

# Statistical Pattern Recognition Day 24 in 2021

## From EM to Variational Inference

Readings: Bishop PRML Sec 10.1  
" " Sec 10.2

### Outline

(1) Variational methods



(2) Overview of optimization methods  
for probabilistic models

(3) GMM case study:

EM vs. EE side-by-side

(4) Choosing a family of approximate posteriors  $Q$

# Variational Methods

Consider any probabilistic model with random variables

$X$  : observed data

$Z$  : hidden quantity

discrete or continuous

e.g. means in GMM  
weights in linear regression

The model defines joint  $p(X, Z)$ .

Our goal is to estimate posterior  $p(Z|X)$  } {

- evaluate PDF
- draw samples

} so we can

Usually, this is difficult due to intractable sum/integral

$$p(z|x) = \frac{p(x, z)}{p(x)}$$

*known by definition of model*      *difficult*       $p(x) = \int p(x, z) dz$

Solution: Find a good approximation among an easier to work with set of distributions  $Q$  over the domain of r.v.  $Z$ .

Define set  $Q$  of easy-to-use distributions  $q(z|w)$  with <sup>possible</sup> parameters  $w \in Q$ .  
We'll assume we can easily evaluate  $q(z|w)$  as a PDF.

We want to find parameter  $w^*$  that makes  $q(z|w^*) \approx p(z|x)$

$$w^* = \min_{w \in Q} KL(q(z|w) || p(z|x))$$

How can we solve this?

Turns out, the following are equivalent optimization problems:

$$\operatorname{arg\,min}_{w \in Q} \text{KL}(q(z|w) \parallel p(z|x))$$

$$\operatorname{arg\,max}_{w \in Q} \underbrace{\mathbb{E}_{q(z|w)} [\log p(x, z) - \log q(z|w)]}_{\text{evidence lower bound ("ELBO")}}$$

evidence lower bound ("ELBO")

$$d(x, w) \leq \log p(x)$$

with equality iff  $q(z|w) = p(z|x)$

Thus, we can solve the tractable problem

$$w^* = \operatorname{arg\,max}_{w \in Q} d(x, w)$$

using our usual optimization toolkit  
(coordinate ascent, gradient ascent)

And use  $q(z|w^*)$  as an "approximate posterior"  
however we might use  $p(z|x)$

Proof the two problems are equivalent :

$$\begin{aligned}
 \log p(x) &= \int q(z) \log p(x) dz \\
 &= \mathbb{E}_{q(z)} [\log p(x)] && \text{defn of expectation} \\
 &= \mathbb{E}_{q(z)} \left[ \log \frac{p(x, z)}{p(z|x)} \right] && \text{Bayes rule} \\
 &= \mathbb{E}_{q(z)} \left[ \log \frac{p(x, z)}{p(z|x)} \frac{q(z)}{q(z)} \right] && \text{multiply inside by } 1 = \frac{q(z)}{q(z)} \\
 &= \mathbb{E}_{q(z)} \left[ \log \frac{p(x, z)}{q(z)} \right] + \mathbb{E}_{q(z)} \left[ \log \frac{q(z)}{p(z|x)} \right] \\
 &= d(x, w) + \underbrace{KL(q(z|w) \parallel p(z|x))}_{\text{always } \geq 0 \text{ for any } z, \text{ discrete or continuous}} \\
 &\geq d(x, w)
 \end{aligned}$$

When sum of  $a(w)$  and  $b(w)$  is fixed, maximizing  $a(w)$  is equivalent to minimizing  $b(w)$  for any functions  $a$  and  $b$ .

# Overview of Optimization Strategies for Probabilistic Models

Assume a model  $p(x, z, \theta)$  with:

- $X$  : observed data
- $Z$  : hidden variable
- $\theta$  : hidden variable

how to treat  $\theta$ ?

