

DSC190: Machine Learning with Few Labels

Supervised Learning, Unsupervised Learning

Zhiting Hu

Lecture 3, September 30, 2021

UC San Diego

HALICIOĞLU DATA SCIENCE INSTITUTE

Outline

- Recap: EM algorithm
- Variational Inference

KL Divergence

- Kullback-Leibler (KL) divergence: measures the closeness of two distributions $p(\mathbf{x})$ and $q(\mathbf{x})$

$$\text{KL}(q(\mathbf{x}) \parallel p(\mathbf{x})) = \sum_{\mathbf{x}} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})}$$

- a.k.a. Relative entropy
- $\text{KL} \geq 0$ (Jensen's inequality)
- Intuitively:
 - If q is high and p is high, then we are happy (i.e. low KL divergence)
 - If q is high and p is low then we pay a price (i.e. high KL divergence).
 - If q is low then we don't care (i.e. also low KL divergence, regardless of p)
- not a true "distance":
 - not commutative (symmetric) $\text{KL}(p \parallel q) \neq \text{KL}(q \parallel p)$
 - doesn't satisfy triangle inequality

KL Divergence

- Kullback-Leibler (KL) divergence: measures the closeness of two distributions $p(\mathbf{x})$ and $q(\mathbf{x})$

$$\text{KL}(q(\mathbf{x}) \parallel p(\mathbf{x})) = \sum_{\mathbf{x}} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})}$$

- a.k.a. Relative entropy
- MLE is minimizing the KL divergence between the empirical distribution and the model distribution

$$\text{KL}(\tilde{p}(\mathbf{x}) \parallel p_{\theta}(\mathbf{x})) = -\mathbb{E}_{\tilde{p}(\mathbf{x})}[\log p_{\theta}(\mathbf{x})] + H(\tilde{p}(\mathbf{x}))$$



Cross entropy

EM Algorithm

- Observed variables \mathbf{x} , latent variables \mathbf{z}
- To learn a model $p(\mathbf{x}, \mathbf{z}|\theta)$, we want to maximize the marginal log-likelihood

$$\ell(\theta; \mathbf{x}) = \log p(\mathbf{x}|\theta) = \log \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}|\theta)$$

- But it's too difficult
- EM algorithm:
 - maximize a lower bound of $\ell(\theta; \mathbf{x})$
 - Or equivalently, minimize an upper bound of $\ell(\theta; \mathbf{x})$
- Key equation:

$$\ell(\theta; \mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right] + \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta))$$

(Whiteboard)

- Marginal log-likelihood $\ell(\theta; \mathbf{x}) = \log p(\mathbf{x}|\theta) = \log \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}|\theta)$

$$\ell(\theta; \mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right] + \text{KL}(q(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x}, \theta))$$

EM Algorithm

- Observed variables \mathbf{x} , latent variables \mathbf{z}
- To learn a model $p(\mathbf{x}, \mathbf{z}|\theta)$, we want to maximize the marginal log-likelihood

$$\ell(\theta; \mathbf{x}) = \log p(\mathbf{x}|\theta) = \log \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}|\theta)$$

- But it's too difficult
- EM algorithm:
 - maximize a lower bound of $\ell(\theta; \mathbf{x})$
 - Or equivalently, minimize an upper bound of $\ell(\theta; \mathbf{x})$
- Key equation:

$$\ell(\theta; \mathbf{x}) = \underbrace{\mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right]}_{\text{Evidence Lower Bound (ELBO)}} + \text{KL}(q(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x}, \theta))$$

$$= \underbrace{-F(q, \theta)}_{\text{Variational free energy}} + \text{KL}(q(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x}, \theta))$$

