

Variational inference

Consider a prob. dist over observed + latent variables

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

The inference problem is to compute the posterior

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

$$= \frac{p(y, z)}{p(y)}$$

$$\int_z p(x, z) dz \leftarrow \text{this is pretty hard!}$$

Variational inference :

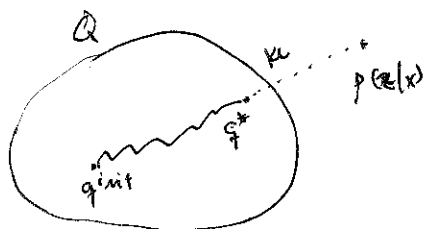
Optimize.

First: Assume there is a family of distributions Q over z .

Second: Find $q \in Q$ so that

$$q^*(z) = \underset{q \in Q}{\operatorname{argmin}} KL(q(z) \parallel p(z|x))$$

Third: Approximate the posterior with $q^*(\cdot)$.



Notation $p(x)$ is also known as the evidence. $p(x) = \int_z p(z, x) dx.$

The example we will use this entire note is that of a mixture of Gaussians is a GMM.

To elaborate:

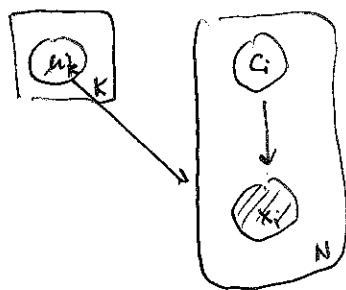
a Gaussian mixture model has K mixture components, each cluster a Gaussian w/ mean μ_k , $k=1, \dots, K$.

By assumption, each mean is sampled $\mu_k \sim p(\mu_k) = N(0, \sigma^2)$
 \uparrow
 hyperparameter.

As a generative model, a GMM is

$$\begin{aligned} \mu_k &\sim N(0, \sigma^2) \\ c_i &\sim \text{Categorical}(\frac{1}{K}, \dots, \frac{1}{K}) = \text{Uniform}(K) \\ x_i | c_i, \mu &\sim N(c_i^T \mu, 1) \end{aligned}$$

As a probabilistic graphical network



where the boxes are plates.

$$p(\mu, c, x) = p(\mu) \prod_{i=1}^n p(c_i) p(x_i | c_i, \mu)$$

Here, $z = \{\mu, c\}$ are the latent variables.

What is the ~~model~~ evidence?

$$\begin{aligned} p(x) &= \int p(\mu) \prod_{i=1}^n \sum_{c_i} p(c_i) p(x_i | c_i, \mu) d\mu \\ &= \sum_{c_i} p(c_i) \int p(\mu) \prod_{i=1}^n p(x_i | c_i, \mu) d\mu \end{aligned}$$

Oh god. ~~So much stuff~~

So $p(x)$ is hard, $p(z|x)$ is hard. let's try VI.

We want

$$q^*(z) = \underset{q \in \mathcal{Q}}{\operatorname{argmin}} \underbrace{KL(q(z) \parallel p(z|x))}_{\text{So what's this?}}$$

$$KL(q(z) \parallel p(z|x)) = \mathbb{E}_{z \sim q(z)} [\log q(z) - \log p(z|x)]$$

$$= \underbrace{\mathbb{E}_{z \sim q(z)} [\log q(z)]}_{\text{yes}} - \underbrace{\mathbb{E}_{z \sim q(z)} [\log p(x, z)]}_{\text{yes}} + \log p(x) \quad \underline{\underline{\text{NO}}}$$

What's the point if we gotta compute $p(x)$ anyway!?

KEY IDEA: But note that KL divergence is always positive!

$$0 \leq \mathbb{E}_{z \sim q(z)} [\log q(z)] - \mathbb{E}_{z \sim q(z)} [\log p(x, z)] + \log p(x)$$

So rearranging we get

$$\log p(x) \geq \underbrace{\mathbb{E}_{z \sim q(z)} [\log p(x, z)] - \mathbb{E}_{z \sim q(z)} [\log q(z)]}_{\text{ELBO}(q)}$$

"evidence lower bound"

KEY: Maximizing this = Minimizing KL

This is totally tractable. We will seek to optimize this!

Aside The ELBO can be written as

$$\begin{aligned} \text{ELBO}(q) &= \mathbb{E}_{z \sim q(z)} [\log p(x, z)] - \mathbb{E}_{z \sim q(z)} [\log q(z)] \\ &= \mathbb{E}_{z \sim q(z)} [\log p(x|z)] - \mathbb{E}_{z \sim q(z)} [\log p(z)] - \mathbb{E}_{z \sim q(z)} [\log q(z)] \end{aligned}$$

$$\boxed{\text{ELBO}(q) = \mathbb{E}_{z \sim q(z|x)} [\log p(x|z)] - \text{KL}(q(z) \parallel p(z))}$$

↑
reconstruction
error

Writing it this way expresses z (latent variables) as a "code" with encoder given by $q(z)$ and $p(x|z)$ as "decoder".

