

# Variational Bayes

①

Consider a model with observed data  $x$  and hidden variables  $z$ .

$$p(x, z) = p(x|z) p(z)$$

And the marginal likelihood is

$$p(x) = \int_Z p(x, z) dz = E_{p(z)} [p(x|z) \cdot p(z)]$$

Goal: find posterior  $p(z|x)$ .  
~~Bayesian~~

Idea: - posit an approximation from some family  $Q$ .

- Find the member of the family that minimizes the KL divergence to the true posterior.

$$\begin{aligned} q^*(z) &= \arg \min_{q(z) \in Q} KL [q(z) \parallel p(z|x)] = E_{q(z)} \left( \log \frac{q(z)}{p(z|x)} \right) \\ &= E_z (\log q(z)) - E_z (\log p(z|x)) \end{aligned}$$

Reminder:  ~~$KL(p \parallel q)$~~

$$\begin{aligned} KL(P_L \parallel P_R) &= \int_{\Theta} P_L(\theta) \cdot \log \frac{P_L(\theta)}{P_R(\theta)} d\theta \\ &= \text{expected log ratio under } P_L \end{aligned}$$

2)

- We find this  $q^*(z)$  by optimization.
- It only ever approximates the true posterior... underestimates variance
- less precise; usually faster

In general, we have

$$p(z|x) = \frac{p(z)p(x|z)}{p(x)} = \frac{p(z)p(x|z)}{\int_Z p(x,z) dz} \leftarrow \text{hard}$$

$z$  could include "parameters" as well as "latent variables"

Example  $(x_i | c_i, \mu) \sim N(c_i^T \mu, 1)$

$c_i \sim \text{Multinomial}(M=1, w = (\frac{1}{K}, \dots, \frac{1}{K}))$  Mixture of normals

$\mu_k \sim N(0, \tau^2)$

→ "one-hot" encoding of the mixture component

"Params":  $\mu, \sigma^2$

"Latents":  $c_i$  for  $i=1, \dots, n$

The joint dist of (data, pars) is

$$p(\mu, \sigma^2, c, x) = p(\mu, \sigma^2) \cdot \prod_{i=1}^N p(c_i) p(x_i | c_i, \mu)$$

The marginal likelihood ("evidence") is

$$p(x) = \int p(\mu, \sigma^2) \prod_{i=1}^N \sum_{c_i} p(c_i) p(x_i | c_i, \mu) d\mu$$

## The evidence lower bound (ELBO)

⑤

Family of possible approximations:  $Q$ .

Each  $q(z) \in Q$  is a candidate.

Goal: compute  $\arg \min_{q(z) \in Q} KL [q(z) \parallel p(z|x)]$

Problem: this is not even computable:

$$KL(q(z) \parallel p(z|x)) = E_q \left( \log \frac{q(z)}{p(x|z)} \right)$$

$$* = E_q [\log q(z)] - E_q [\log p(x|z)]$$

$$= E_q [\log q(z)] - E_q [\log p(x, z)] + \log p(x)$$

This depends on  $\log p(x)$ , which we can't compute.

So instead consider the functional

$$\text{ELBO}(q) = E_q [\log p(x, z)] - E_q [\log q(z)]$$

$$\text{Clearly } ELBO(q) = -KL [q(z) \parallel p(z|x)] + \log p(x)$$

~~So maximizing~~ (hence ELBO)

So maximizing ELBO is equivalent to minimizing the KL divergence in (\*)

Let's rewrite the ELBO:

$$\text{ELBO}(q) = E[\log p(z)] + E[\log p(x|z)] - E[\log q(z)]$$

$$= \underbrace{E[\log p(x|z)]}_{\text{Expected log-likelihood of data under } q} - \underbrace{\text{KL}[q(z) \parallel p(z)]}_{\text{KL b/t } q(z) \text{ and prior } p(z)}$$

Expected log-likelihood of data under  $q$

under  $q$

"Explain the data"

KL b/t  $q(z)$  and prior  $p(z)$

"Don't stray from the prior."

## Mean-field family

Key question: what is  $Q$ ?

The mean-field family are all  $q(z)$  of the form

$$q(z) = \prod_{j=1}^M q_j(z_j)$$

Note: not a model of the observed data (no  $x$ !)

