

10 Variational inference

10.1 Introduction

In this chapter, we discuss **variational inference**, which reduces posterior inference to optimization. Note that VI is a large topic; this chapter just gives a high level overview. For more details, see e.g., [Jor+98; JJ00; Jaa01; WJ08; SQ05; TLG08; Zha+19b; Bro18].

10.1.1 The variational objective

Consider a model with unknown (latent) variables \mathbf{z} , known variables \mathbf{x} , and fixed parameters $\boldsymbol{\theta}$. (If the parameters are unknown, they can be added to \mathbf{z} , as we discuss later.) We assume the prior is $p_{\boldsymbol{\theta}}(\mathbf{z})$ and the likelihood is $p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})$, so the unnormalized joint is $p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}) = p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})p_{\boldsymbol{\theta}}(\mathbf{z})$, and the posterior is $p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}) = p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})/p_{\boldsymbol{\theta}}(\mathbf{x})$. We assume that it is intractable to compute the normalization constant, $p_{\boldsymbol{\theta}}(\mathbf{x}) = \int p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})d\mathbf{z}$, and hence intractable to compute the normalized posterior. We therefore seek an approximation to the posterior, which we denote by $q(\mathbf{z})$, such that we minimize the following loss:

$$q = \underset{q \in \mathcal{Q}}{\operatorname{argmin}} D_{\text{KL}}(q(\mathbf{z}) \parallel p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})) \quad (10.1)$$

Since we are minimizing over functions (namely distributions q), this is called a **variational method**.

In practice we pick a parametric family \mathcal{Q} , where we use $\boldsymbol{\psi}$, known as the **variational parameters**, to specify which member of the family we are using. We can compute the best variational parameters (for given \mathbf{x}) as follows:

$$\boldsymbol{\psi}^* = \underset{\boldsymbol{\psi}}{\operatorname{argmin}} D_{\text{KL}}(q_{\boldsymbol{\psi}}(\mathbf{z}) \parallel p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})) \quad (10.2)$$

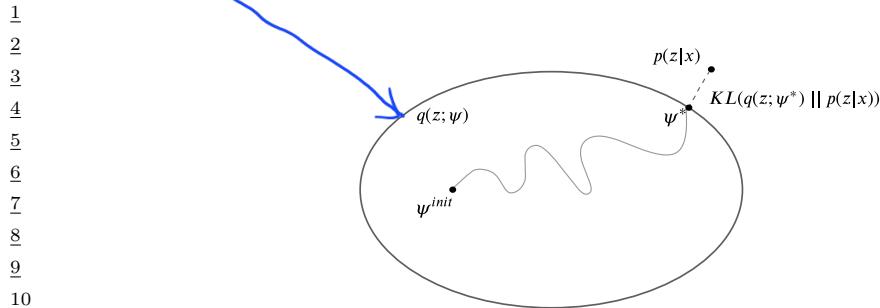
$$= \underset{\boldsymbol{\psi}}{\operatorname{argmin}} \mathbb{E}_{q_{\boldsymbol{\psi}}(\mathbf{z})} \left[\log q_{\boldsymbol{\psi}}(\mathbf{z}) - \log \left(\frac{p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})p_{\boldsymbol{\theta}}(\mathbf{z})}{p_{\boldsymbol{\theta}}(\mathbf{x})} \right) \right] \quad (10.3)$$

$$= \underset{\boldsymbol{\psi}}{\operatorname{argmin}} \underbrace{\mathbb{E}_{q_{\boldsymbol{\psi}}(\mathbf{z})} [\log q_{\boldsymbol{\psi}}(\mathbf{z}) - \log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z}) - \log p_{\boldsymbol{\theta}}(\mathbf{z})]}_{\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\psi}|\mathbf{x})} + \log p_{\boldsymbol{\theta}}(\mathbf{x}) \quad (10.4)$$

The final term $\log p_{\boldsymbol{\theta}}(\mathbf{x}) = \log(\int p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})d\mathbf{z})$ is generally intractable to compute. Fortunately, it is independent of $\boldsymbol{\psi}$, so we can drop it. This leaves us with the first term, which we write as follows:

