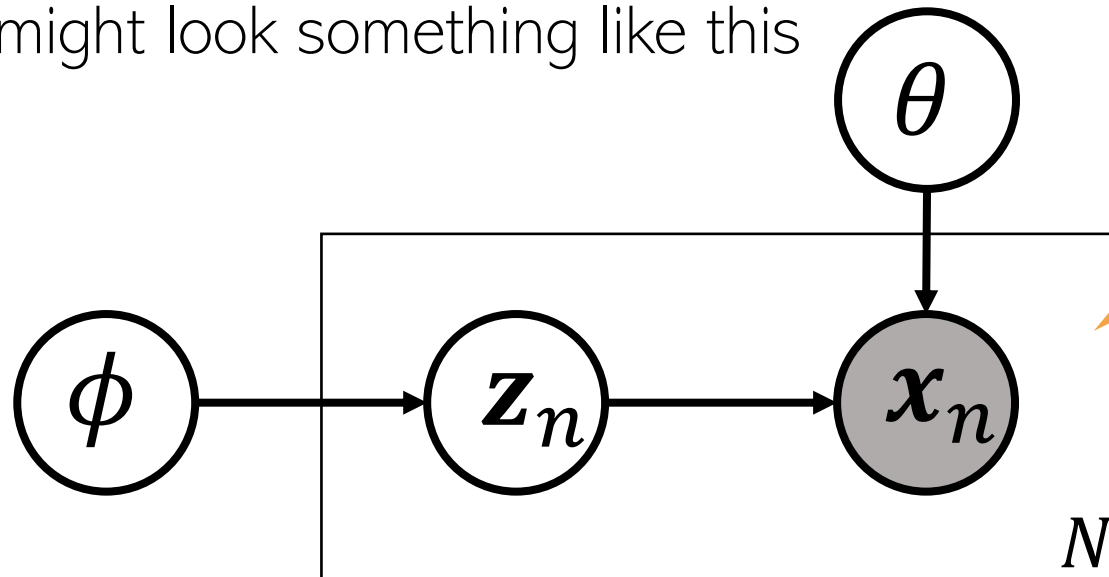# Variational Inference

CS772A: Probabilistic Machine Learning

Piyush Rai

# Variational Inference (VI)

- Assume a latent variable model with data $\boldsymbol{\mathcal{D}}$ and latent variables $\boldsymbol{Z}$

- 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 $p(\boldsymbol{\mathcal{D}}|\boldsymbol{Z}, \Theta)$ and prior is $p(\boldsymbol{Z}|\Theta)$. Want posterior over $\boldsymbol{Z}$

- $\Theta = (\theta, \phi)$ denotes the other parameters that define the likelihood and the prior

- For now, assume $\Theta$ is known and only $\boldsymbol{Z}$ is unknown (the $\Theta$ unknown case later)

- Assume CP $p(\boldsymbol{Z}|\boldsymbol{\mathcal{D}}, \Theta)$ 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 $\phi$ which makes our approximation $q(\mathbf{Z}|\phi)$ as closed to the true as possible to the true posterior $p(\mathbf{Z}|\mathcal{D})$
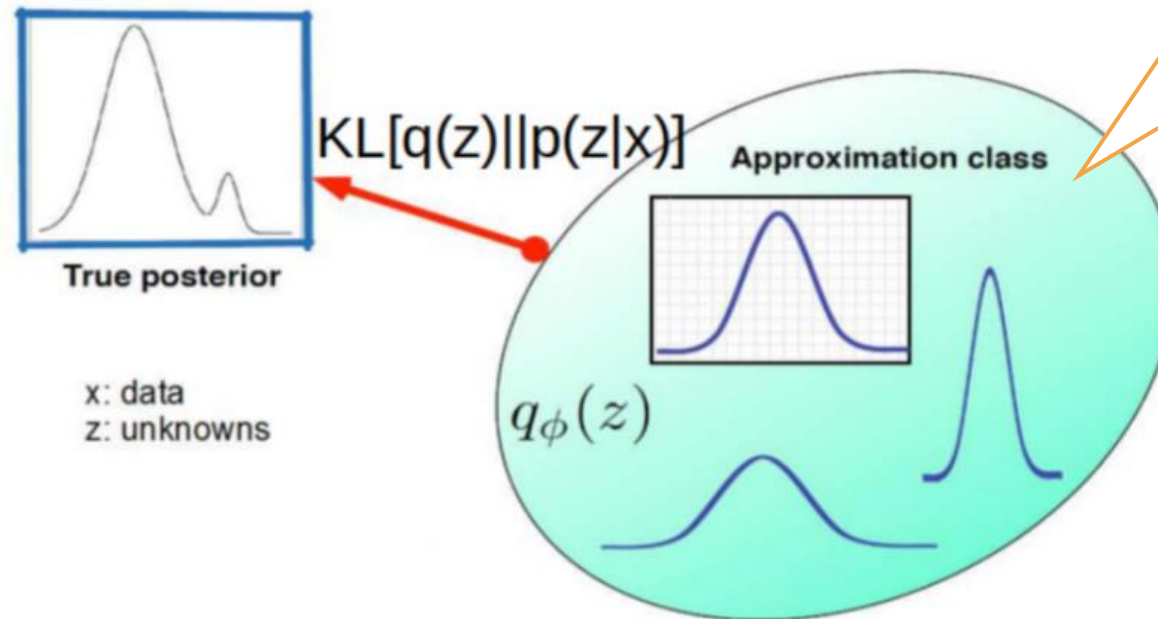
Kullback Leibler divergence $\mathrm{KL}[q||p]$ between $q$ and $p$