The ELBO connects  $x$  to the  $q(z_j)$ 's.

- In principle, can take any parametric form to match desired form of each parameter

- Sometimes, we can find an "optimal" form for  $q_j$  based on the model.

## Back to mixture of normals example

(5)

The mean-field family consists of models of this form:

$$g(\mu, \sigma^2, c) = g(\sigma^2) \cdot \underbrace{\prod_{k=1}^K g(\mu_k | m_k, s_k^2)}_{\text{Gaussian prior distn for } \mu_k} \cdot \prod_{i=1}^N g(c_i | \phi_i)$$

K vector of probabilities

- We have simply asserted that these are Gaussians / categorical
- But in fact these choices are provably optimal (later)

Algorithm = coordinate ascent variational inference.

Consider the  $j^{\text{th}}$  latent variable  $z_j$ .

It's complete conditional is  $p(z_j | z_{-j}, x)$ ,

Fact: With  $z_{-j}$  fixed, the optimal  $q_j(z_j)$  is of the form

(prove later)

$$q_j^*(z_j) \propto \exp \left\{ E_{-j} [\log p(z_j | z_{-j}, x)] \right\}$$

or equivalently, working with the joint,

$$q_j^*(z_j) \propto \exp \left\{ E_{-j} [\log p(z_j, z_{-j}, x)] \right\}$$

This expectation is w.r.t. the variational density over  $z_{-j}$ ,  
~~with~~ ie.  $\prod_{l \neq j} q_l(z_l)$

Remember: all latent variables are independent. So these expectations don't involve the  $j^{\text{th}}$  ~~variable~~ variational factor.

$$\log q_j^*(z_j) = E_{-j} [\log p(z_j | z_{-j}, x)]$$

So coordinate ascent cycles through these variational factors, updating them one at a time. (Remember, we're updating the variational parameters)

Derivation of the key fact:

$$\text{ELBO}(q) = E_q[\log p(x, z)] - E_q[\log q(z)]$$

$$\text{And } q(z) = \prod_{j=1}^M q_j(z_j)$$

So write this isolating  $q_j(z_j)$ :

$$\text{ELBO}(q_j) = E_{q_j} \left[ \underbrace{E_{q_{-j}}[\log p(x, z_j, z_{-j})]}_{\text{(Iterated expectations)}} \right] - E_{q_j}[\log q_j(z_j)] + \text{constant}$$

all  $E_{q_{-j}}$ 's in here

Focus on this. This is, up to constant,

$$\text{ELBO}(q_j) = E_{q_j}[\log q_j^*(z_j)] - E_{q_j}[\log q_j(z_j)]$$

$$\text{where } q_j^*(z_j) = \exp \left\{ E_{q_{-j}}[\log p(x, z_j, z_{-j})] \right\}$$

$$= -KL(q_j \| q_j^*)$$

So clearly we maximize  $\text{ELBO}(q_j)$  by minimizing

$KL(q_j \| q_j^*)$ , which happens when  $q_j(z_j) = q_j^*(z_j)$

QED

①

# Return to the Gaussian mixture-model example

$x_i \in \mathbb{R}$ , data point

$\mu_k$ :  $k^{\text{th}}$  mean

$c_i$ : indicator vector w/  $c_{ik} = 1$  if  $x_i$  came from component  $k$

$\tau^2$ : assume fixed  $\mu_k \sim N(0, \tau^2)$

Mean-field family:  $q(z) = \prod_{j=1}^M q_j(z_j)$

$$= \prod_{k=1}^K N(\mu_k | m_k, s_k^2) \cdot \prod_{i=1}^N p(c_i | \phi_i)$$

$$ELBO(q) = E_q \left[ \log \left( \prod_{i=1}^N N(x_i | \mu, c) \right) p(\mu, c) \right] - E_q \left[ \log \left[ \prod_{k=1}^K q(\mu_k | m_k, s_k^2) \right] + \log \prod_{i=1}^N p(c_i | \phi_i) \right]$$

(credit) This is a prior density... a number, whose expectation is taken

$$= \sum_{k=1}^K E_q \left[ \log p(\mu_k) ; m_k, s_k^2 \right] + \sum_{i=1}^N \left\{ E \left( \log p(c_i) ; \phi_i \right) + E \left[ \log p(x_i | c_i, \mu) ; \phi_i, m, s^2 \right] \right\} - \sum_{i=1}^N E \left\{ \log q(c_i | \phi_i) \right\} - \sum_{k=1}^K E \left\{ \log q(\mu_k ; m_k, s_k^2) \right\}$$

## ① Update for $c_i$ (cluster assignment)

From our earlier key fact, we have

$$q^*(c_i | \phi_i) \propto \exp \left\{ \log p(c_i) + E \left[ \log p(x_i | c_i, \mu) ; m, s^2 \right] \right\}$$

all other terms not involving  $c_i$  are constants.