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\psi}|\mathbf{x}) = \mathbb{E}_{q_{\boldsymbol{\psi}}(\mathbf{z})} [-\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}) + \log q_{\boldsymbol{\psi}}(\mathbf{z})] \quad (10.5)$$

set of distributions = $q(z|\psi)$, $p(z|x)$ = true distribution ψ^* best approximation to p , closest KL divergence
 434 $\psi^{in} \rightsquigarrow \psi^* \rightarrow$ optimization procedure



11
 12 Figure 10.1: Illustration of variational inference. The large oval represents the set of variational distributions
 13 $Q = \{q_\psi(z) : \psi \in \Theta\}$, where Θ is the set of possible variational parameters. The true distribution is the
 14 point $p(z|x)$, which we assume lies outside the set. Our goal is to find the best approximation to p within our
 15 variational family; this is the point ψ^* which is closest in KL divergence. We find this point by starting an
 16 optimization procedure from the random initial point ψ^{init} . Adapted from a figure by David Blei.
 17

18 Minimizing this objective will minimize the KL divergence, causing our approximation to approach
 19 the true posterior. See Figure 10.1 for an illustration. In the sections below, we give two different
 20 interpretations of this objective function.

22 10.1.1.1 The view from physics: minimize the variational free energy

23 If we define $\mathcal{E}_\theta(z) = -\log p_\theta(z, x)$ as the energy, then we can rewrite the loss in Equation (10.5)

$$24 \quad \mathcal{L}(\theta, \psi|x) = \mathbb{E}_{q_\psi(z)} [\mathcal{E}_\theta(z)] - \mathbb{H}(q_\psi) = \text{expected energy} - \text{entropy} \quad (10.6)$$

25 In physics, this is known as the **variational free energy** (VFE). This is an upper bound on the
 26 **free energy** (FE), $-\log p_\theta(x)$, which follows from the fact that

$$27 \quad D_{KL}(q_\psi(z) || p_\theta(z|x)) = \mathcal{L}(\theta, \psi|x) + \log p_\theta(x) \geq 0 \quad (10.7)$$

$$28 \quad \underbrace{\mathcal{L}(\theta, \psi|x)}_{\text{VFE}} \geq \underbrace{-\log p_\theta(x)}_{\text{FE}} \quad (10.8)$$

29 Variational inference is equivalent to minimizing the VFE. If we reach the minimum value of
 30 $-\log p_\theta(x)$, then the KL divergence term will be 0, so our approximate posterior will be exact.

31 10.1.1.2 The view from statistics: maximize the evidence lower bound (ELBO)

32 The negative of the VFE is known as the **evidence lower bound** or **ELBO** function [BKM16]:

$$33 \quad \mathcal{L}(\theta, \psi|x) \triangleq \mathbb{E}_{q_\psi(z)} [\log p_\theta(x, z) - \log q_\psi(z)] = \text{ELBO} \quad (10.9)$$

34 The name ‘‘ELBO’’ arises because

$$35 \quad \mathcal{L}(\theta, \psi|x) \leq \log p_\theta(x) \quad \text{evidence} \quad (10.10)$$

36 where $\log p_\theta(x)$ is called the ‘‘evidence’’. The inequality follows from Equation (10.8). Therefore
 37 maximizing the ELBO wrt ψ will minimize the original KL, since $\log p_\theta(x)$ is a constant wrt ψ .

38

(Note: we use the symbol \mathbb{L} for the ELBO, rather than \mathcal{L} , since we want to maximize \mathbb{L} but minimize \mathcal{L} .)

We can rewrite the ELBO as follows:

$$\mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\psi} | \mathbf{x}) = \mathbb{E}_{q_{\boldsymbol{\psi}}(\mathbf{z})} [\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})] + \mathbb{H}(q_{\boldsymbol{\psi}}(\mathbf{z})) \quad (10.11)$$

We can interpret this

$$\text{ELBO} = \text{expected log joint} + \text{entropy} \quad (10.12)$$

The second term encourages the posterior to be maximum entropy, while the first term encourages it to be a joint MAP configuration.

