

VI Talk

roko

June 2023

1 Introduction

Suppose we believe something follows a latent variable model (LVM).

Idea: Bayesian Inference to the rescue!

Start with:

Prior: $p_\theta(z)$

Likelihood: $p_\theta(x|z)$

End with:

$$\text{Posterior: } \frac{p_\theta(x|z) * p_\theta(z)}{\int p_\theta(x|z) * p_\theta(z) dx}$$

Oh... this is bad (intractable).

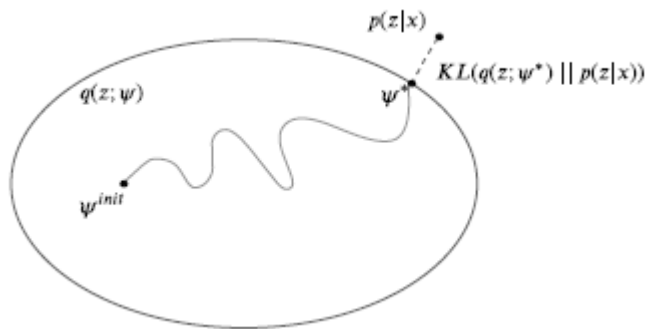


Figure 1: Caption

Approach: Use the best q distribution out of the distribution family \mathcal{Q} such that:

$$q = \underset{q \in \mathcal{Q}}{\operatorname{argmin}} D_{\mathbb{KL}}(q(z) \| p_{\boldsymbol{\theta}}(z | \mathbf{x}))$$

Solution:

$$\begin{aligned} \psi^* &= \underset{\psi}{\operatorname{argmin}} D_{\mathbb{KL}}(q_{\psi}(z) \| p_{\boldsymbol{\theta}}(z | \mathbf{x})) \\ &= \underset{\psi}{\operatorname{argmin}} \mathbb{E}_{q_{\psi}(z)} \left[\log q_{\psi}(z) - \log \left(\frac{p_{\boldsymbol{\theta}}(\mathbf{x} | z) p_{\boldsymbol{\theta}}(z)}{p_{\boldsymbol{\theta}}(\mathbf{x})} \right) \right] \\ &= \underset{\psi}{\operatorname{argmin}} \mathbb{E}_{q_{\psi}(z)} \underbrace{[\log q_{\psi}(z) - \log p_{\boldsymbol{\theta}}(\mathbf{x} | z) - \log p_{\boldsymbol{\theta}}(z)]}_{\mathcal{L}(\boldsymbol{\theta}, \psi | \mathbf{x})} + \log p_{\boldsymbol{\theta}}(\mathbf{x}) \\ &= \underset{\psi}{\operatorname{argmin}} \mathbb{E}_{q_{\psi}(z)} [\log q_{\psi}(z) - \log p_{\boldsymbol{\theta}}(\mathbf{x}, z)] \\ &= \underset{\psi}{\operatorname{argmin}} \mathbb{E}_{q_{\psi}(z)} [-\log p_{\boldsymbol{\theta}}(\mathbf{x}, z) + \log q_{\psi}(z)] \end{aligned}$$

Notice that the last term in line 3 is constant in terms of ψ so we can drop it. Note that because a KL divergence is always above 0, we have that (from line 3)

$$D_{\mathbb{KL}}(q_{\psi}(z) \| p_{\boldsymbol{\theta}}(z | \mathbf{x})) = \mathcal{L}(\boldsymbol{\theta}, \psi | \mathbf{x}) + \log p_{\boldsymbol{\theta}}(\mathbf{x}) \geq 0$$

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) \geq -\mathcal{L}(\boldsymbol{\theta}, \psi | \mathbf{x}) := \text{ELBO}$$

$$\text{ELBO} = \mathbb{E}_{q_{\psi}(z)} [\log p_{\boldsymbol{\theta}}(\mathbf{x}, z) - \log q_{\psi}(z)]$$

