

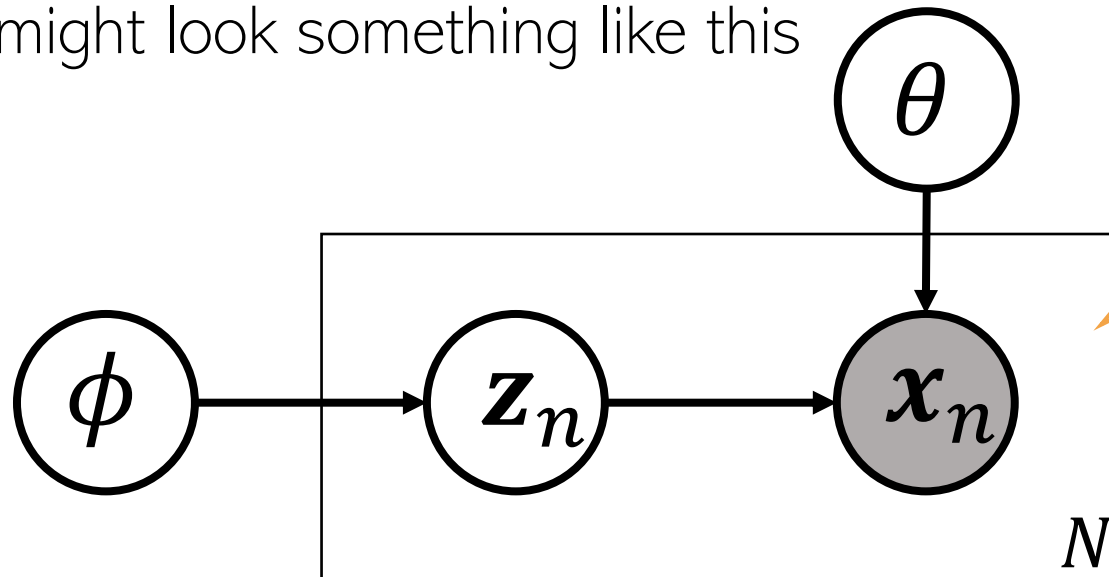
Variational Inference

CS772A: Probabilistic Machine Learning

Piyush Rai

Variational Inference (VI)

- Assume a latent variable model with data \mathcal{D} and latent variables \mathbf{Z}
- A simple setting might look something like this



This setting is just one example. VI is applicable in more general and more complex probabilistic models with and without latent variables

- Assume the likelihood is $p(\mathcal{D}|\mathbf{Z}, \Theta)$ and prior is $p(\mathbf{Z}|\Theta)$. **Want posterior over \mathbf{Z}**
- $\Theta = (\theta, \phi)$ denotes the other parameters that define the likelihood and the prior
- For now, assume Θ is known and only \mathbf{Z} is unknown (the Θ unknown case later)
- Assume CP $p(\mathbf{Z}|\mathcal{D}, \Theta)$ is intractable



Variational Inference (VI)

- Assuming $p(\mathbf{Z}|\mathcal{D}, \Theta)$ is intractable, VI approximates it by a distr $q(\mathbf{Z}|\phi)$ or $q_\phi(\mathbf{Z})$

Find the optimal ϕ which makes our approximation $q(\mathbf{Z}|\phi)$ as closed as possible to the true posterior $p(\mathbf{Z}|\mathcal{D})$