Take each term:

$$\begin{aligned}\log p(c_i) &= \log \text{prior} \\ &= \log \left( \frac{1}{K} \right) = -\log K \quad \text{regardless of } c_i\end{aligned}$$

2nd term  $E[\log p(x_i | c_i, \mu) ; \mu, s^2]$

Well, we can write  $p(x_i | c_i, \mu)$  as follows:

$$p(x_i | c_i, \mu) = \prod_{k=1}^K p(x_i | \mu_k)^{c_{ik}}$$

$$\begin{aligned}\text{So } \log p(x_i | c_i, \mu) &= \sum_{k=1}^K c_{ik} \underbrace{\log p(x_i | \mu_k)}_{= -\frac{(x_i - \mu_k)^2}{2} + \text{constant}} \\ &= -\frac{(x_i - \mu_k)^2}{2} + \text{constant}\end{aligned}$$

So our second term is

$$\begin{aligned}E[\log p(x_i | c_i, \mu) ; \mu, s^2] &= E\left[ \sum_{k=1}^K c_{ik} \frac{(x_i - \mu_k)^2}{2} \right] \\ &= \sum_k c_{ik} E\left[ -\frac{(x_i - \mu_k)^2}{2} ; m_k, s_k^2 \right] + \text{constant}\end{aligned}$$

$$= \sum_k c_{ik} \left[ E(\mu_k ; m_k, s_k^2) x_i - E\left( \frac{\mu_k^2}{2} ; m_k, s_k^2 \right) \right] + \text{const}$$

Thus we need  $E(\mu_k)$  and  $E(\mu_k^2)$  under the variational Gaussian for each component:

$$E(\mu_k ; m_k, s_k^2) = m_k$$

$$E(\mu_k^2 ; m_k, s_k^2) = s_k^2 + m_k^2$$

~~$E(x^2) = \text{var}(x) + E(x)^2$~~   
 $\text{var}(x) = E(x^2) - E(x)^2$

Note:  
 $\sum_{k=1}^K c_{ik} = 1$   
 independence  
 of  $c_i$



(9)

So our optimal  $q$  looks like a multinomial:

$$q^*(c_i; \phi_i) \propto \prod_{k=1}^K [e^{\phi_{ik}}]^{C_{ik}}$$

where  ~~$\phi_i$~~   $\psi_{ik} = \lambda_i E[\mu_k; m_k, s_k^2] - \frac{E[\mu_k^2; m_k, s_k^2]}{2}$

So  $\phi_{ik} \propto e^{\psi_{ik}}$

② Update for the mixture-component means.

Recall  $q(\mu_k; m_k, s_k^2) = N(\mu_k; m_k, s_k^2)$

Again, use our basic fact:

$$q^*(\mu_k) \propto \exp \left\{ \log p(\mu_k) + \sum_{i=1}^N E \left[ \log p(x_i | c_i, \mu) ; m_{-k}, s_{-k}^2 \right] \right\}$$

$$\begin{aligned} \text{So } \log q^*(\mu_k) &= \log p(\mu_k) + \sum_{i=1}^N E \left[ \log p(x_i | c_i, \mu) ; \phi_i m_{-k}, s_{-k}^2 \right] + \text{const.} \\ &= \frac{-\mu_k^2}{2\tau^2} + \sum_i E \left[ \sum_{k=1}^K C_{ik} \log p(x_i | \mu_k) ; \phi_i m_{-k}, s_{-k}^2 \right] + \text{constant} \\ &= \frac{-\mu_k^2}{2\tau^2} + \sum_i E \left[ C_{ik} \log p(x_i | \mu_k) ; \phi_i, \mu_k, s_{-k}^2 \right] + \text{constant} \\ &= \frac{-\mu_k^2}{2\tau^2} + \sum_i \log p(x_i | \mu_k) \cdot E(C_{ik}; \phi_i) + \text{constant} \\ &= \frac{-\mu_k^2}{2\tau^2} + \sum_{i=1}^N \phi_{ik} \left( \frac{-(x_i - \mu_k)^2}{2} \right) + \text{constant} \end{aligned}$$