Variational free energy

EM Algorithm

- The EM algorithm is coordinate-decent on $F(q, \theta)$
 - E-step: $q^{t+1} = \arg \min_q F(q, \theta^t) = p(\mathbf{z}|\mathbf{x}, \theta^t)$
 - the posterior distribution over the latent variables given the data and the current parameters
 - M-step: $\theta^{t+1} = \arg \min_{\theta} F(q^{t+1}, \theta) = \operatorname{argmax}_{\theta} \sum_{\mathbf{z}} q^{t+1}(\mathbf{z}|\mathbf{x}) \log p(\mathbf{x}, \mathbf{z}|\theta)$

$$\begin{aligned} \ell(\theta; \mathbf{x}) &= \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right] + \text{KL}(q(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x}, \theta)) \\ &= -F(q, \theta) + \text{KL}(q(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x}, \theta)) \end{aligned}$$

Example: Gaussian Mixture Models (GMMs)

- Consider a mixture of K Gaussian components:

- Z is a latent class indicator vector:

$$p(z_n) = \text{multi}(z_n : \pi) = \prod_k (\pi_k)^{z_n^k}$$

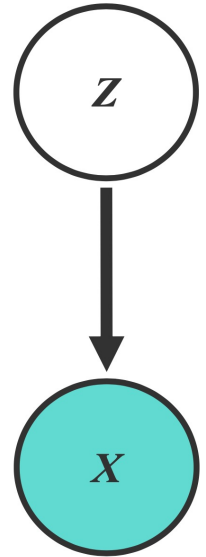
- X is a conditional Gaussian variable with a class-specific mean/covariance

$$p(x_n | z_n^k = \mathbf{1}, \mu, \Sigma) = \frac{1}{(2\pi)^{m/2} |\Sigma_k|^{1/2}} \exp\left\{-\frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k)\right\}$$

- The likelihood of a sample:

$$\begin{aligned} p(x_n | \mu, \Sigma) &= \sum_k p(z^k = \mathbf{1} | \pi) p(x, | z^k = \mathbf{1}, \mu, \Sigma) \\ &= \sum_{z_n} \prod_k \left((\pi_k)^{z_n^k} N(x_n : \mu_k, \Sigma_k)^{z_n^k} \right) = \sum_k \pi_k N(x, | \mu_k, \Sigma_k) \end{aligned}$$

mixture proportion mixture component



GMM E-step:

- Compute the posterior of z given x , using Bayes rule

GMM E-step:

- Compute the posterior of z given \mathbf{x} , using Bayes rule

$$\begin{aligned} p(z^k = 1 \mid \mathbf{x}) &= \frac{p(z^k = 1)p(\mathbf{x} \mid z^k = 1)}{p(\mathbf{x})} \\ &= \frac{p(z^k = 1)p(\mathbf{x} \mid z^k = 1)}{\sum_{j=1}^K p(z^j = 1)p(\mathbf{x} \mid z^j = 1)} \\ &= \frac{\pi_k \mathcal{N}(\mathbf{x} \mid \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x} \mid \mu_j, \Sigma_j)} \\ &:= \gamma_k \end{aligned}$$

GMM M-step:

- Once we have $q^{t+1}(z^k|x) = p(z^k|x, \theta^t) = \gamma^k$, we can compute the expected likelihood:

$$\theta^{t+1} = \operatorname{argmax}_{\theta} \sum_k q^{t+1}(z^k = 1|x) \log p(x, z^k = 1|\theta)$$

$$\begin{aligned} & \mathbb{E}_{q^{t+1}} [\log (p(\mathbf{x}, z | \boldsymbol{\theta}))] \\ &= \sum_k \gamma_k (\log p(z^k = 1|\boldsymbol{\theta}) + \log P(\mathbf{x} | z^k = 1, \boldsymbol{\theta})) \\ &= \sum_k \gamma_k \log \pi_k + \sum_k \gamma_k \log \mathcal{N}(\mathbf{x}; \mu_k, \Sigma_k) \end{aligned}$$

- We need to fit K Gaussians, just need to weight examples by γ_k

EM Algorithm for GMM

- Initialize the means μ_k , covariances Σ_k and mixing coefficients π_k
- Iterate until convergence:
 - E-step: Evaluate the posterior given current parameters

$$p(z^k = 1 | \mathbf{x}) = \frac{\pi_k \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x} | \mu_j, \Sigma_j)} := \gamma_k$$