Also possible to use $\mathrm{KL}[p||q]$ or divergences other than KL

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

$q_\phi$ defines a class of distributions parametrized by $\phi$ 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)



KL[q(z)||p(z|x)]

True posterior

x: data
z: unknowns

Approximation class

$q_\phi(z)$

# Variational Inference (VI)

- The optimization problem

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

$$= \text{argmin}_\phi \ \mathbb{E}_{q_\phi(\boldsymbol{Z})}\left[\log q_\phi(\boldsymbol{Z}) - \log \frac{p(\boldsymbol{\mathcal{D}}|\boldsymbol{Z}, \Theta)p(\boldsymbol{Z}|\Theta)}{p(\boldsymbol{\mathcal{D}}|\Theta)}\right]$$

$$= \text{argmin}_\phi \ \mathbb{E}_{q_\phi(\boldsymbol{Z})}\left[\log q_\phi(\boldsymbol{Z}) - \log p(\boldsymbol{\mathcal{D}}|\boldsymbol{Z}, \Theta) - \log p(\boldsymbol{Z}|\Theta)\right] + \log p(\boldsymbol{\mathcal{D}}|\Theta)$$

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

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

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

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

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

# The ELBO

- The ELBO is defined as

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

$$= \mathbb{E}_{q_\phi(\mathbf{Z})}[\log p(\mathcal{D}, \mathbf{Z}|\Theta)] + \mathrm{H}[q_\phi(\mathbf{Z})]$$

- Thus maximizing the ELBO w.r.t. $\phi$ 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

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

- Thus maximizing the ELBO w.r.t. $\phi$ 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 $\mathrm{KL}[q_\phi(\mathbf{Z})||p(\mathbf{Z}|\Theta))$

# Maximizing the ELBO

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

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

- The general approach to maximize ELBO is based on gradient-based methods
  - Assume some suitable/convenient form for $q_\phi(\boldsymbol{Z})$, e.g., $\mathcal{N}(\boldsymbol{Z}|\mu, \Sigma)$ so $\phi = (\mu, \Sigma)$
  - Maximize the ELBO w.r.t. $\phi$ 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. $\phi$ may not be easy
  - Will see methods to handle such issues later
  - Assuming simple forms for $q_\phi(\boldsymbol{Z})$ also helps (we can use random variable transformation methods to transform the simple form to more expressive ones – will see later)
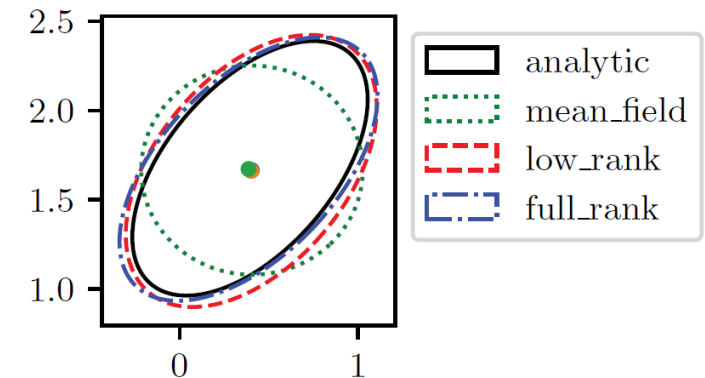
# A Simple Illustration for VI