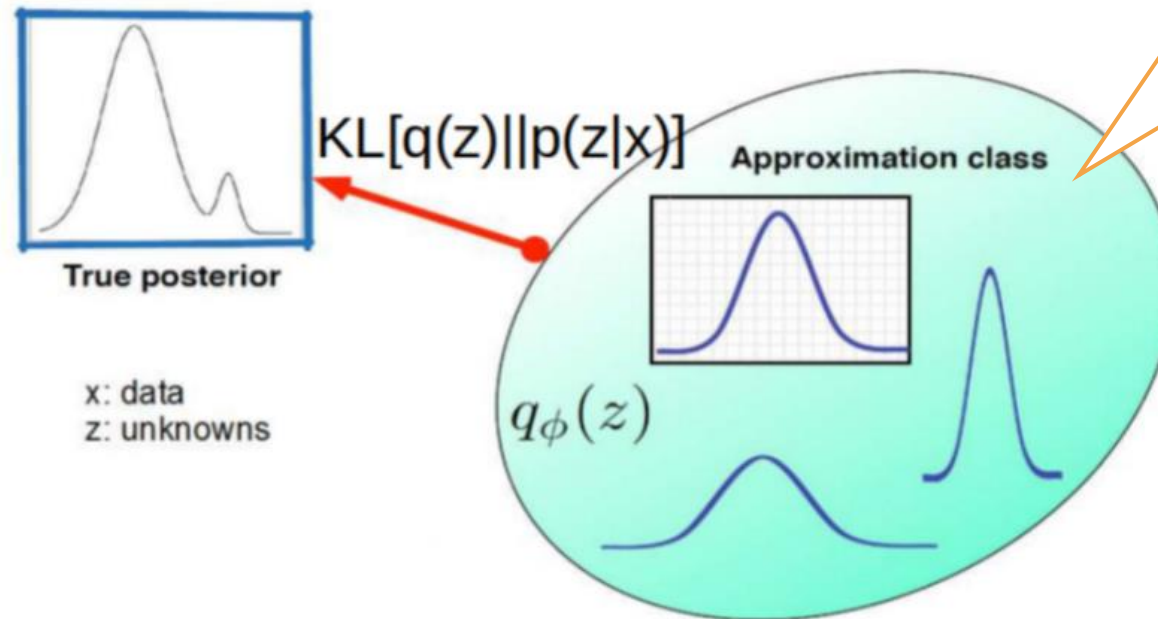
Kullback Leibler divergence $KL[q||p]$ between q and p

Also possible to use $KL[p||q]$ or divergences other than KL

$$\phi^* = \operatorname{argmin}_{\phi} KL[q_{\phi}(\mathbf{Z})||p(\mathbf{Z}|\mathcal{D}, \Theta)]$$

q_{ϕ} defines a class of distributions parametrized by ϕ sometimes called “variational parameters”

Name “variational” comes from Physics and refers to problems where we are optimizing functions of distributions (here the function is the KL divergence)



Variational Inference (VI)

- The optimization problem

$$\begin{aligned}\phi^* &= \operatorname{argmin}_{\phi} \operatorname{KL}[q_{\phi}(\mathbf{Z}) || p(\mathbf{Z}|\mathcal{D}, \Theta)] \\ &= \operatorname{argmin}_{\phi} \mathbb{E}_{q_{\phi}(\mathbf{Z})} \left[\log q_{\phi}(\mathbf{Z}) - \log \frac{p(\mathcal{D}|\mathbf{Z}, \Theta)p(\mathbf{Z}|\Theta)}{p(\mathcal{D}|\Theta)} \right] \\ &= \operatorname{argmin}_{\phi} \mathbb{E}_{q_{\phi}(\mathbf{Z})} [\log q_{\phi}(\mathbf{Z}) - \log p(\mathcal{D}|\mathbf{Z}, \Theta) - \log p(\mathbf{Z}|\Theta)] + \log p(\mathcal{D}|\Theta)\end{aligned}$$