- M-step: Re-estimate the parameters given current posterior

$$\begin{aligned} & \mathbb{E}_{q^{t+1}} [\log (p(\mathbf{x}, z | \boldsymbol{\theta}))] \\ &= \sum_k \gamma_k (\log p(z^k = 1 | \boldsymbol{\theta}) + \log P(\mathbf{x} | z^k = 1, \boldsymbol{\theta})) \\ &= \sum_k \gamma_k \log \pi_k + \sum_k \gamma_k \log \mathcal{N}(\mathbf{x}; \mu_k, \Sigma_k) \end{aligned}$$

Key Takeaways of EM Algorithm

- The EM algorithm is coordinate-decent on $F(q, \theta)$
 - E-step: $q^{t+1} = \arg \min_q F(q, \theta^t) = p(\mathbf{z}|\mathbf{x}, \theta^t)$
 - M-step: $\theta^{t+1} = \arg \min_{\theta} F(q^{t+1}, \theta) = \operatorname{argmax}_{\theta} \sum_{\mathbf{z}} q^{t+1}(\mathbf{z}|\mathbf{x}) \log p(\mathbf{x}, \mathbf{z}|\theta)$

$$\begin{aligned} \ell(\theta; \mathbf{x}) &= \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right] + \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta)) \\ &= -F(q, \theta) + \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta)) \end{aligned}$$

- Limitation: need to be able to compute $p(\mathbf{z}|\mathbf{x}, \theta)$, not possible for more complicated models --- solution: Variational inference

Variational Inference

Inference

- Given a model, the goals of inference can include:
 - Computing the likelihood of observed data $p(\mathbf{x}^*)$
 - Computing the marginal distribution over a given subset of variables in the model $p(\mathbf{x}_A)$
 - Computing the conditional distribution over a subsets of nodes given a disjoint subset of nodes $p(\mathbf{x}_A|\mathbf{x}_B)$
 - Computing a mode of the density (for the above distributions) $\operatorname{argmax}_x p(\mathbf{x})$
 -

Variational Inference

- Observed variables \mathbf{x} , latent variables \mathbf{z}
- Variational (Bayesian) inference, a.k.a. **variational Bayes**, is most often used to **approximately** infer the conditional distribution over the latent variables given the observations (and parameters)
 - i.e., the **posterior distribution** over the latent variables

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

Motivating Example

- Why do we often need to use an approximate inference methods (such as variational Bayes) to compute the posterior distribution?
- It's because we cannot directly compute the posterior distribution for many interesting models
 - I.e. the posterior density is in an intractable form (often involving integrals) which cannot be easily analytically solved.
- As a motivating example, we will try to compute the posterior for a (Bayesian) mixture of Gaussians.

Bayesian mixture of Gaussians

- The mean μ_k is treated as a (latent) random variable

$$\mu_k \sim \mathcal{N}(0, \tau^2) \text{ for } k = 1, \dots, K$$

- For each data $i = 1, \dots, n$

$$z_i \sim \text{Cat}(\pi).$$

$$x_i \sim \mathcal{N}(\mu_{z_i}, \sigma^2).$$

- We have
 - observed variables $x_{1:n}$
 - latent variables $\mu_{1:k}$ and $z_{1:n}$
 - parameters $\{\tau^2, \pi, \sigma^2\}$

- $p(x_{1:n}, z_{1:n}, \mu_{1:k} | \tau^2, \pi, \sigma^2) = \prod_{k=1}^K p(\mu_k) \prod_{i=1}^n p(z_i) p(x_i | z_i, \mu_{1:K})$

Bayesian mixture of Gaussians

- We can write the posterior distribution as

$$p(\mu_{1:K}, z_{1:n} | x_{1:n}) = \frac{\prod_{k=1}^K p(\mu_k) \prod_{i=1}^n p(z_i) p(x_i | z_i, \mu_{1:K})}{\int_{\mu_{1:K}} \sum_{z_{1:n}} \prod_{k=1}^K p(\mu_k) \prod_{i=1}^n p(z_i) p(x_i | z_i, \mu_{1:K})}$$

- The numerator can be computed for any choice of the latent variables
- The problem is the denominator (the marginal probability of the observations)
 - This integral cannot easily be computed analytically
- We need some approximation..