call other terms in  
Sum constant in  $\mu_k$

$$= -\frac{\mu_k^2}{2\tau^2} + \sum_i \left[ \phi_{ik} x_i \mu_k - \phi_{ik} \frac{\mu_k^2}{2} \right]$$

$$= \mu_k \left( \sum_{i=1}^n \phi_{ik} x_i \right) - \mu_k^2 \left[ \frac{1}{2\tau^2} + \sum_i \frac{\phi_{ik}}{2} \right]$$

This is the log density of an exponential family

with "natural sufficient statistics"  $(\mu_k, \mu_k^2)$

and "natural parameters"  $\left( \sum_i \phi_{ik} x_i, -\frac{1}{2} \left[ \frac{1}{\tau^2} + \sum_i \phi_{ik} \right] \right)$

This is a Gaussian  $(\mu_k, s_k^2)$

$$\frac{\mu_k}{s_k^2} = \sum_i \phi_{ik} x_i$$

$$-\frac{1}{2s_k^2} = -\frac{1}{2} \left[ \frac{1}{\tau^2} + \sum_i \phi_{ik} \right]$$

$$\text{So } s_k^2 = \frac{1}{\frac{1}{\tau^2} + \sum_i \phi_{ik}}$$

$$\mu_k = \frac{\sum_i \phi_{ik} x_i}{\frac{1}{\tau^2} + \sum_i \phi_{ik}}$$

How conjugate looking?

weights/precisions on each data point = variational probability  $(\phi_{ik})$  of being assigned to each cluster

Calculating the ELBO:

①

$$\begin{aligned} \text{ELBO}(q) = & \sum_{k=1}^K E_q \left[ \log p(\mu_k) ; m_k, s_k^2 \right] \\ & + \sum_{i=1}^N \left[ E \left( \log p(c_i) ; \phi_i \right) + E \left[ \log p(x_i | c_i, \mu) ; \phi_i, m, s^2 \right] \right] \\ & - \sum_{i=1}^N E \left\{ \log q(c_i ; \phi_i) \right\} - \sum_{k=1}^K E \left\{ \log q(\mu_k ; m_k, s_k^2) \right\} \end{aligned}$$

$$\begin{aligned} \textcircled{1} \quad & E_q \left[ \log p(\mu_k) ; m_k, s_k^2 \right] \\ & = E_q \left[ -\frac{\mu_k^2}{2\tau^2} ; m_k, s_k^2 \right] \\ & = -\frac{1}{2\tau^2} E_q(\mu_k^2) \end{aligned}$$

Well, under  $q$ ,  $\mu_k \sim N(m_k, s_k^2)$

$$\begin{aligned} \text{So } E_q(\mu_k^2) &= \text{var}_q(\mu_k) + E_q(\mu_k)^2 \\ &= s_k^2 + m_k^2 \end{aligned}$$

So the whole term is

$$\sum_{k=1}^K \left( -\frac{1}{2\tau^2} \right) (s_k^2 + m_k^2)$$

$$\textcircled{2} \quad E \left\{ \log p(c_i) ; \phi_i \right\} = \log \frac{1}{K} \quad \dots \text{constant}$$

$$\textcircled{7} \quad E \left\{ \log p(x_i | c_i, \mu) ; \phi_i, m, s^2 \right\} =$$

$$\sum_K \left( \underbrace{E(\mu_k ; m_k, s_k^2)}_{\mu_k} \cdot x_i - \frac{E(\mu_k^2 ; m_k, s_k^2)}{2} \right)$$

$$\begin{aligned}
 \textcircled{4} \quad & E \left\{ \log g(c_i; \phi_i) \right\} \\
 &= E \left[ \sum_{k=1}^K c_{ik} \log(\phi_{ik}) \right] \\
 &= \sum_{k=1}^K (\log \phi_{ik}) E(c_{ik}) \\
 &= \sum_{k=1}^K \phi_{ik} \log \phi_{ik}
 \end{aligned}$$