- Since $\log p(\mathcal{D}|\Theta)$ is independent of ϕ , the optimization problem becomes

$$\phi^* = \operatorname{argmin}_{\phi} \mathbb{E}_{q_{\phi}(\mathbf{Z})} [\log q_{\phi}(\mathbf{Z}) - \log p(\mathcal{D}|\mathbf{Z}, \Theta) - \log p(\mathbf{Z}|\Theta)]$$

$$\phi^* = \operatorname{argmin}_{\phi} \mathbb{E}_{q_{\phi}(\mathbf{Z})} [\log q_{\phi}(\mathbf{Z}) - \log p(\mathcal{D}, \mathbf{Z}|\Theta)]$$

$$\phi^* = \operatorname{argmax}_{\phi} \mathbb{E}_{q_{\phi}(\mathbf{Z})} [\log p(\mathcal{D}, \mathbf{Z}|\Theta) - \log q_{\phi}(\mathbf{Z})] = \operatorname{argmax} \mathcal{L}(\phi, \Theta)$$

- Note that $\mathcal{L}(\phi, \Theta) \leq \log p(\mathcal{D}|\Theta)$ and is called “Evidence Lower Bound” (ELBO)



The ELBO

- The ELBO is defined as

$$\begin{aligned}\mathcal{L}(\phi, \Theta) &= \mathbb{E}_{q_{\phi}(\mathbf{Z})} [\log p(\mathcal{D}, \mathbf{Z} | \Theta) - \log q_{\phi}(\mathbf{Z})] \\ &= \mathbb{E}_{q_{\phi}(\mathbf{Z})} [\log p(\mathcal{D}, \mathbf{Z} | \Theta)] + H[q_{\phi}(\mathbf{Z})]\end{aligned}$$

- Thus maximizing the ELBO w.r.t. ϕ gives us a $q_{\phi}(\mathbf{Z})$ which
 - Maximizes the expected joint probability of data and latent variables
 - Has a high entropy
- We can also write the ELBO as follows

$$\mathcal{L}(\phi, \Theta) = \mathbb{E}_{q_{\phi}(\mathbf{Z})} [\log p(\mathcal{D} | \mathbf{Z}, \Theta)] - \text{KL}[q_{\phi}(\mathbf{Z}) || p(\mathbf{Z} | \Theta)]$$

- Thus maximizing the ELBO w.r.t. ϕ will give us a $q_{\phi}(\mathbf{Z})$ which
 - Explains the data \mathcal{D} well, i.e., gives it large expected probability $\mathbb{E}_q[\log p(\mathcal{D} | \mathbf{Z}, \Theta)]$
 - Is close to the prior $p(\mathbf{Z})$, i.e. is simple/regularized (small $\text{KL}[q_{\phi}(\mathbf{Z}) || p(\mathbf{Z} | \Theta)]$)



Maximizing the ELBO

Unknown Θ case later

- We need to maximize the ELBO w.r.t. ϕ (for now, assuming Θ is known)

$$\mathcal{L}(\phi, \Theta) = \mathbb{E}_{q_{\phi}(\mathbf{Z})}[\log p(\mathcal{D}|\mathbf{Z}, \Theta)] - \text{KL}[q_{\phi}(\mathbf{Z})||p(\mathbf{Z}|\Theta)]$$

- The general approach to maximize ELBO is based on gradient-based methods
 - Assume some suitable/convenient form for $q_{\phi}(\mathbf{Z})$, e.g., $\mathcal{N}(\mathbf{Z}|\mu, \Sigma)$ so $\phi = (\mu, \Sigma)$
 - Maximize the ELBO w.r.t. ϕ using gradient ascent

$$\phi_{t+1} = \phi_t + \eta_t \nabla_{\phi_t} \mathcal{L}(\phi, \Theta)$$

- Note: Expectations in ELBO and ELBO's gradients w.r.t. ϕ may not be easy
 - Will see methods to handle such issues later
 - Assuming simple forms for $q_{\phi}(\mathbf{Z})$ also helps (we can use random variable transformation methods to transform the simple form to more expressive ones – will see later)