Variational Inference

The main idea behind variational inference:

- Choose a family of distributions over the latent variables $z_{1:m}$ with its own set of variational parameters ν , i.e.

$$q(z_{1:m}|\nu)$$

- Then, we find the setting of the parameters that makes our approximation q closest to the posterior distribution.
 - This is where optimization algorithms come in.
- Then we can use q with the fitted parameters in place of the posterior.
 - E.g. to form predictions about future data, or to investigate the posterior distribution over the hidden variables, find modes, etc.

Variational Inference

- We want to minimize the KL divergence between our approximation $q(\mathbf{z}|\mathbf{x})$ and our posterior $p(\mathbf{z}|\mathbf{x})$

$$\text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}))$$

- But we can't actually minimize this quantity w.r.t q because $p(\mathbf{z}|\mathbf{x})$ is unknown

Evidence Lower Bound (ELBO)

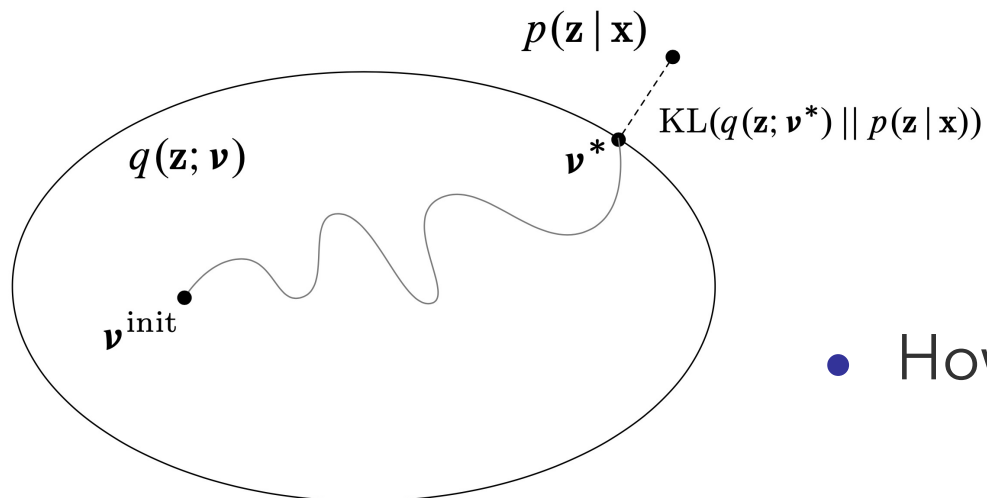
$$\ell(\theta; \mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right] + \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta))$$

- The ELBO is equal to the negative KL divergence up to a constant $\ell(\theta; \mathbf{x})$
- We maximize the ELBO over q to find an "optimal approximation" to $p(\mathbf{z}|\mathbf{x})$

Variational Inference

- Choose a family of distributions over the latent variables \mathbf{z} with its own set of variational parameters ν , i.e. $q(\mathbf{z}|\mathbf{x}, \nu)$
- We maximize the ELBO over q to find an “optimal approximation” to $p(\mathbf{z}|\mathbf{x})$

$$\begin{aligned} & \operatorname{argmax}_{\nu} \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \nu)} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x}, \nu)} \right] \\ & = \operatorname{argmax}_{\nu} \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \nu)} [\log p(\mathbf{x}, \mathbf{z}|\theta)] - \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \nu)} [\log q(\mathbf{z}|\mathbf{x}, \nu)] \end{aligned}$$



- How do we choose the variational family $q(\mathbf{z}|\mathbf{x}, \nu)$?