The equation in blue is VERY important. It's important that we keep this equation in mind as we perform VI going forward. Keep in mind that the joint likelihood can ALWAYS be decomposed into the product of the likelihood and prior. It's also nice that you are maximizing a lower bound of the evidence. So higher ELBO implies a higher MINIMUM evidence.

$$\begin{aligned} \mathcal{L}(\psi | \boldsymbol{\theta}, \mathbf{x}) &= \mathbb{E}_{q_{\psi}(z)} [\log p_{\boldsymbol{\theta}}(\mathbf{x} | z) + \log p_{\boldsymbol{\theta}}(z) - \log q_{\psi}(z)] \\ &= \mathbb{E}_{q_{\psi}(z)} [\log p_{\boldsymbol{\theta}}(\mathbf{x} | z)] + \int q_{\psi}(z) (\log p_{\boldsymbol{\theta}}(z) - \log q_{\psi}(z)) dz \\ &= \mathbb{E}_{q_{\psi}(z)} [\log p_{\boldsymbol{\theta}}(\mathbf{x} | z)] + \int q_{\psi}(z) \frac{\log p_{\boldsymbol{\theta}}(z)}{\log q_{\psi}(z)} dz \\ &= \mathbb{E}_{q_{\psi}(z)} [\log p_{\boldsymbol{\theta}}(\mathbf{x} | z)] - \int q_{\psi}(z) \frac{\log q_{\psi}(z)}{\log p_{\boldsymbol{\theta}}(z)} dz \\ &= \mathbb{E}_{q_{\psi}(z)} [\log p_{\boldsymbol{\theta}}(\mathbf{x} | z)] - D_{\mathbb{KL}}(q_{\psi}(z) \| p_{\boldsymbol{\theta}}(z)) \end{aligned}$$

ELBO = expected log likelihood - KL from posterior to prior

2 Variational Posterior Forms

2.1 Fixed Form VI

Fixed-Form VI means that you pick a convenient functional form (i.e. MV Gaussian) for the posterior $q_\psi(z|x)$ and the ELBO is optimized with gradient based methods.

2.2 Free Form VI

Free-Form VI is VI that is optimized one parameter at a time in a coordinate ascent manner. This is typically done with the mean field assumption by assuming the posterior factorizes:

$$q_\phi(z|x) = \prod_{j=1}^J q_j(z_j)$$

3 Parameter Estimation

3.1 MLE for LVMs

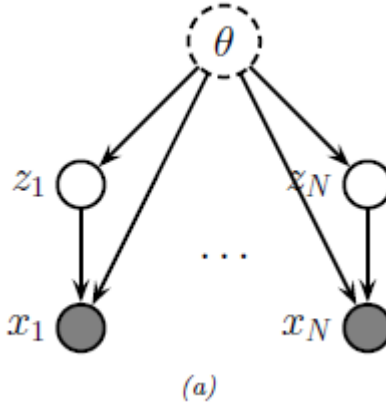


Figure 2: Caption

Suppose we have a latent variable model of the form.

$$p(\mathcal{D}, z_{1:N} | \theta) = \prod_{n=1}^N p(z_n | \theta) p(x_n | z_n, \theta)$$

How do you find the MLE of θ given the observable data? With unobservable latent variables, we can't take them as given in the likelihood function and need to marginalize them out.

$$\log p(\mathbf{x}_n | \boldsymbol{\theta}) = \log \left[\int p(\mathbf{x}_n | z_n, \boldsymbol{\theta}) p(z_n | \boldsymbol{\theta}) dz_n \right]$$

This integral is intractable, but the ELBO is not!