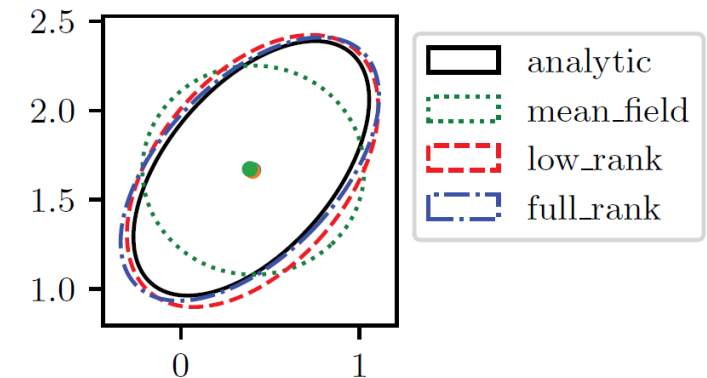


A Simple Illustration for VI

- Assume a simple likelihood model

$$p(\mathcal{D}|\mathbf{z}) = \prod_{n=1}^N \mathcal{N}(\mathbf{x}_n|\mathbf{z}, \mathbf{\Sigma}) \propto \mathcal{N}(\bar{\mathbf{x}}|\mathbf{z}, \frac{1}{N} \mathbf{\Sigma})$$

- Suppose we want to estimate the posterior of the mean \mathbf{z}
- Assuming a Gaussian prior on \mathbf{z} and assuming $\mathbf{\Sigma}$ is known, the posterior can be computed analytically (because of conjugacy)
- Let's still try VI to see how well it does
- Figure shows VI result for three Gaussian forms for $q(\mathbf{z})$
 - Low-rank: $q(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mu_{\mathbf{z}}, \mathbf{\Sigma}_{\mathbf{z}})$ where $\mathbf{\Sigma}_{\mathbf{z}} = \mathbf{L}\mathbf{L}^T$
 - Full-rank: $q(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mu_{\mathbf{z}}, \mathbf{\Sigma}_{\mathbf{z}})$ with no constraint on $\mathbf{\Sigma}_{\mathbf{z}}$
 - Mean-field: $q(\mathbf{z}) = q(z_1)q(z_2) = \mathcal{N}(z_1|\mu_{z_1}, \sigma_{z_1}^2) \mathcal{N}(z_2|\mu_{z_2}, \sigma_{z_2}^2)$



Detour

- Consider a scalar transformation of a scalar random variable \mathbf{u} as $\boldsymbol{\theta} = T(\mathbf{u})$
- Probability distributions of random variables \mathbf{u} and $\boldsymbol{\theta}$ are related as

$$p(\boldsymbol{\theta}) = p(\mathbf{u}) \left| \frac{d\mathbf{u}}{d\boldsymbol{\theta}} \right|$$

- Similarly, for multivariate random variables (of same size) related as $\boldsymbol{\theta} = T(\mathbf{u})$

$$p(\boldsymbol{\theta}) = p(\mathbf{u}) \left| \det \left(\frac{\partial \mathbf{u}}{\partial \boldsymbol{\theta}} \right) \right|$$

Absolute value of the determinant of the Jacobian
(note that $\mathbf{u} = T^{-1}(\boldsymbol{\theta})$)

- We can use such transformations for VI by using a simple distribution for $q(\mathbf{Z})$ and then transform it to a more expressive/appropriate distribution (more on this later)



Mean-Field VI

- A special way to maximize the ELBO is via the mean-field approximation
- Doesn't require specifying the form of $q(\mathbf{Z}|\phi)$ or computing ELBO's gradients
- The idea: Assumes unknowns \mathbf{Z} can be partitioned into M groups $\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_M$, s.t.,

As a shorthand, often written as
 $q = \prod_{i=1}^M q_i$ where $q_i = q(\mathbf{Z}_i|\phi_i)$

$$q(\mathbf{Z}|\phi) = \prod_{i=1}^M q(\mathbf{Z}_i|\phi_i)$$

For models with **local conjugacy**,
 it becomes super easy!