Mean Field Variational Inference

- A popular family of variational approximations
- In this type of variational inference, we assume the variational distribution over the latent variables **factorizes** as

$$q(\mathbf{z}) = q(z_1, \dots, z_m) = \prod_{j=1}^m q(z_j)$$

- (where we omit variational parameters for ease of notation)
 - We refer to $q(z_j)$, the variational approximation for a single latent variable, as a “local variational approximation”
- In the above expression, the variational approximation $q(z_j)$ over each latent variable z_j is independent

Mean Field Variational Inference

- Note that this is a fairly general setup; we can also partition the latent variables z_1, \dots, z_m into R groups z_{G_1}, \dots, z_{G_R} , and use the approximation:

$$q(z_1, \dots, z_m) = q(z_{G_1}, \dots, z_{G_R}) = \prod_{r=1}^R q(z_{G_r})$$

- Often called “generalized mean field” versus (the above) “naïve mean field”.

Mean Field Variational Inference

- Note that this is a fairly general setup; we can also partition the latent variables z_1, \dots, z_m into R groups z_{G_1}, \dots, z_{G_R} , and use the approximation:

$$q(z_1, \dots, z_m) = q(z_{G_1}, \dots, z_{G_R}) = \prod_{r=1}^R q(z_{G_r})$$

- Often called “generalized mean field” versus (the above) “naïve mean field”.
- Typically, this approximation does not contain the true posterior (because the latent variables are dependent).
 - E.g.: in the (Bayesian) mixture of Gaussians model, all of the cluster assignments z_i for $i = 1, \dots, n$ are dependent on each other and on the cluster locations $\mu_{1:K}$ given data.

Optimizing the ELBO in Mean Field Variational Inference

How do we optimize the ELBO in mean field variational inference?

- Typically, we use coordinate ascent optimization.
- I.e. we optimize each latent variable's variational approximation $q(z_j)$ in turn while holding the others fixed.
 - At each iteration we get an updated “local” variational approximation.
 - And we iterate through each latent variable until convergence.

Optimizing the ELBO in Mean Field Variational Inference

- Recall that the ELBO is defined as:

$$\mathcal{L} = \mathbb{E}_q[\log p(\mathbf{x}, \mathbf{z})] - \mathbb{E}_q[\log q(\mathbf{z})]$$

- Note that we can decompose the entropy term of the ELBO (using the mean field variational approximation) as:

$$\mathbb{E}_q[\log q(z_{1:m})] = \sum_{j=1}^m \mathbb{E}_{q_j}[\log q(z_j)]$$

- Therefore, under the mean field approximation, the ELBO can be written:

$$\mathcal{L} = \mathbb{E}_{q_j} \mathbb{E}_{q_{-j}}[\log p(\mathbf{x}, \mathbf{z})] - \sum_{j=1}^m \mathbb{E}_{q_j}[\log q(z_j)]$$

Optimizing the ELBO in Mean Field Variational Inference

- Therefore, under the mean field approximation, the ELBO can be written:

$$\mathcal{L} = \mathbb{E}_{q_j} \mathbb{E}_{q_{-j}} [\log p(\mathbf{x}, \mathbf{z})] - \sum_{j=1}^m \mathbb{E}_{q_j} [\log q(z_j)]$$

- Next, we want to derive the coordinate ascent update for a latent variable z_j , keeping all other latent variables fixed.
 - i.e. we want the $\operatorname{argmax}_{q_j} \mathcal{L}$.
- Removing the parts that do not depend on $q(z_j)$, we can write:

$$\mathcal{L} = \mathbb{E}_{q_j} \mathbb{E}_{q_{-j}} [\log p(\mathbf{x}, \mathbf{z})] - \mathbb{E}_{q_j} [\log q(z_j)] + \text{const.}$$

Optimizing the ELBO in Mean Field Variational Inference

- To find this argmax, we take the derivative of \mathcal{L} w.r.t $q(z_j)$ and set the derivative to zero :

$$\frac{d\mathcal{L}}{dq(z_j)} = \mathbb{E}_{q_j} \mathbb{E}_{q_{-j}} [\log p(\mathbf{x}, \mathbf{z})] - \log q(z_j) - 1 = 0$$