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\psi}_{1:N} | \mathcal{D}) = \sum_{n=1}^N \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\psi}_n | \mathbf{x}_n) \leq \log p(\mathcal{D} | \boldsymbol{\theta})$$

We can optimize this using VEM:

Algorithm 10.1: Amortized stochastic variational EM

```

1 Initialize  $\boldsymbol{\theta}, \boldsymbol{\phi}$ 
2 repeat
3   | Sample  $\mathbf{x}_n \sim p_{\mathcal{D}}$ 
4   | E step:  $\boldsymbol{\phi} = \operatorname{argmax}_{\boldsymbol{\phi}} \mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\phi} | \mathbf{x}_n)$ 
5   | M step:  $\boldsymbol{\theta} = \operatorname{argmax}_{\boldsymbol{\theta}} \mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\phi} | \mathbf{x}_n)$ 
6 until converged

```

Figure 3: Caption

Note that unlike the traditional EM algorithm, where the E-step involves taking an expectation for missing/unknown values, VEM's E-step is not an expectation. An EM algorithm might be a borderline misuse of the term. The important thing is that we update variational parameters ψ , and then use them for sampling purposes when learning model parameters θ .

3.2 Empirical Bayes

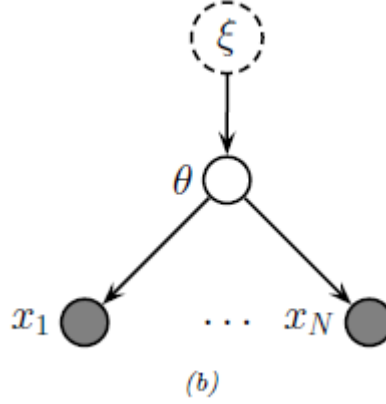


Figure 4: Caption

If you have a model with no local latents, we can try for an (empirical) Bayes approach. In this case, we assume the model has the following graphical decomposition

$$p(\mathcal{D}, \boldsymbol{\theta} \mid \boldsymbol{\xi}) = p(\boldsymbol{\theta} \mid \boldsymbol{\xi}) \prod_{n=1}^N p(\mathbf{x}_n \mid \boldsymbol{\theta})$$

Obviously the goal is to compute the posterior:

$$p(\boldsymbol{\theta} \mid \mathcal{D}, \boldsymbol{\xi}) = \frac{p(\mathcal{D} \mid \boldsymbol{\theta})p(\boldsymbol{\theta} \mid \boldsymbol{\xi})}{p(\mathcal{D} \mid \boldsymbol{\xi})}$$

Option 1: Use standard Bayesian inference. This assumes we "know" (more like have a reasonable guess at) the values of $\boldsymbol{\xi}$.

Option 2: Use empirical Bayes if hyperprior parameters are not known. Basically, maximize the hyperprior parameters s.t. the maximize the likelihood.

$$\hat{\boldsymbol{\xi}} = \underset{\boldsymbol{\xi}}{\operatorname{argmax}} \log p(\mathcal{D} \mid \boldsymbol{\xi})$$

We can use variational EM to compute this. The parameter to be estimated are $\boldsymbol{\xi}$, the latent variables are the shared global parameters $\boldsymbol{\theta}$, and the observations are \mathcal{D} . We then get the lower bound

$$\log p(\mathcal{D} \mid \boldsymbol{\xi}) \geq L(\boldsymbol{\xi}, \boldsymbol{\psi} \mid \mathcal{D}) = \mathbb{E}_{q_{\boldsymbol{\psi}}(\boldsymbol{\theta})} \left[\sum_{n=1}^N \log p(\mathbf{x}_n \mid \boldsymbol{\theta}) \right] - D_{\text{KL}}(q_{\boldsymbol{\psi}}(\boldsymbol{\theta}) \parallel p(\boldsymbol{\theta} \mid \boldsymbol{\xi}))$$

If ξ is fixed, then we only need to optimize the variational parameters and this puts in the case of Option 1, called Variational Bayes.

4 VI in Practice

4.1 Stochastic VI

Now that we have seen the ways with which we can leverage variational inference to find an approximating distribution in theory, how does it work in practice?