$$\begin{aligned}
 \textcircled{1} \quad & E \left\{ \log g(\mu_k; m_k, s_k^2) \right\} \\
 &= E \left\{ -\frac{1}{2s_k^2} (\mu_k - m_k)^2 \right\} \\
 &= -\frac{1}{2s_k^2} E(\mu_k^2 - 2\mu_k m_k + m_k^2) \\
 &= -\frac{1}{2s_k^2} \left[ E(\mu_k^2) - 2m_k^2 + m_k^2 \right] \\
 &= -\frac{1}{2s_k^2} \left[ s_k^2 + m_k^2 - 2m_k^2 + m_k^2 \right] = 0
 \end{aligned}$$

## Five-minute math: exponential families

Suppose that  $f(x|\theta)$  takes the form

$$f(x|\theta) = \cancel{h(x) \cdot \exp\{\eta(\theta)^T T(x) - A(\theta)\}} \\ h(x) \cdot \exp\{\eta(\theta)^T T(x) - A(\theta)\}$$

where  $\theta \in \mathbb{R}^D$

$$\eta: \mathbb{R}^D \rightarrow \mathbb{R}^D$$

$$T(x): \mathbb{R} \rightarrow \mathbb{R}^D$$

Exponential family

If  $\eta(\theta)$  is the identity, we have

$$f(x|\eta) = h(x) \cdot \exp\{\eta^T T(x) - A(\eta)\}$$

$\eta$ : "natural parameter"

$T(x)$ : "natural sufficient statistics"

$A(\eta)$  = log of normalization factor

Example:  $X \sim N(\mu, \sigma^2)$

$$\theta = (\mu, \sigma^2)$$

$$\text{Then } p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\}$$

$$= \frac{1}{\sqrt{2\pi}} \cdot \exp\left\{\eta^T \cancel{\left(\frac{x-\mu}{\sigma}\right)} - A(\eta)\right\}$$

$$\text{where } \eta = \left(\frac{\mu}{\sigma^2}, -\frac{1}{2\sigma^2}\right)^T, T(x) = (x, x^2)^T$$

## Variational inference in exponential families

(1)

Consider our generic model  $p(z, x)$ , where  $x = \text{data}$  and  $z = \text{hidden parameters}$ .

Suppose that each complete conditional is in ~~an~~ an exponential family; with sufficient statistic  $z_j$ :

$$p(z_j | z_{-j}, x) = h(z_j) \exp \left\{ \eta_j(z_{-j}, x)^T z_j - a(\eta_j(z_{-j}, x)) \right\}$$

$\uparrow$                        $\uparrow$   
"parameters"        "statistic"  
= function of  
conditioning variables

Now consider a mean-field variational approximation,

$$\text{where } q(z) = \prod_j q_j(z_j).$$

~~The exponential~~

And recall our key fact, that in coordinate ascent, the optimal  $q_j$  takes the form

$$\log q_j^*(z_j) = \mathbb{E}_{-j} [\log p(z_j | z_{-j}, x)] \quad \text{under } q_{-j}(z) = \prod_{l \neq j} q_l(z_l)$$

The exponential-family assumption makes this especially simple: (2)

$$\begin{aligned}\log q_j^*(z_j) &= E(\log p(z_j | z_{-j}, x)) \\ &= \log h(z_j) + E[\eta_j(z_{-j}, x)]^T z_j - \underbrace{E[a(\eta_j(z_{-j}, x))]}_{\text{constant in } z_j}\end{aligned}$$

$$\text{So } q_j^*(z_j) \propto h(z_j) \cdot \exp\{E[\eta_j(z_{-j}, x)] \cdot z_j\}$$

- In the same exponential family as its complete conditional  $p(z_j | z_{-j}, x)$ .
- Each update; set parameter 
$$\eta_j = E[\eta_j(z_{-j}, x)]$$

---

Any "conditionally conjugate" model in Bayes looks like this.

Suppose that :-  $\beta$  is a vector of "global" parameters

-  $z$  is a set of local parameters for each "context"

$$\text{then } p(\beta, z, x) = p(\beta) \cdot \prod_{i=1}^N p(z_i, x_i | \beta)$$

Mixture of Gaussians is a clear example.