Maximizing the ELBO means we want to maximize the log likelihood of ~~the~~ getting the original representation back, while if given a good prior, penalizing representations that merely copy the data.

This gives the theoretical backing behind variational autoencoders.

Now, $\log p(x) \geq \text{ELBO}(q)$ for any $q(z)$.

ex) What is the ELBO in the GMM ~~the~~ model?

To solve this question, first we need a variational family to pull q from. Note that the complexity of the family determines the complexity of the optimization.

Usually we impose ~~these~~ stronger independence relations for distributions in Q , as these might be more computable.

A common family to use is the mean-field variational family.

$$q(z) = \prod_{j=1}^m q_j(z_j), \quad m = \# \text{ of latent variables.}$$

In the mean-field variational family, the latent variables are decoupled.

What are the q_j then? Depends on the latent variable.

If z_j is ~~discrete~~ $\left\{ \begin{array}{l} \text{discrete} \rightarrow \text{categorical, multinomial, ...} \\ \text{continuous} \rightarrow \text{Gaussian, ...} \end{array} \right.$

Note If we add ~~the~~ dependencies between the variables, we get families for structured variational inference.

In the GMM model, the mean-field VF comes as the form

$$q(\mu, c) = \underbrace{\prod_{k=1}^K q(\mu_k; m_k, s_k^2)}_{\text{Gaussians}} \underbrace{\prod_{i=1}^n q(c_i; \varphi_i)}_{\text{distribution of } x_i \text{'s mixture assignment.}}$$

$\mu_k \sim N(\mu_k; m_k, s_k^2)$

$\varphi_i = K\text{-vector of assignment probabilities.}$

The total parameter space is given by

$$\{m_k, s_k^2 : k=1, \dots, K, \varphi_i \in \mathbb{R}^K; i=1, \dots, N\}$$

$$q_{\{m_k, s_k^2, \varphi_i\}}(\mu, c) \in \mathcal{Q}_{\text{mean-field}}.$$

The ELBO here is thus given by

$$\begin{aligned} \text{ELBO}(q_{\{m_k, s_k^2, \varphi_i\}}) &= \sum_{k=1}^K \mathbb{E}_{z \sim q(z)} [\log p(\mu_k); m_k, s_k^2] \\ &+ \sum_{i=1}^n \left(\mathbb{E} [\log p(c_i); \varphi_i] + \mathbb{E} [\log p(x_i | c_i, \mu); \varphi_i, m, s^2] \right) \\ &- \sum_{i=1}^n \mathbb{E} [\log q(c_i; \varphi_i)] - \sum_{k=1}^K \mathbb{E} [\log q(\mu_k; m_k, s_k^2)]. \end{aligned}$$

Now to optimize. We use a strategy called coordinate-ascent mean-field variational inference (CAVI).

CAVI iteratively optimizes each factor of the mean-field variational density, while holding the others fixed.

How to do it : let z_j be a latent variable.

The complete conditional of z_j is $p(z_j | z_{\neq j}, x)$

Fix the other variational factors $q_l(z_l)$, $l \neq j$.

We want the optimal $q_j(z_j)$, called $q_j^*(z_j)$ to maximize ELBO.

Claim : $q_j^*(z_j) \propto \exp(\mathbb{E}_{z_{\neq j} \sim q(z)} [\log p(z_j | z_{\neq j}, x)])$

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

$$= \mathbb{E}_{z_j \sim q_j(z_j)} [\mathbb{E}_{z_{\neq j} \sim q(z)} [\log p(z_j, z_{\neq j}, x)]] - \mathbb{E}_{z_j \sim q_j(z_j)} [\log q_j(z_j)] + \text{const.}$$

$$= \text{const} - \text{KL}(q_j(z_j) \parallel \exp(\mathbb{E}_{z_{\neq j} \sim q(z)} [\log p(z_j, z_{\neq j}, x)]))$$

↑
minimize this

so take $q_j^*(z_j) \propto \exp(\mathbb{E}_{z_{\neq j} \sim q(z)} [\log p(z_j | z_{\neq j}, x)])$. □.

ex) What are the optimal updates then for the GMM?