- Assume a simple likelihood model

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

- Suppose we want to estimate the posterior of the mean $z$

- Assuming a Gaussian prior on $z$ and assuming $\Sigma$ 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 $q(z)$
    - Low-rank: $q(z) = \mathcal{N}(z|\mu_z, \Sigma_z)$ where $\Sigma_z = LL^{\mathsf{T}}$
    - Full-rank: $q(z) = \mathcal{N}(z|\mu_z, \Sigma_z)$ with no constraint on $\Sigma_z$
    - Mean-field: $q(z) = q(z_1)q(z_2) = \mathcal{N}(z_1|\mu_{z_1}, \sigma_{z_1}^2)\,\mathcal{N}(z_2|\mu_{z_2}, \sigma_{z_2}^2)$



(Example courtesy: PML-2 (Murphy))

# Detour

Transformed random variable

A one-to-one transformation function

- Consider a scalar transformation of a scalar random variable $u$ as $\theta = T(u)$

- Probability distributions of random variables $u$ and $\theta$ are related as

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

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

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

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

- 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 $q(\boldsymbol{Z}|\phi)$ or computing ELBO's gradients

- The idea: Assumes unknowns $\boldsymbol{Z}$ can be partitioned into $M$ groups $\boldsymbol{Z}_1, \boldsymbol{Z}_2, \ldots, \boldsymbol{Z}_M$, s.t.,

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

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

For models with local conjugacy, it becomes super easy!

- Learning the optimal $q(\boldsymbol{Z}|\phi)$ reduces to learning the optimal $q_1, q_2, \ldots, q_M$

- Can select groupsbased on model's structure, e.g., in Bayesian neural net for regression

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

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 $\mathbf{argmax}_\phi \, \mathcal{L}(\phi, \Theta)$. We are just writing optimization w.r.t. $q$ directly

- With $q = \prod_{i=1}^{M} q_i$, what's the optimal $q_i$ when we do $\mathbf{argmax}_q \, \mathcal{L}(q)$?

- Note that under this mean-field assumption, the ELBO simplifies to

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

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

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

$$= \int q_j \log \hat{p}(\mathbf{D}, Z_j|\Theta) \, dZ_j - \int q_j \log q_j Z_j$$

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

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

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

# 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 $p(Z_j|\Theta)$
  - 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: $\mathcal{L}(q) = \int q_j \log \hat{p}(\mathcal{D}, Z_j|\Theta) \, Z_j - \int q_j \log q_j Z_j$ 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 $q_j(\mathbf{Z}_j), j = 1, 2, \dots M$

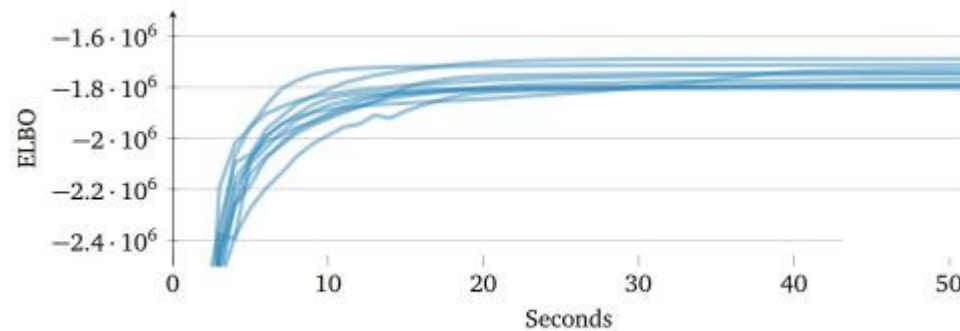- While the ELBO has not converged
  - For each $j = 1, 2, \dots M$, set