After all, we have assumed thus far that we have access to all N samples and that it's computationally feasible to use all of them (spoiler alert, it's not). Luckily,

$$L(\psi, \theta \mid \mathbf{x}) = L(\psi_{1:N}, \theta \mid \mathbf{x}_n) = \sum_{n=1}^N L(\psi_n, \theta \mid \mathbf{x}) = \frac{N}{B} \sum_{x_n \in B} L(\psi_{x_n}, \theta \mid \mathbf{x}_n)$$

Similar to how we do not need to process all of the data to make an update in gradient descent, we don't need all of the samples to make an update in maximizing the ELBO.

4.2 Amortized VI

First of all, what does it mean for something to be amortized?

amortized: optimized in a shared way (usually via parameter sharing)

This can be useful because we can optimize several parameters at once. By allowing our ψ_n values to be predicted all at once as outputs of a neural network (the inference network in our case), we can avoid having to optimize them separately.

$$q(\mathbf{z}_n \mid \psi_n) = q(\mathbf{z}_n \mid f_\phi^{\text{inf}}(\mathbf{x}_n)) = q_\phi(\mathbf{z}_n \mid \mathbf{x}_n)$$

So, if we created the inference net for amortized inference and combine it with Stochastic VI, we get an ELBO formula of...

$$L(\psi, \theta \mid \mathcal{D}) = \frac{N}{B} \sum_{x_n \in B} \mathbb{E}_{q_\phi(\mathbf{z}_n \mid \mathbf{x}_n)} [\log p_\theta(\mathbf{x}_n, \mathbf{z}_n) - \log q_\phi(\mathbf{z}_n \mid \mathbf{x}_n)]$$

4.2.1 Semi-amortized VI

Amortized inference is no free lunch though (what is though really)...

There exists a tradeoff between joint optimization and optimizing each variational parameter in isolation: accuracy vs. computational expense. The gap

in accuracy is called the amortization gap. A possible solution is to use the inference net as a warm start for the initial guess of each variational parameter before individually learning optimal values.

5 VI and the Reparameterization Trick

How exactly do we optimize over distributions exactly? In this section we tackle the ability to take gradients over *seemingly* random variables.

The key is to rewrite the random variable $\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})$ as a differentiable AND invertible transformation g in terms of another random variable $\epsilon \sim p(\epsilon)$ which does not depend on ϕ .

If we can create such a reparameterization, we can create move the gradient under the expectation and create a tractable gradient update.

$$\nabla_\phi \mathbb{E}_{q_\phi(z|x)}[f(z)] = \nabla_\phi \mathbb{E}_{p(\epsilon)}[f(z)] = \mathbb{E}_{p(\epsilon)}[\nabla_\phi f(z)]$$

5.1 Example: Isotropic Gaussian

$$\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$z = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon}$$

$$(\boldsymbol{\mu}, \log \boldsymbol{\sigma}) = f_\phi^{\text{inf}}(\mathbf{x})$$

Recall that right now, the ELBO formula for a single sample is:

$$\log p_\theta(\mathbf{x}_n, \mathbf{z}_n) - \log q_\phi(\mathbf{z}_n|\mathbf{x}_n)$$

The likelihood is well defined in our model. In order to evaluate the term $q_\phi(\mathbf{z}_n|\mathbf{x}_n)$, we need to leverage the change of variables formula:

$$\begin{aligned} q_\phi(z|x) &= p_\epsilon(g^{-1}(z)) * \left| \det \left(\frac{\partial}{\partial z} g^{-1}(z) \right) \right| = p_\epsilon((z - \mu)/\sigma) * \left| \det \left(\frac{\partial}{\partial z} (z - \mu) ./ \sigma \right) \right| = p_\epsilon(\epsilon) \left| \det \left(\frac{\partial}{\partial z} \epsilon \right) \right| \\ &= p_\epsilon(\epsilon) \left| \det \left(\frac{\partial \epsilon}{\partial z} \right) \right| = p_\epsilon(\epsilon) \left| \det \left(\frac{\partial \epsilon}{\partial z} \right) \right| = p_\epsilon(\epsilon) \frac{1}{\left| \det \left(\frac{\partial z}{\partial \epsilon} \right) \right|} \end{aligned}$$