We can also rewrite the ELBO as

$$\mathbb{L}(\boldsymbol{\psi} | \boldsymbol{\theta}, \mathbf{x}) = \mathbb{E}_{q_{\boldsymbol{\psi}}(\mathbf{z})} [\log p_{\boldsymbol{\theta}}(\mathbf{x} | \mathbf{z}) + \log p_{\boldsymbol{\theta}}(\mathbf{z}) - \log q_{\boldsymbol{\psi}}(\mathbf{z})] \quad (10.13)$$

$$= \mathbb{E}_{q_{\boldsymbol{\psi}}(\mathbf{z})} [\log p_{\boldsymbol{\theta}}(\mathbf{x} | \mathbf{z})] - D_{\text{KL}}(q_{\boldsymbol{\psi}}(\mathbf{z}) || p_{\boldsymbol{\theta}}(\mathbf{z})) \quad (10.14)$$

We can interpret this as follows:

$$\text{ELBO} = \text{expected log likelihood} - \text{KL from posterior to prior} \quad (10.15)$$

The KL term acts like a regularizer, preventing the posterior from diverging too much from the prior.

10.1.2 Form of the variational posterior

There are two main approaches for choosing the form of the variational posterior, $q_{\boldsymbol{\psi}}(\mathbf{z} | \mathbf{x})$. In the first approach, we pick a convenient functional form, such as multivariate Gaussian, and then optimize the ELBO using gradient-based methods. This is called **fixed-form VI**, and is discussed in Section 10.2.

An alternative is to make the **mean field** assumption, namely that the posterior factorizes:

$$q_{\boldsymbol{\psi}}(\mathbf{z}) = \prod_{j=1}^J q_j(z_j) \quad (10.16)$$

where $q_j(z_j) = q_{\psi_j}(z_j)$ is the posterior over the j 'th group of variables. We don't need to specify the functional form for each q_j . Instead, the optimal distributional form can be derived by maximizing the ELBO wrt each group of variational parameters one at a time, in a coordinate ascent manner. This is therefore called **free-form VI**, and is discussed in Section 10.3.

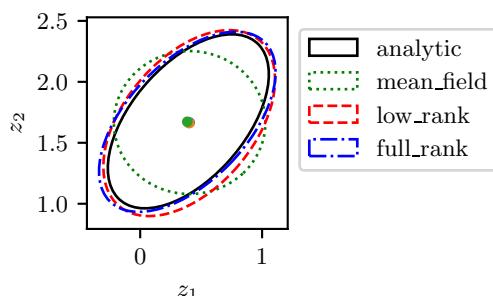
We now give a simple example of variational inference applied to a 2d latent vector \mathbf{z} , representing the mean of a Gaussian. The prior is $\mathcal{N}(\mathbf{z} | \bar{\mathbf{m}}, \bar{\mathbf{V}})$, and the likelihood is

$$p(\mathcal{D} | \mathbf{z}) = \prod_{n=1}^N \mathcal{N}(\mathbf{x}_n | \mathbf{z}, \Sigma) \propto \mathcal{N}(\bar{\mathbf{x}} | \mathbf{z}, \frac{1}{N} \Sigma) \quad (10.17)$$

The exact posterior, $p(\mathbf{z} | \mathcal{D}) = \mathcal{N}(\mathbf{z} | \hat{\mathbf{m}}, \hat{\mathbf{V}})$, can be computed analytically, as discussed in Section 3.4.4.1. In Figure 10.2, we compare three Gaussian variational approximations to the posterior. If q uses a full covariance matrix, it matches the exact posterior; however, this is intractable in high

1
2
3
4
5
6
7
8
9
10

take a look at code !!!



11 *Figure 10.2: Variational approximation to the exact (Gaussian) posterior for the mean of a 2d Gaussian*
 12 *likelihood with a Gaussian prior. We show 3 Gaussian approximations to the posterior, using a full covari-*
 13 *ance (blue), a diagonal covariance (green), and a diagonal plus rank one covariance (red). Generated by*
 14 *gaussian_2d_vi.ipynb.*

15

16

17 dimensions. If q uses a diagonal covariance matrix (corresponding to the mean field approximation),
 18 we see that the approximation is over confident, which is a well-known flaw of variational inference,
 19 due to the mode-seeking nature of minimizing $D_{\text{KL}}(q \parallel p)$ (see Section 5.1.4.3 for details). Finally, if
 20 q uses a rank-1 plus diagonal approximation, we get a much better approximation; furthermore, this
 21 can be computed quite efficiently, as we discuss in Section 10.2.1.3.
 22

