p(x_i | \theta) = \sum_{k=1}^k p(x_i | k, \theta) p(z_i = k | \theta),$$

and E-step becomes:

$$q^{(t)}(z_i = k) = \frac{p(x_i | k, \theta^{(t)}) p(z_i = k | \theta^{(t)})}{\sum_{\ell=1}^k p(x_i | \ell, \theta^{(t)}) p(z_i = \ell | \theta^{(t)})}$$

while M-step is: 
$$\theta^{(t)} = \sum_{i=1}^n \sum_{k=1}^k q^{(t-1)}(z_i = k) \log p(x_i, k | \theta^{(t-1)}).$$

4) if  $z_i$  is continuous, E-step can be done in closed form only in case of conjugate distributions.

(because we compute  $q^{(t)}(z_i) = \frac{p(x_i | z_i, \theta^{(t)}) p(z_i | \theta^{(t)})}{\int p(x_i | z, \theta) p(z | \theta) dz}$ )

5) To deal with non-conjugate priors we use Variational Bayes:

instead of  $q^{(t)} = \arg \min_q D_{KL}(q || p(z | x, \theta^{(t)}, H))$

we approximate

it with  $q^{(t)} = \arg \min_{q \in Q} D_{KL}(q || p(z | x, \theta^{(t)}, H))$

Inference becomes optimization.

$$\arg \max_{q \in Q} L_V(q, \theta^{(t)}).$$

$$\underset{q \in Q}{\arg \max} L_V(q, \theta^{(t)}) = \int q(z | \phi) \log \frac{p(x, z | \theta)}{q(z | \phi)} dz \rightarrow \max_{\phi} - \text{Variational E-step.}$$

4

c) In classical EM:

d) Variational Bayes with classical EM:

$$\varphi_n = \arg \max_{\varphi} \hat{L}_V(\varphi, \theta_{n-1}), \text{ E-step}$$

$$\theta_n = \arg \max_{\theta} \hat{L}_V(\varphi_n, \theta), \text{ M-step}$$

If we cannot solve this problem analytically we'd better use:

$$\varphi_n = \varphi_{n-1} + \eta \nabla_{\varphi} \hat{L}_V(\varphi, \theta_{n-1}), \text{ Variational E-step}$$

$$\theta_n = \theta_{n-1} + \varepsilon \nabla_{\theta} \hat{L}_V(\varphi_n, \theta), \text{ Variational M-step}$$

We can even use the stochastic gradients instead of the full  $\hat{L}_V$  - ELBO (Evidence Lower-Bound)

ELBO is often intractable, but:

$$\log p(x, z | \theta) = \sum_{i=1}^n \log p(x_i, z_i | \theta) = \sum_{i=1}^n [\log p(x_i | z_i, \theta) + \log p(z_i | \theta)]$$

if  $\{x_i, z_i\}_{i=1}^n$  is i.i.d.  $\Rightarrow$  can use mini-batching for unbiased gradient estimations.

Example (Ada-gram)

Word2Vec with Hierarchical softmax:

$$p(y | x, \theta) = \prod_{c \in \text{Path}(y)} \sigma(d_{c,y} \ln(x)^T \text{Out}(c))$$

$$p(y | x, \theta) \rightarrow \max_{\theta}, \quad \theta = [\ln, \text{Out}]$$

Path(y),  $d_{c,y}$  are obtained from the Huffman tree, build for our dictionary.

How to account different meanings of the same word?  
 Let's define a latent variable  $z_i$  that indicates the meaning of  $x_i$ :

$$p(y_i | x_i, z_i, \theta) = \prod_{c \in \text{path}(y_i)} \sigma(d_{c, y_i} \ln(x_i, z_i)^T \text{Out}(c)),$$

$$p(z_i = k | x_i) = \frac{1}{K(x_i)}, \quad K(x_i) - \text{total number of meanings for } x_i.$$

Now we can use a standard EM-algorithm with discrete latent variables (as in 3)...

But that we'll require a lot of time.

To make it scalable - use Variational EM:

M-step:

$$\begin{aligned} \nabla_{\theta} \mathbb{E}_z \log p(y, z | x, \theta) &= \nabla_{\theta} \mathbb{E}_z \sum_{i=1}^n (\log p(y_i | z_i, x_i, \theta) + \log p(z_i | x_i)) = \\ &= \sum_{i=1}^n \mathbb{E}_{z_i} (\nabla_{\theta} \log p(y_i | z_i, x_i, \theta) + \nabla_{\theta} \log p(z_i | x_i)) = \end{aligned}$$

doesn't depend on  $\theta$ .

$$= \sum_{i=1}^n \mathbb{E}_{z_i} (\nabla_{\theta} \log p(y_i | z_i, x_i, \theta))$$

Its unbiased estimate is  $\frac{1}{n} \sum_{j=1}^{K(x_i)} \underbrace{p(z_i = k | y_i, x_i, \theta)}_{\text{we know from E-step}} \nabla_{\theta} \log p(y_i | k, x_i, \theta)$

~~During M-step, we update~~

During E-step we update ~~only~~ probabilities only for single  $x_i$ .

To deal with different # of meanings for each  $x_i$  use DP:

$$p(z = k | x_i, \vec{\beta}) = \beta_{ik} \prod_{r=2}^{k-1} (1 - \beta_{ir}), \quad p(\beta_{ik} | \mathcal{L}) = \text{Beta}(\beta_{ik} | 1, \mathcal{L})$$

$$p(y, z, \beta | x, \mathcal{L}, \theta) = \prod_{i: k=1}^{|V|} \prod_{k=2}^{\infty} p(\beta_{ik} | \mathcal{L}) \cdot \prod_{i=1}^n \left[ p(z_i | x_i, \vec{\beta}) \prod_{j=1}^c p(y_{ij} | z_i, x_i, \theta) \right]$$