Taking the logarithm of both sides gives up

$$\log q_\phi(z|x) = \log p_\epsilon(\epsilon) - \log \left(\left| \det \left(\frac{\partial z}{\partial \epsilon} \right) \right| \right)$$

where the partial derviative is the Jacobian.

5.2 Example: Full Rank Gaussian

$$\begin{aligned}\epsilon &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ \Sigma &= LL^T \text{ (Cholesky decomposition)} \\ z &= \boldsymbol{\mu} + \mathbf{L}\boldsymbol{\epsilon} \\ (\boldsymbol{\mu}, \log \sigma, \mathbf{L}') &= f_{\phi}^{\text{inf}}(\mathbf{x})\end{aligned}$$

5.3 Example: Low Rank Gaussian

$$\begin{aligned}\epsilon_1, \epsilon_2 &\stackrel{i.i.d.}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ z &= \boldsymbol{\mu} + \mathbf{B}\boldsymbol{\epsilon}_1 + \mathbf{C}\boldsymbol{\epsilon}_2 \\ B &\in \mathbb{R}^{d \times f} \\ C &\in \mathbb{R}^{d \times d} \text{ (Diagonal)} \\ \Sigma &= BB^T + C^2 \\ (\mathbf{B}, \text{diag}(\mathbf{C}), \boldsymbol{\mu}) &= f_{\phi}^{\text{inf}}(\mathbf{x})\end{aligned}$$

It may not be immediately clear why this paradigm makes sense until you break it down to properties of covariance.

$$\begin{aligned}\text{Cov}(\mathbf{z}) &= \text{Cov}(\boldsymbol{\mu} + \mathbf{B}\boldsymbol{\epsilon}_1 + \mathbf{C}\boldsymbol{\epsilon}_2) = \text{Cov}(\mathbf{B}\boldsymbol{\epsilon}_1 + \mathbf{C}\boldsymbol{\epsilon}_2) = \mathbf{B}\text{Cov}(\boldsymbol{\epsilon}_1)\mathbf{B}^T + \mathbf{C}\text{Cov}(\boldsymbol{\epsilon}_2)\mathbf{C}^T \\ &= \mathbf{B}\mathbf{B}^T + \mathbf{C}\mathbf{C}^T = \mathbf{B}\mathbf{B}^T + \mathbf{C}^2\end{aligned}$$

This type of setup is useful if a full Gaussian rank requires learning TOO many parameters. You go from learning d (diagonals) + $d(d+1)/2$ (off diagonals) in the full case to fd (\mathbf{B}) + d (diag \mathbf{C}) + d ($\boldsymbol{\mu}$) = $(f+2)d$ in the low rank case.

6 ADVI

Notice that VI seems like a tool that has free range and for the most part does not have many restrictions (barring valid transformations in the reparameterization trick).

But variances need to be positive, and probabilities need to be in $[0, 1]$. So, in order to use gradient based methods to learn certain parameters of our variational distributions, we need to be careful not to output parameter values that would violate constraints. ENTER ADVI!

Approach: Map the constrained support space to the real line (in the dimension of the parameter ofc) and use THAT space for gradient algorithms.