	$\theta \in \Theta$ point estimate	$q(\theta \lambda)$ for $\lambda \in \Lambda$ approx. posterior
point estimate $z \in \Omega$	M / M $\arg \max_{z \in \Omega, \theta \in \Theta} \log p(x, z, \theta)$	M / E $\arg \max_{z \in \Omega, \lambda \in \Lambda} \mathbb{E}_{q(\theta \lambda)} \left[ \log \frac{p(x, z, \theta)}{q(\theta \lambda)} \right]$
approx posterior $q(z r)$ $r \in \mathcal{R}$	E / M $\arg \max_{r \in \mathcal{R}, \theta \in \Theta} \mathbb{E}_{q(z r)} \left[ \log \frac{p(x, z, \theta)}{q(z r)} \right]$	E / E $\arg \max_{r \in \mathcal{R}, \lambda \in \Lambda} \mathbb{E} \left[ \log \frac{p(x, z, \theta)}{q(z r)q(\theta \lambda)} \right]$

how to treat  $z$ ?

EM is one of many possible variational approaches

all □ use variational objectives based on ELBO

# EM vs EE side by side (simpler version of PRML 10.2)

GMM with fixed variance  $\sigma^2$

$$p(x, z, \pi, \mu) = \text{Dir}(\pi | a_0 \dots a_0) \prod_{k=1}^K N(\mu_k | m_0, s_0^2) \prod_{n=1}^N [\text{Cat}(z_n | \pi) N(x_n | \mu_{z_n}, \sigma^2)]$$

E/M

E/E

Update to  $z$

or  $q(z|r) = \text{Cat}(r_1 \dots r_k)$

$$\tilde{r}_{nk} = \tilde{\pi}_k N(x_n | \mu_k, \sigma^2)$$

$$r_{nk} = \frac{\tilde{r}_{nk}}{\sum_{l=1}^K \tilde{r}_{nl}}$$

$$\log \tilde{r}_{nk} = \mathbb{E}_q[\log \pi_k] + \mathbb{E}_q[\log p(x_n | \mu_k, \sigma^2)]$$

$$r_{nk} = \frac{\tilde{r}_{nk}}{\sum_{l=1}^K \tilde{r}_{nl}}$$

Update to  $\mu$

or  $q(\mu_k) = N(\mu_k | m_k, s_k^2)$

ML estimate adjust for MAP

$$\mu_k = \frac{\sum_n r_{nk} x_n}{\sum_n r_{nk}}$$

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

$$m_k^* = f(m_0, s_0, \sum_n r_{nk} x_n) \text{ see 10.61}$$

$$s_k^* = f(\underbrace{m_0, s_0}_{\text{prior}}, \underbrace{\sum_n r_{nk}}_{\text{soft. stat}}) \text{ see 10.62}$$

Update to  $\pi$

or  $q(\pi) = \text{Dir}(a_1 \dots a_k)$

ML estimate adjust for MAP  $\pi_k = \frac{N_k + a_0 - 1}{N + K a_0 - K}$

$$\pi_k = \frac{\sum_n r_{nk}}{N} = \frac{N_k}{N}$$

$$q(\pi) = \text{Dir}(a_1, \dots, a_k)$$

$$a_k^* = a_0 + \underbrace{\sum_n r_{nk}}_{\text{sufficient statistic}}$$

prior term

see 10.58

# Side-by-Side view in terms of ELBO objective

Where do updates come from?

All variational methods set up an objective function to maximize: the evidence lower bound (ELBO).

If we do coordinate ascent and try to optimally solve at each "coordinate"'s update, we get the updates on previous page.

	E/M	E/E
ELBO objective	$d(r, \pi, \mu)$ $= \mathbb{E}_{q(z r)} \left[ \log \frac{p(x, z   \pi, \mu) p(\pi, \mu)}{q(z r)} \right]$	$d(r, a, m, s)$ $= \mathbb{E}_{\substack{q(z r) \\ q(\pi a) \\ q(\mu m, s)}} \left[ \log \frac{p(x, z   \pi, \mu) p(\pi, \mu)}{q(z) q(\pi) q(\mu)} \right]$
$z/q(z)$ update	$\max_r d(r, \pi^{t-1}, \mu^{t-1})$	$\max_r d(r, a^{t-1}, m^{t-1}, s^{t-1})$
$\pi/q(\pi)$ update	$\max_{\pi} d(r^t, \pi, \mu^{t-1})$	$\max_a d(r^t, a, m^{t-1}, s^{t-1})$
$\mu/q(\mu)$ update	$\max_{\mu} d(r^t, \pi^t, \mu)$	$\max_{m, s} d(r^t, a^t, m, s)$