23

24 10.1.3 Parameter estimation using variational EM

25

26 So far, we have assumed the model parameters θ are known. However, we can try to estimate them
 27 by maximizing the log marginal likelihood of the dataset, $\mathcal{D} = \{\mathbf{x}_n : n = 1 : N\}$,

28

$$29 \quad \log p(\mathcal{D}|\theta) = \sum_{n=1}^N \log p(\mathbf{x}_n|\theta) \quad (10.18)$$

30

31

32 In general, this is intractable to compute, but we discuss approximations below.

33

34 10.1.3.1 MLE for latent variable models

35

36 Suppose we have a latent variable model of the form

37

$$38 \quad p(\mathcal{D}, \mathbf{z}_{1:N}|\theta) = \prod_{n=1}^N p(\mathbf{z}_n|\theta)p(\mathbf{x}_n|\mathbf{z}_n, \theta) \quad (10.19)$$

39

40

41 as shown in Figure 10.3a. Furthermore, suppose we want to compute the MLE for θ given the dataset
 42 $\mathcal{D} = \{\mathbf{x}_n : n = 1 : N\}$. Since the local latent variables \mathbf{z}_n are hidden, we must marginalize them out
 43 to get the local (per example) log marginal likelihood:

44

$$45 \quad \log p(\mathbf{x}_n|\theta) = \log \left[\int p(\mathbf{x}_n|\mathbf{z}_n, \theta)p(\mathbf{z}_n|\theta)d\mathbf{z}_n \right] \quad (10.20)$$

46

47

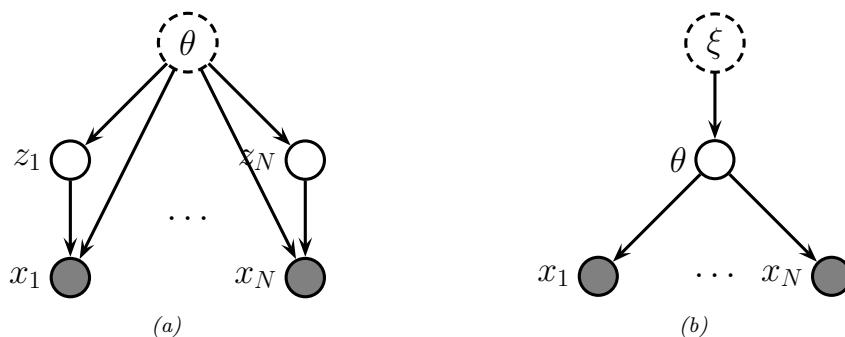


Figure 10.3: Graphical models with: (a) Local stochastic latent variables z_n and global deterministic latent parameter θ . (b) Global stochastic latent parameter θ and global deterministic latent hyper-parameter ξ . The observed variables x_n are shown by shaded circles.

Unfortunately, computing this integral is usually intractable, since it corresponds to the normalization constant of the exact posterior. Fortunately, the ELBO is a lower bound on this:

$$\mathcal{L}(\boldsymbol{\theta}, \psi_n | \mathbf{x}_n) \leq \log p(\mathbf{x}_n | \boldsymbol{\theta}) \quad (10.21)$$

We can thus optimize the model parameters by maximizing

$$\mathcal{L}(\boldsymbol{\theta}, \psi_{1:N} | \mathcal{D}) \triangleq \sum_{n=1}^N \mathcal{L}(\boldsymbol{\theta}, \psi_n | \mathbf{x}_n) \leq \log p(\mathcal{D} | \boldsymbol{\theta}) \quad (10.22)$$

This is the basis of the **variational EM** algorithm. We discuss this in more detail in Section 6.5.6.1, but the basic idea is to alternate between maximizing the ELBO wrt the variational parameters $\{\psi_n\}$ in the E step, to give us $q_{\psi_n}(z_n)$, and then maximizing the ELBO (using the new ψ_n) wrt the model parameters $\boldsymbol{\theta}$ in the M step. (We can also use SGD and amortized inference to speed this up, as we explain in Sections 10.1.4 to 10.1.5.)