- Learning the optimal $q(\mathbf{Z}|\phi)$ reduces to learning the optimal q_1, q_2, \dots, q_M
- Can select groups based on model's structure, e.g., in Bayesian neural net for regression

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \lambda, \beta) \approx q(\mathbf{w}|\phi) = \prod_{\ell=1}^L q(\mathbf{w}^{(\ell)}|\phi_\ell)$$

Assuming a network with L
 layers, mean-field across layers

- Mean-field has limitations. Factorized form ignores the correlations among unknowns
 - Variants such as "**structured mean-field**" exist where some correlations can be modeled



Deriving Mean-Field VI Updates

Writing this is the same as $\operatorname{argmax}_{\phi} \mathcal{L}(\phi, \Theta)$. We are just writing optimization w.r.t. q directly

0

- With $q = \prod_{i=1}^M q_i$, what's the optimal q_i when we do $\operatorname{argmax}_q \mathcal{L}(q)$?
- Note that under this mean-field assumption, the ELBO simplifies to

$$\mathcal{L}(q) = \int q(\mathbf{Z}) \log \left[\frac{p(\mathcal{D}, \mathbf{Z} | \Theta)}{q(\mathbf{Z})} \right] d\mathbf{Z} = \int \prod_i q_i \left[\log p(\mathcal{D}, \mathbf{Z} | \Theta) - \sum_i \log q_i \right] d\mathbf{Z}$$

- Suppose we wish to find the optimal q_j given all other q_i 's ($i \neq j$) as fixed, then

$$\mathcal{L}(q) = \int q_j \left[\int \log p(\mathcal{D}, \mathbf{Z} | \Theta) \prod_{i \neq j} q_i d\mathbf{Z}_i \right] d\mathbf{Z}_j - \int q_j \log q_j d\mathbf{Z}_j + \text{const w.r.t. } q_j$$

$$= \int q_j \log \hat{p}(\mathcal{D}, \mathbf{Z}_j | \Theta) d\mathbf{Z}_j - \int q_j \log q_j d\mathbf{Z}_j$$

$$= -\text{KL}(q_j || \hat{p}) \quad \log \hat{p}(\mathcal{D}, \mathbf{Z}_j | \Theta) = \mathbb{E}_{i \neq j} [\log p(\mathcal{D}, \mathbf{Z} | \Theta)] + \text{const}$$

$$q_j^* = \frac{\exp(\mathbb{E}_{i \neq j} [\log p(\mathcal{D}, \mathbf{Z} | \Theta)])}{\int \exp(\mathbb{E}_{i \neq j} [\log p(\mathcal{D}, \mathbf{Z} | \Theta)] d\mathbf{Z}_j}$$

- Thus $q_j^* = \operatorname{argmax}_{q_j} \mathcal{L}(q) = \operatorname{argmin}_{q_j} \text{KL}(q_j || \hat{p}) = \hat{p}(\mathcal{D}, \mathbf{Z}_j | \Theta)$



Deriving Mean-Field VI Updates

- So we saw that the optimal q_j when doing mean-field VI is

$$q_j^*(\mathbf{Z}_j) = \frac{\exp(\mathbb{E}_{i \neq j}[\log p(\mathcal{D}, \mathbf{Z}|\Theta)])}{\int \exp(\mathbb{E}_{i \neq j}[\log p(\mathcal{D}, \mathbf{Z}|\Theta)] d\mathbf{Z}_j}$$

- Note: Can often just compute the numerator and recognize denominator by inspection
- **Important:** For locally conj models, $q_j^*(\mathbf{Z}_j)$ will have the same form as prior $p(\mathbf{Z}_j|\Theta)$
 - Only the distribution parameters will be different
- **Important:** For estimating q_j the required expectation depends on other $\{q_i\}_{i \neq j}$
 - Thus we use an alternating update scheme for these
- Guaranteed to converge (to a local optima)
 - We are basically solving a sequence of **concave maximization** problems
 - Reason: $\mathcal{L}(q) = \int q_j \log \hat{p}(\mathcal{D}, \mathbf{Z}_j|\Theta) d\mathbf{Z}_j - \int q_j \log q_j d\mathbf{Z}_j$ is concave in q_j