- From this, we arrive at the coordinate ascent update:

$$q^*(z_j) \propto \exp \left\{ \mathbb{E}_{q_{-j}} [\log p(\mathbf{x}, \mathbf{z})] \right\}$$

Optimizing the ELBO in Mean Field Variational Inference

- The coordinate ascent update:

$$q^*(z_j) \propto \exp \left\{ \mathbb{E}_{q_{-j}} [\log p(\mathbf{x}, \mathbf{z})] \right\}$$

- The optimal solution for factor $q(z_j)$ is obtained simply by considering the log of the joint distribution over all observed and latent variables and then taking the expectation with respect to all of the other factors $q(z_k)$, $k \neq j$, then taking exponential and normalizing
- Note that the only assumption we made so far is the mean-field factorization:
$$q(\mathbf{z}) = q(z_1, \dots, z_m) = \prod_{j=1}^m q(z_j)$$

- We haven't yet made any assumptions on the form of $q(z_j)$

Simple example:

- Consider a univariate Gaussian distribution $p(x) = \mathcal{N}(x|\mu, \tau^{-2})$, given a dataset $\mathcal{D} = \{x_1, \dots, x_N\}$:

$$p(\mathcal{D}|\mu, \tau) = \left(\frac{\tau}{2\pi}\right)^{N/2} \exp\left\{-\frac{\tau}{2} \sum_{n=1}^N (x_n - \mu)^2\right\}$$

$$p(\mu|\tau) = \mathcal{N}(\mu|\mu_0, (\lambda_0\tau)^{-1})$$

$$p(\tau) = \text{Gam}(\tau|a_0, b_0)$$

- $\text{Gam}(\tau|a_0, b_0) = \frac{1}{\Gamma(a)} b^a \lambda^{a-1} \exp(-b\lambda)$: gamma distribution
- For this simple problem the posterior distribution can be found exactly. But we use it as an example for tutorial anyway

$$q^*(z_j) \propto \exp \left\{ \mathbb{E}_{q_{-j}} [\log p(\mathbf{x}, \mathbf{z})] \right\}$$

Simple example:

$$p(\mathcal{D}|\mu, \tau) = \left(\frac{\tau}{2\pi} \right)^{N/2} \exp \left\{ -\frac{\tau}{2} \sum_{n=1}^N (x_n - \mu)^2 \right\}$$

$$p(\mu|\tau) = \mathcal{N}(\mu|\mu_0, (\lambda_0\tau)^{-1})$$

$$p(\tau) = \text{Gam}(\tau|a_0, b_0)$$

- Introduce the factorized variational approximation: $q(\mu, \tau) = q_\mu(\mu)q_\tau(\tau)$
- Solution to q_μ :

$$\begin{aligned} \ln q_\mu^*(\mu) &= \mathbb{E}_\tau [\ln p(\mathcal{D}|\mu, \tau) + \ln p(\mu|\tau)] + \text{const} \\ &= -\frac{\mathbb{E}[\tau]}{2} \left\{ \lambda_0(\mu - \mu_0)^2 + \sum_{n=1}^N (x_n - \mu)^2 \right\} + \text{const.} \end{aligned}$$

- We can see q_μ^* is a Gaussian $\mathcal{N}(x|\mu_N, \lambda_N^{-1})$:

$$\mu_N = \frac{\lambda_0\mu_0 + N\bar{x}}{\lambda_0 + N}$$

$$\lambda_N = (\lambda_0 + N)\mathbb{E}[\tau]$$

$$q^*(z_j) \propto \exp \left\{ \mathbb{E}_{q_{-j}} [\log p(\mathbf{x}, \mathbf{z})] \right\}$$

Simple example:

$$p(\mathcal{D}|\mu, \tau) = \left(\frac{\tau}{2\pi} \right)^{N/2} \exp \left\{ -\frac{\tau}{2} \sum_{n=1}^N (x_n - \mu)^2 \right\} \quad \begin{array}{l} p(\mu|\tau) = \mathcal{N}(\mu|\mu_0, (\lambda_0\tau)^{-1}) \\ p(\tau) = \text{Gam}(\tau|a_0, b_0) \end{array}$$