10.1.3.2 Empirical Bayes for fully observed models

Suppose we have a fully observed model (with no local latent variables) of the form

$$p(\mathcal{D}, \boldsymbol{\theta} | \boldsymbol{\xi}) = p(\boldsymbol{\theta} | \boldsymbol{\xi}) \prod_{n=1}^N p(\mathbf{x}_n | \boldsymbol{\theta}) \quad (10.23)$$

as shown in Figure 10.3b. In the context of Bayesian parameter inference, our goal is to compute the parameter posterior:

$$p(\boldsymbol{\theta} | \mathcal{D}, \boldsymbol{\xi}) = \frac{p(\mathcal{D} | \boldsymbol{\theta}) p(\boldsymbol{\theta} | \boldsymbol{\xi})}{p(\mathcal{D} | \boldsymbol{\xi})} \quad \begin{matrix} \text{θ - latent variable} \\ \text{ξ - hyper parameters for prior} \end{matrix} \quad (10.24)$$

where $\boldsymbol{\theta}$ are the global unknown model parameters (latent variables), and $\boldsymbol{\xi}$ are the hyper-parameters for the prior. If the hyper-parameters are unknown, we can estimate them using empirical Bayes (see

1 Section 3.7) by computing
2

3
4 $\hat{\xi} = \underset{\xi}{\operatorname{argmax}} \log p(\mathcal{D}|\xi)$ (10.25)
5

6 We can use variational EM to compute this, similar to Section 10.1.3.1, except now the parameters
7 to be estimated are ξ , the latent variables are the shared global parameters θ , and the observations
8 are the entire dataset, \mathcal{D} . We then get the lower bound
9

10
11 $\log p(\mathcal{D}|\xi) \geq \mathbb{L}(\xi, \psi|\mathcal{D}) = \mathbb{E}_{q_\psi(\theta)} \left[\sum_{n=1}^N \log p(\mathbf{x}_n|\theta) \right] - D_{\text{KL}}(q_\psi(\theta) \parallel p(\theta|\xi))$ (10.26)
12

13 We optimize this wrt the parameters of the variational posterior, ψ , and wrt the prior hyper-parameters
14 ξ .

15 If the prior ξ is fixed, we just need to optimize the variational parameters ψ to compute the
16 posterior, $q_\psi(\theta|\mathcal{D})$. This is known as **variational Bayes**. See Section 10.3.3 for more details.
17

18 10.1.4 Stochastic VI

19 In Section 10.1.3, we saw that parameter estimation requires optimizing the ELBO for the entire
20 dataset, which is defined as the sum of the ELBOs for each of the N data samples \mathbf{x}_n . Computing
21 this objective can be slow if N is large. Fortunately, we can replace this objective with a stochastic
22 approximation, which is faster to compute, and provides an unbiased estimate. In particular, at each
23 step, we can draw a random minibatch of $B = |\mathcal{B}|$ examples from the dataset, and then make the
24 approximation
25

26
27
28 $\mathbb{L}(\theta, \psi_{1:N}|\mathcal{D}) = \sum_{n=1}^N \mathbb{L}(\theta, \psi_n|\mathbf{x}_n) \approx \frac{N}{B} \sum_{\mathbf{x}_n \in \mathcal{B}} \left[\mathbb{E}_{q_{\psi_n}(\mathbf{z}_n)} [\log p_\theta(\mathbf{x}_n|\mathbf{z}_n) + \log p_\theta(\mathbf{z}_n) - \log q_{\psi_n}(\mathbf{z}_n)] \right]$
29
30 (10.27)
31

32 This can be used inside of a stochastic optimization algorithm, such as SGD. This is called **stochastic**
33 **variational inference** or **SVI** [Hof+13], and allows VI to scale to large datasets.
34

35 10.1.5 Amortized VI

problem with SVI

36 In Section 10.1.4, we saw that in each iteration of SVI, we need to optimize the local variational
37 parameters ψ_n for each example n in the minibatch. This nested optimization can be quite slow.