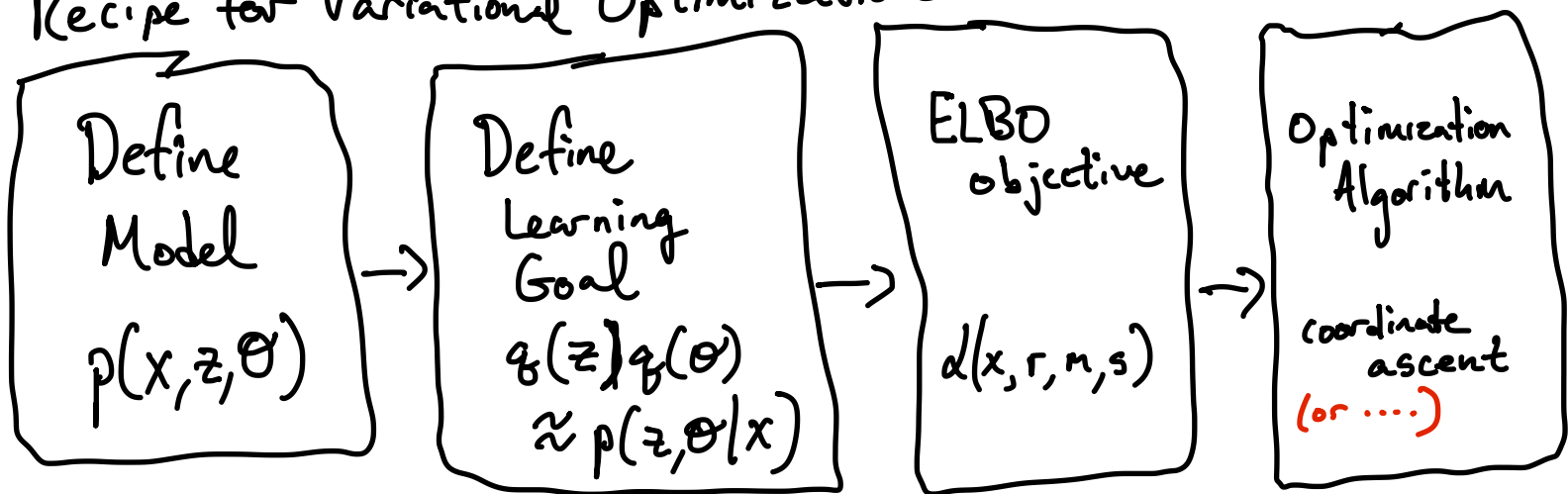
Note that we can use any algorithm we like to optimize our chosen objective function.

Given an ELBO optimization problem, the common approach (e.g. EM for GMMs) is to use coordinate ascent with closed-form updates. When possible, this is likely to be best choice (see CP4).

However, we could:

- (1) do coordinate updates that just step uphill (esp. if closed-form optima unknown)
- (2) do gradient ascent for all variables simultaneously
- (3) use any other algorithm to optimize

Recipe for Variational Optimization





How to select the "approximate posterior" set of possible distributions  $Q$ ?

Usually our goal is to balance accurate approximation with tractability.

- need  $q(z)$  easy to evaluate
- need ELBO easy to evaluate
- need ELBO easy to optimize

↳ ideally,  
 $q(z) = p(z|x)$   
"perfect" approximation

## Common strategies

(1) Impose independence assumptions on  $q(z_1, \dots, z_v)$

Full Independence:  $q(z_{1:v}) = q(z_1)q(z_2) \dots q(z_v)$

for historical reasons, often called "mean field"

Some Structured Independence:  $q(z_{1:v}) = q(z_1) \prod_{v=2}^v q(z_v | z_1)$

(2) Take advantage of conjugate model structure

e.g. in GMM, prior on  $\pi$  is  $\text{Dir}(a_0, \dots, a_0)$   
make posterior on  $\pi$  also  $\text{Dir}(a_1, \dots, a_k)$

(3) To get tractable continuous r.v. density, use a Gaussian either diagonal or low rank or full rank covariance matrix