The Mean-Field VI Algorithm

- Also known as **Co-ordinate Ascent Variational Inference** (CAVI) Algorithm
- Input: Model in form of priors and likelihood, or joint $p(\mathcal{D}, \mathbf{Z}|\Theta)$, Data \mathcal{D}
- Output: A variational distribution $q(\mathbf{Z}) = \prod_{j=1}^M q_j(\mathbf{Z}_j)$
- Initialize: Variational distributions $q_j(\mathbf{Z}_j)$, $j = 1, 2, \dots, M$
- While the ELBO has not converged
 - For each $j = 1, 2, \dots, M$, set

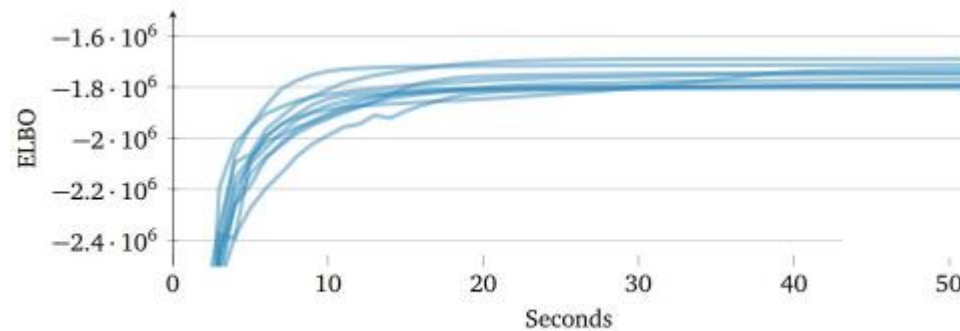
$$q_j(\mathbf{Z}_j) \propto \exp(\mathbb{E}_{i \neq j}[\log p(\mathcal{D}, \mathbf{Z}|\Theta)])$$

- Compute ELBO $\mathcal{L}(q) = \mathbb{E}_q[\log p(\mathcal{D}, \mathbf{Z}|\Theta)] - \mathbb{E}_q[\log q(\mathbf{Z})]$
- NOTE: We can also use mean-field assumption for $q(\mathbf{Z})$ and optimize the ELBO using gradient based methods if we don't have local conjugacy



VI and Convergence

- VI is guaranteed to converge to a local optima (just like EM)
- Therefore proper initialization is important (just like EM)
 - Can sometimes run multiple times with different initializations and choose the best run



Different initializations may lead to different optima

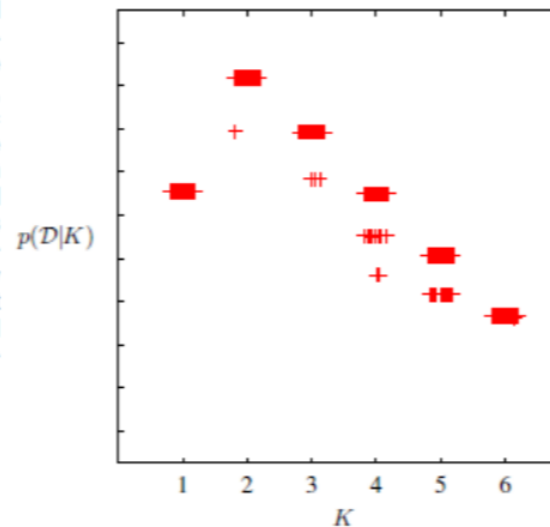
- ELBO increases monotonically with iterations
 - Can thus monitor the ELBO to assess convergence