38 An alternative approach is to train a model, known as an **inference network** or **recognition**
39 **network**, to predict ψ_n from the observed data, \mathbf{x}_n , using $\psi_n = f_\phi^{\text{inf}}(\mathbf{x}_n)$. This technique is known
40 as **amortized variational inference** [GG14], or **inference compilation** [LBW17], since we are
41 reducing the cost of per-example time inference by training a model that is shared across all examples.
42 (See also [Amo22] for a general discussion of amortized optimization.) For brevity, we will write the
43 amortized posterior as
44

45
46 $q(\mathbf{z}_n|\psi_n) = q(\mathbf{z}_n|f_\phi^{\text{inf}}(\mathbf{x}_n)) = q_\phi(\mathbf{z}_n|\mathbf{x}_n)$ (10.28)
47

The corresponding ELBO becomes

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi} | \mathcal{D}) = \sum_{n=1}^N [\mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_n | \mathbf{x}_n)} [\log p_{\boldsymbol{\theta}}(\mathbf{x}_n, \mathbf{z}_n) - \log q_{\boldsymbol{\phi}}(\mathbf{z} | \mathbf{x}_n)]] \quad (10.29)$$

If we combine this with SVI we get an amortized version of Equation (10.27). For example, if we use a minibatch of size 1, we get

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi} | \mathbf{x}_n) \approx N [\mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_n | \mathbf{x}_n)} [\log p_{\boldsymbol{\theta}}(\mathbf{x}_n, \mathbf{z}_n) - \log q_{\boldsymbol{\phi}}(\mathbf{z} | \mathbf{x}_n)]] \quad (10.30)$$

We can optimize this as shown in Algorithm 10.1. Note that the (partial) maximization wrt $\boldsymbol{\theta}$ in the M step is usually done with a gradient update, but the maximization wrt $\boldsymbol{\phi}$ in the E step is trickier, since the loss uses $\boldsymbol{\phi}$ to define an expectation operator, so we can't necessarily push the gradient operator inside; we discuss ways to optimize the variational parameters in Section 10.2 and Section 10.3.