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

  - 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
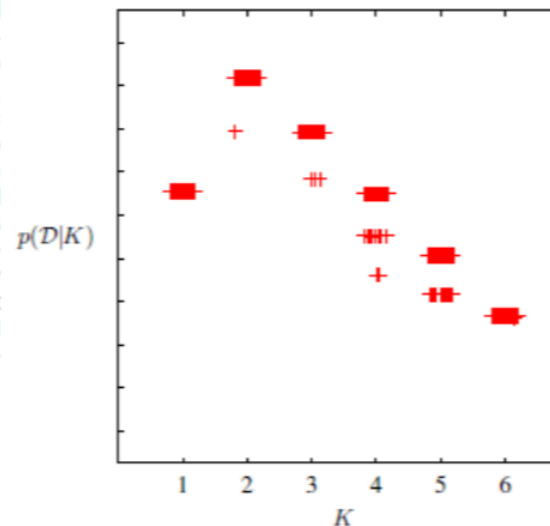


Different initializations may lead to different optima

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

# ELBO for Model Selection

- Recall that ELBO is a <u>lower bound</u> on log of model evidence $\log p(X|m)$
- Can compute ELBO for each model $m$ 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.

$P(\mathcal{D}|K)$



Each value of $K$ represents a different model

- Some criticism since we are using a lower-bound but often works well in practice
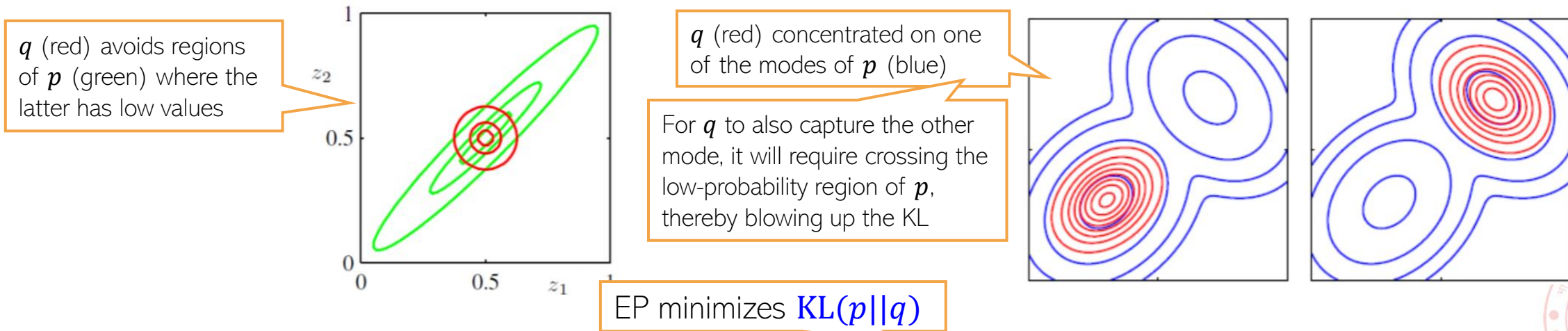
# VI might <u>under-estimate</u> posterior's variance

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

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

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



$q$ (red) avoids regions of $p$ (green) where the latter has low values

$q$ (red) concentrated on one of the modes of $p$ (blue)

For $q$ to also capture the other mode, it will require crossing the low-probability region of $p$, thereby blowing up the KL

EP minimizes KL$(p||q)$

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

# Variational EM

- If the parameters $\Theta$ 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 $\boldsymbol{Z}$

- VEM alternates between the following two steps

  - Maximize the ELBO w.r.t. $\boldsymbol{\phi}$ (gives the variational approximation $q(\boldsymbol{Z})$ of CP of $\boldsymbol{Z}$)

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

  - Maximize the ELBO w.r.t. $\Theta$ (gives us point estimate of $\Theta$)

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

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

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

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

# Extra Slides - Mean-Field VI: A Simple Example

- Consider data $\mathbf{X} = \{x_1, x_2, \ldots, x_N\}$ from a one-dim Gaussian $\mathcal{N}(\mu, \tau^{-1})$

- Assume the following normal-gamma prior on $\mu$ and $\tau$

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

- 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 $q(\mu, \tau) = q_\mu(\mu) q_\tau(\tau)$

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

- 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

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

# 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 $q_\tau$

- 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 $q_\mu$

- Note: Updates of $q_\mu^*$ and $q_\tau^*$ depend on each other (hence alternating updates needed)