Ans : $q^*(c_i; \psi_i) \propto \exp(\underbrace{\log p(c_i)}_{\log \frac{1}{K}} + \mathbb{E}_{p(x_i, z_i)} [\log p(x_i | c_i, \mu)] ; m, s^2)$

$$p(x_i | c_i, \mu) = \prod_{k=1}^K p(x_i | \mu_k)^{c_{ik}} \quad (\text{as } c_i \text{ is a one-hot vector})$$

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

ie the variational update for the i th cluster assignment is

$$\varphi_{ik} \propto \exp \left(\mathbb{E} [\mu_k ; m_k, s_k^2] x_i - \mathbb{E} [\mu_k^2 ; m_k, s_k^2] / 2 \right)$$

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

Now, as $\varphi_{ik} = \mathbb{E}_{\varphi_i} [c_{ik}]$ (c_{ik} is an indicator random variable),

$$\begin{aligned}
 \log q^*(\mu_k ; m_k, s_k^2) &= \log p(\mu_k) + \sum_{i=1}^n \mathbb{E} [\log p(x_i | c_i, \mu) ; \varphi_i, m_{\neq k}, s_{\neq k}^2] + \text{const} \\
 &= -\mu_k^2 / 2\sigma^2 + \sum_{i=1}^n \mathbb{E} [c_{ik}] \log p(x_i | \mu_k) + \text{const} \\
 &= -\mu_k^2 / 2\sigma^2 + \sum_{i=1}^n \varphi_{ik} \cdot \left(- (x_i - \mu_k)^2 / 2 \right) + \text{const} \\
 &= \left(\sum_{i=1}^n \varphi_{ik} x_i \right) \mu_k - \left(\frac{1}{2\sigma^2} + \sum_{i=1}^n \varphi_{ik} / 2 \right) \mu_k^2 + \text{const.}
 \end{aligned}$$

ie $q^*(\mu_k ; m_k, s_k^2)$ is a Gaussian sufficient statistics.

Normalizing, we get the updates are

$$m_k = \frac{\sum_{i=1}^n \varphi_{ik} x_i}{\frac{1}{\sigma^2} + \sum_{i=1}^n \varphi_{ik}}, \quad s_k^2 = \frac{1}{\frac{1}{\sigma^2} + \sum_{i=1}^n \varphi_{ik}}$$

$$\log p(\mu, c, x) = \sum_{k=1}^K \log p(\mu_k) + \sum_{i=1}^n \underbrace{(\log p(c_i) + \log p(x_i | c_i, \mu))}_{-\log K}$$

$$\propto \sum_{k=1}^K \left(-\frac{\mu_k^2}{2\sigma^2} \right) + \sum_{i=1}^n \sum_{k=1}^K c_{ik} \left(-\frac{(x_i - \mu_k)^2}{2} \right)$$

$\mu_k \sim N(0, \sigma^2)$
 $x_i | c_i, \mu \sim N(c_i^T \mu, 1)$

$$\log p(\mu, c, x) \propto \sum_{k=1}^K \left(-\frac{\mu_k^2}{2\sigma^2} \right) + \sum_{i=1}^n \sum_{k=1}^K c_{ik} \left(-\frac{(x_i - \mu_k)^2}{2} \right).$$

$$\mathbb{E}_{z \sim q(z)} [\log p(\mu, c, x)] = \sum_{k=1}^K -\frac{1}{2\sigma^2} \mathbb{E}_q[\mu_k^2] + \sum_{i=1}^n \sum_{k=1}^K \mathbb{E}[c_{ik}] \left(-\frac{(x_i - \mu_k)^2}{2} \right)$$

$$= \sum_{k=1}^K -\frac{(m_k^2 + s_k^2)}{2\sigma^2} - \sum_{i=1}^n \sum_{k=1}^K \cancel{\varphi_{ik}} \mathbb{E}_q \left[\frac{(x_i - \mu_k)^2}{2} \right]$$

$$\mathbb{E}_{z \sim q(z)} [\log q(\mu, s^2, \varphi)] \propto \sum_{i=1}^n \sum_{k=1}^K \mathbb{E} \log \varphi_{ik} + \sum_{k=1}^K \mathbb{E} \left(-\frac{1}{2} \log s_k^2 - \frac{(\mu_k - m_k)^2}{2s_k^2} \right)$$

$$= \sum_{i=1}^n \sum_{k=1}^K \log \varphi_{ik} - \sum_{k=1}^K \left(\frac{1}{2} \log s_k^2 + \frac{1}{2} \right)$$

constant.

$$\propto \sum_{i=1}^n \sum_{k=1}^K \log \varphi_{ik} - \frac{1}{2} \sum_{k=1}^K \log s_k^2.$$

$$\bullet \mathbb{E}_q[(x_i - \mu_k)^2] = x_i^2 - 2x_i m_k + (m_k^2 + s_k^2)$$

$$\bullet \log p(x_i | c_i, \mu) = \log N(x_i | c_i^T \mu, 1) = -\frac{1}{2} \left(x_i - \sum_{k=1}^K c_{ik} \mu_k \right)^2$$

~~Not~~