Let $T : \Theta \rightarrow \mathbb{R}^D$ be a bijective mapping that maps from the constrained to the unconstrained space. Because this transformation is invertible, we can always get back our variational parameter value (u): $u = T^{-1}(\theta)$

$$p(\mathbf{u}) = p(T^{-1}(\mathbf{u})) |\det(\mathbf{J}_{T^{-1}}(\mathbf{u}))|$$

where $\mathbf{J}_{T^{-1}}$ is the Jacobian of the inverse mapping $\mathbf{u} \rightarrow \boldsymbol{\theta}$. Hence the ELBO becomes

$$\mathcal{L}(\boldsymbol{\psi}) = E_{\mathbf{u} \sim q_{\boldsymbol{\psi}}(\mathbf{u})} [\log p(\mathcal{D} | T^{-1}(\mathbf{u})) + \log p(T^{-1}(\mathbf{u})) + \log |\det(\mathbf{J}_{T^{-1}}(\mathbf{u}))|] + \mathbb{H}(\boldsymbol{\psi})$$

In the equation above, the first term is $\log p(x|z)$, the second and third terms make up $\log p(z)$ and the last term is the entropy: $E_{\mathbf{u} \sim q_{\boldsymbol{\psi}}(\mathbf{u})} (-\log q_{\boldsymbol{\psi}}(\mathbf{u}))$. Example transformations include the natural logarithm to map from positives to reals for variances.

Example transformations include the inverse hyperbolic tangent to map from probabilities $[0, 1]$ to reals (may need to use epsilon to prevent overflow errors). $\tanh^{-1}(2x - 1)$ to be exact.

7 CAVI

So far, we discussed ways to make VI work with gradient descent. However, this need not be the only way we can optimize variational parameters. Recall the mean field approximation and free-form VI.

$$q_{\phi}(z|x) = \prod_{j=1}^J q_{\psi_j}(z_j) \stackrel{\text{def}}{=} \prod_{j=1}^J q_j(z_j)$$

The overall idea is summarized in the algorithm below.

Algorithm 10.4: Coordinate ascent variational inference (CAVI).

```

1 Initialize  $q_j(z_j)$  for  $j = 1 : J$ 
2 foreach  $t = 1 : T$  do
3   foreach  $j = 1 : J$  do
4     Compute  $g_j(z_j) = \mathbb{E}_{z_{\text{mb}_j}} [\log \tilde{p}(z_i, z_{\text{mb}_i})]$ 
5     Compute  $q_j(z_j) \propto \exp(g_j(z_j))$ 
```

Figure 5: Caption

7.1 Derivation

10.3.1 Derivation of CAVI algorithm

In this section, we derive the coordinate ascent variational inference (CAVI) procedure.

To derive the update equations, we initially assume there are just 3 discrete latent variables, to simplify notation. In this case the ELBO is given by

$$\mathbb{L}(q_1, q_2, q_3) = \sum_{z_1} \sum_{z_2} \sum_{z_3} q_1(z_1) q_2(z_2) q_3(z_3) \log \tilde{p}(z_1, z_2, z_3) + \sum_{j=1}^3 \mathbb{H}(q_j) \quad (10.77)$$

where we define $\tilde{p}(z) = p_{\theta}(z, x)$ for brevity. We will optimize this wrt each q_i , one at a time, keeping the others fixed.

Let us look at the objective for q_3 :

$$\mathbb{L}_3(q_3) = \sum_{z_3} q_3(z_3) \left[\sum_{z_1} \sum_{z_2} q_1(z_1) q_2(z_2) \log \tilde{p}(z_1, z_2, z_3) \right] + \mathbb{H}(q_3) + \text{const} \quad (10.78)$$

$$= \sum_{z_3} q_3(z_3) [g_3(z_3) - \log q_3(z_3)] + \text{const} \quad (10.79)$$

where

$$g_3(z_3) \triangleq \sum_{z_1} \sum_{z_2} q_1(z_1) q_2(z_2) \log \tilde{p}(z_1, z_2, z_3) = \mathbb{E}_{\mathbf{z}_{-3}} [\log \tilde{p}(z_1, z_2, z_3)] \quad (10.80)$$

where $\mathbf{z}_{-3} = (z_1, z_2)$ is all variables except z_3 . Here $g_3(z_3)$ can be interpreted as an expected negative energy (log probability). We can convert this into an unnormalized probability distribution by defining

$$\tilde{f}_3(z_3) = \exp(g_3(z_3)) \quad (10.81)$$

which we can normalize to get

$$f_3(z_3) = \frac{\tilde{f}_3(z_3)}{\sum_{z'_3} \tilde{f}_3(z'_3)} \propto \exp(g_3(z_3)) \quad (10.82)$$