- Introduce the factorized variational approximation: $q(\mu, \tau) = q_\mu(\mu)q_\tau(\tau)$

- Solution to q_τ :
$$\begin{aligned} \ln q_\tau^*(\tau) &= \mathbb{E}_\mu [\ln p(\mathcal{D}|\mu, \tau) + \ln p(\mu|\tau)] + \ln p(\tau) + \text{const} \\ &= (a_0 - 1) \ln \tau - b_0\tau + \frac{N}{2} \ln \tau \\ &\quad - \frac{\tau}{2} \mathbb{E}_\mu \left[\sum_{n=1}^N (x_n - \mu)^2 + \lambda_0(\mu - \mu_0)^2 \right] + \text{const} \end{aligned}$$

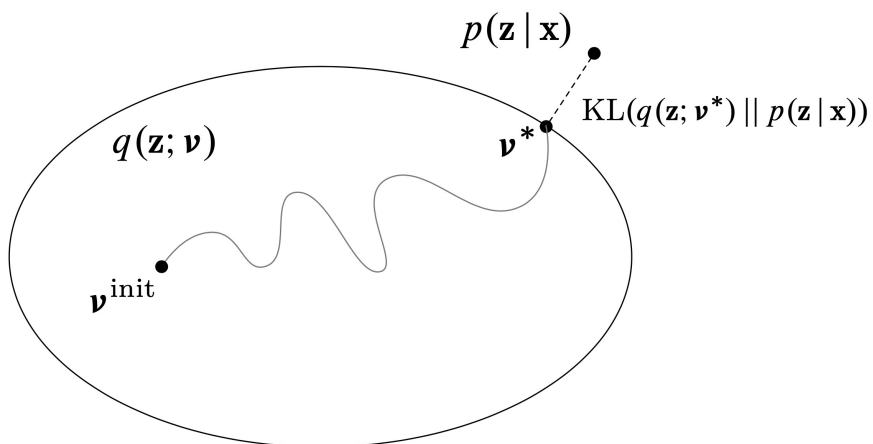
- We can see q_τ^* is a gamma distribution $\text{Gam}(\tau|a_N, b_N)$:

$$a_N = a_0 + \frac{N}{2}$$

$$b_N = b_0 + \frac{1}{2} \mathbb{E}_\mu \left[\sum_{n=1}^N (x_n - \mu)^2 + \lambda_0(\mu - \mu_0)^2 \right]$$

Quick Recap

- We often cannot compute posteriors, and so we need to approximate them, using variational methods.
- In variational Bayes, we'd like to find an approximation within some family that minimizes the KL divergence to the posterior, but we can't directly minimize this
- Therefore, we defined the ELBO, which we can maximize, and this is equivalent to minimizing the KL divergence.



Evidence Lower Bound (ELBO)

$$\ell(\theta; \mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right] + \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta))$$

Quick Recap

- We defined a family of approximations called “mean field” approximations, in which there are no dependencies between latent variables

$$q(\mathbf{z}) = q(z_1, \dots, z_m) = \prod_{j=1}^m q(z_j)$$

- We optimize the ELBO with coordinate ascent updates to iteratively optimize each local variational approximation under mean field assumptions

$$q^*(z_j) \propto \exp \left\{ \mathbb{E}_{q_{-j}} [\log p(\mathbf{x}, \mathbf{z})] \right\}$$

Key Takeaways

- KL Divergence $\text{KL}(q(\mathbf{x}) \parallel p(\mathbf{x})) = \sum_{\mathbf{x}} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})}$

- The key equation of EM and VI

Evidence Lower Bound (ELBO)

$$\ell(\theta; \mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right] + \text{KL}(q(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x}, \theta))$$

- Free energy $F(q, \theta)$
- EM: E-step and M-step optimizing ELBO w.r.t q and θ
- Mean-field VI: optimizing ELBO w.r.t factorized q components

Questions?