ELBO for Model Selection

- Recall that ELBO is a lower bound on log of model evidence $\log p(\mathbf{X}|\mathbf{m})$
- Can compute ELBO for each model \mathbf{m} and choose the one with largest ELBO

Plot of the variational lower bound \mathcal{L} versus the number K of components in the Gaussian mixture model, for the Old Faithful data, showing a distinct peak at $K = 2$ components. For each value of K , the model is trained from 100 different random starts, and the results shown as '+' symbols plotted with small random horizontal perturbations so that they can be distinguished. Note that some solutions find suboptimal local maxima, but that this happens infrequently.



Each value of K represents a different model

- Some criticism since we are using a lower-bound but often works well in practice

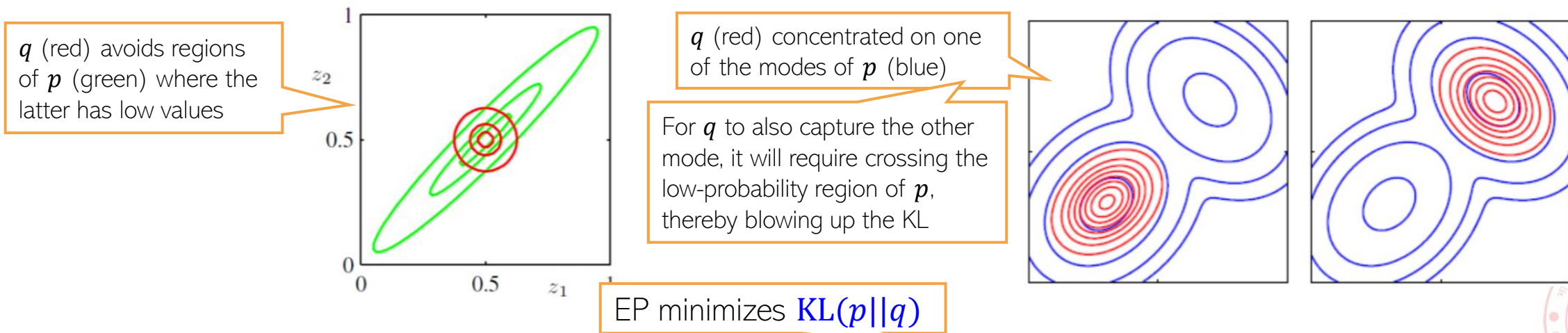


VI might under-estimate posterior's variance

- Recall that VI approximates a posterior p by finding q that minimizes $\text{KL}(q||p)$

$$\text{KL}(q||p) = - \int q(\mathbf{Z}) \log \left\{ \frac{p(\mathbf{Z}|\mathcal{D})}{q(\mathbf{Z})} \right\} d\mathbf{Z}$$

- $q(\mathbf{Z})$ will be small where $p(\mathbf{Z}|\mathcal{D})$ is small otherwise KL will blow up
- Thus $q(\mathbf{Z})$ avoids low-probability regions of the true posterior



- Some methods, e.g., Expectation Propagation (EP), can avoid this behavior

Variational EM

- If the parameters Θ are also unknown then we can use variational EM (VEM)
- VEM is the same as EM except the E step uses VI to approximate the CP of \mathbf{Z}
- VEM alternates between the following two steps
 - Maximize the ELBO w.r.t. ϕ (gives the variational approximation $q(\mathbf{Z})$ of CP of \mathbf{Z})

$$\phi^{(t)} = \operatorname{argmax}_{\phi} \mathbb{E}_{q_{\phi}(\mathbf{Z})} [\log p(\mathcal{D}, \mathbf{Z} | \Theta^{(t-1)}) - \log q_{\phi}(\mathbf{Z})]$$