Draft of “Probabilistic Machine Learning: Advanced Topics”. April 1, 2023

Since $g_3(z_3) \propto \log f_3(z_3)$ we get

$$\mathbb{L}_3(q_3) = \sum_{z_3} q_3(z_3) [\log f_3(z_3) - \log q_3(z_3)] + \text{const} = -D_{\text{KL}}(q_3 \parallel f_3) + \text{const} \quad (10.83)$$

Since $D_{\text{KL}}(q_3 \parallel f_3)$ achieves its minimal value of 0 when $q_3(z_3) = f_3(z_3)$ for all z_3 , we see that $q_3^*(z_3) = f_3(z_3)$.

Now suppose that the joint distribution is defined by a Markov chain, where $z_1 \rightarrow z_2 \rightarrow z_3$, so $z_1 \perp z_3 | z_2$. Hence $\log \tilde{p}(z_1, z_2, z_3) = \log \tilde{p}(z_2, z_3 | z_1) + \log \tilde{p}(z_1)$, where the latter term is independent of $q_3(z_3)$. Thus the ELBO simplifies to

$$\mathbb{L}_3(q_3) = \sum_{z_3} q_3(z_3) \left[\sum_{z_2} q_2(z_2) \log \tilde{p}(z_2, z_3) \right] + \mathbb{H}(q_3) + \text{const} \quad (10.84)$$

$$= \sum_{z_3} q_3(z_3) [\log f_3(z_3) - \log q_3(z_3)] + \text{const} \quad (10.85)$$

where

$$f_3(z_3) \propto \exp \left[\sum_{z_2} q_2(z_2) \log \tilde{p}(z_2, z_3) \right] = \exp [\mathbb{E}_{\mathbf{z}_{\text{mb}_3}} [\log \tilde{p}(z_2, z_3)]] \quad (10.86)$$

8 Variational Bayes

Recall the model with shared latent variable θ .

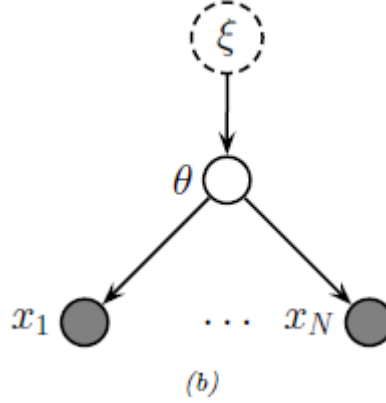


Figure 7: Caption

We discussed the possibility that ξ is fixed. In that case, it cannot be updated and it is independent of θ and can be removed from the graph.

Our new goal is to learn the posterior:

$$p(\boldsymbol{\theta} \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathcal{D})} \propto p(\boldsymbol{\theta}) \prod_{n=1}^N p(\mathcal{D}_n \mid \boldsymbol{\theta})$$

The ELBO takes on the familiar form:

$$L(\phi_{\boldsymbol{\theta}} \mid \mathcal{D}) = \mathbb{E}_{q(\boldsymbol{\theta} \mid \phi_{\boldsymbol{\theta}})} \left[p(\boldsymbol{\theta}) \prod_{n=1}^N p(\mathcal{D}_n \mid \boldsymbol{\theta}) - \sum_{j=1}^J \log(q(\theta_j \mid \phi_{\theta_j})) \right]$$

Now that the ELBO is expressed after taking the mean field approximation into account, we can leverage CAVI to learn the variational parameters!