- Maximize the ELBO w.r.t. Θ (gives us point estimate of Θ)

$$\begin{aligned} \Theta^{(t)} &= \operatorname{argmax}_{\Theta} \mathbb{E}_{q_{\phi^{(t)}}(\mathbf{Z})} [\log p(\mathcal{D}, \mathbf{Z} | \Theta) - \log q_{\phi^{(t)}}(\mathbf{Z})] \\ &= \operatorname{argmax}_{\Theta} \mathbb{E}_{q_{\phi^{(t)}}(\mathbf{Z})} [\log p(\mathcal{D}, \mathbf{Z} | \Theta)] \end{aligned}$$

This looks very similar to the expected CLL with the CP replaced by its variational approximation

- Note: If we want posterior for Θ as well, treat it similar to \mathbf{Z} and apply variational approximation (instead of using VEM) if the posterior isn't tractable



Extra Slides - Mean-Field VI: A Simple Example

- Consider data $\mathbf{X} = \{x_1, x_2, \dots, x_N\}$ from a one-dim Gaussian $\mathcal{N}(\mu, \tau^{-1})$
- Assume the following normal-gamma prior on μ and τ

$$p(\mu|\tau) = \mathcal{N}(\mu|\mu_0, (\lambda_0\tau)^{-1}) \quad p(\tau) = \text{Gamma}(\tau|a_0, b_0)$$

- Posterior is also normal-gamma due to the jointly conjugate prior
- Let's anyway verify this by trying mean-field VI for this model
- With mean-field assumption on the variational posterior $q(\mu, \tau) = q_\mu(\mu)q_\tau(\tau)$

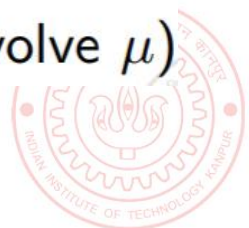
$$\log q_\mu^*(\mu) = \mathbb{E}_{q_\tau} [\log p(\mathbf{X}, \mu, \tau)] + \text{const}$$

$$\log q_\tau^*(\tau) = \mathbb{E}_{q_\mu} [\log p(\mathbf{X}, \mu, \tau)] + \text{const}$$

- In this example, the log-joint $\log p(\mathbf{X}, \mu, \tau) = \log p(\mathbf{X}|\mu, \tau) + \log p(\mu|\tau) + \log p(\tau)$. Thus

$$\log q_\mu^*(\mu) = \mathbb{E}_{q_\tau} [\log p(\mathbf{X}|\mu, \tau) + \log p(\mu|\tau)] + \text{const} \quad (\text{only keeping terms that involve } \mu)$$

$$\log q_\tau^*(\tau) = \mathbb{E}_{q_\mu} [\log p(\mathbf{X}|\mu, \tau) + \log p(\mu|\tau) + \log p(\tau)] + \text{const}$$



Extra Slides - Mean-Field VI: A Simple Example

- Substituting $p(\mathbf{X}|\mu, \tau) = \prod_{n=1}^N p(x_n|\mu, \tau)$ and $p(\mu|\tau)$, we get

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

- (Verify) The above is log of a Gaussian. This $q_\mu^* = \mathcal{N}(\mu|\mu_N, \lambda_N)$ with

$$\mu_N = \frac{\lambda_0 \mu_0 + N \bar{x}}{\lambda_0 + N} \quad \text{and} \quad \lambda_N = (\lambda_0 + N) \mathbb{E}_{q_\tau}[\tau]$$

This update depends on q_τ

- Proceeding in a similar way (verify), we can show that $q_\tau^* = \text{Gamma}(\tau|a_N, b_N)$

$$a_N = a_0 + \frac{N+1}{2} \quad \text{and} \quad b_N = b_0 + \frac{1}{2} \mathbb{E}_{q_\mu} \left[\sum_{n=1}^N (x_n - \mu)^2 + \lambda_0(\mu - \mu_0)^2 \right]$$

This update depends on q_μ

- Note: Updates of q_μ^* and q_τ^* depend on each other (hence alternating updates needed)