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:  $\boldsymbol{\phi} = \operatorname{argmax}_{\boldsymbol{\phi}} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi} | \mathbf{x}_n)$ 
5   M step:  $\boldsymbol{\theta} = \operatorname{argmax}_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi} | \mathbf{x}_n)$ 
6 until converged

```

solution semi-amortized inference

10.1.6 Semi-amortized inference

problem with amortized VI

Amortized SVI is widely used for fitting LVMs, e.g., for VAEs (see Section 21.2), for topic models [SS17a], for probabilistic programming [RHG16], for CRFs [TG18], etc. However, the use of an inference network can result in a suboptimal setting of the local variational parameters ψ_n . This is called the **amortization gap** [CLD18]. We can close this gap by using the inference network to warm-start an optimizer for ψ_n ; this is known as **semi-amortized VI** [Kim+18c]. (See also [MYM18], who propose a closely related method called **iterative amortized inference**.)

An alternative approach is to use the inference network as a proposal distribution. If we combine this with importance sampling, we get the IWAE bound of Section 10.5.1. If we use this with Metropolis-Hastings, we get a VI-MCMC hybrid (see Section 10.4.5).

10.2 Gradient-based VI

In this section, we will choose some convenient form for $q_{\boldsymbol{\phi}}(\mathbf{z})$, such as a Gaussian for continuous \mathbf{z} , or a product of categoricals for discrete \mathbf{z} , and then optimize the ELBO using gradient based methods.

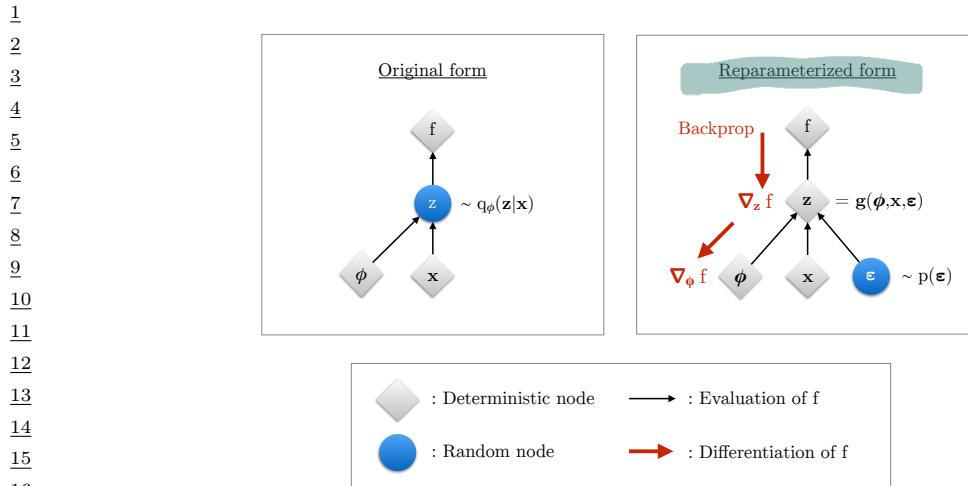


Figure 10.4: Illustration of the reparameterization trick. The objective f depends on the variational parameters ϕ , the observed data x , and the latent random variable $z \sim q_\phi(z|x)$. On the left, we show the standard form of the computation graph. On the right, we show a reparameterized form, in which we move the stochasticity into the noise source ϵ , and compute z deterministically, $z = g(\phi, x, \epsilon)$. The rest of the graph is deterministic, so we can backpropagate the gradient of the scalar f wrt ϕ through z and into ϕ . From Figure 2.3 of [KW19a]. Used with kind permission of Durk Kingma.

The gradient wrt the generative parameters θ is easy to compute, since we can push gradients inside the expectation, and use a single Monte Carlo sample:

$$\nabla_\theta \mathcal{L}(\theta, \phi | x) = \nabla_\theta \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x, z) - \log q_\phi(z|x)] \quad (10.31)$$

$$= \mathbb{E}_{q_\phi(z|x)} [\nabla_\theta \{\log p_\theta(x, z) - \log q_\phi(z|x)\}] \quad (10.32)$$

$$\approx \nabla_\theta \log p_\theta(x, z^s) \quad (10.33)$$

where $z^s \sim q_\phi(z|x)$. This is an unbiased estimate of the gradient, so can be used with SGD.

The gradient wrt the inference parameters ϕ is harder to compute since

$$\nabla_\phi \mathcal{L}(\theta, \phi | x) = \nabla_\phi \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x, z) - \log q_\phi(z|x)] \quad (10.34)$$

$$\neq \mathbb{E}_{q_\phi(z|x)} [\nabla_\phi \{\log p_\theta(x, z) - \log q_\phi(z|x)\}] \quad (10.35)$$

However, we can often use the reparameterization trick, which we discuss in Section 10.2.1. If not, we can use blackbox VI, which we discuss in Section 10.2.3.

10.2.1 Reparameterized VI

In this section, we discuss the **reparameterization trick** for taking gradients wrt distributions over continuous latent variables $z \sim q_\phi(z|x)$. We explain this in detail in Section 6.3.5, but we summarize the basic idea here.

The key trick is to rewrite the random variable $\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})$ as some differentiable (and invertible) transformation g of another random variable $\boldsymbol{\epsilon} \sim p(\boldsymbol{\epsilon})$, which does not depend on ϕ , i.e., we assume we can write

$$\mathbf{z} = g(\phi, \mathbf{x}, \boldsymbol{\epsilon}) \quad (10.36)$$

For example,

$$\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma})) \iff \mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\epsilon} \odot \boldsymbol{\sigma}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (10.37)$$

Using this, we have

$$\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[f(\mathbf{z})] = \mathbb{E}_{p(\boldsymbol{\epsilon})}[f(\mathbf{z})] \quad \text{s.t. } \mathbf{z} = g(\phi, \mathbf{x}, \boldsymbol{\epsilon}) \quad (10.38)$$

where we define

$$f_{\theta, \phi}(\mathbf{z}) = \log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x}) \quad (10.39)$$

Hence

$$\nabla_\phi \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[f(\mathbf{z})] = \nabla_\phi \mathbb{E}_{p(\boldsymbol{\epsilon})}[f(\mathbf{z})] = \mathbb{E}_{p(\boldsymbol{\epsilon})}[\nabla_\phi f(\mathbf{z})] \quad (10.40)$$

which we can approximate with a single Monte Carlo sample. This lets us propagate gradients back through the f function. See Figure 10.4 for an illustration. This is called **reparameterized VI** or **RVI**.

Since we are now working with the random variable $\boldsymbol{\epsilon}$, we need to use the change of variables formula to compute

$$\log q_\phi(\mathbf{z}|\mathbf{x}) = \log p(\boldsymbol{\epsilon}) - \log \left| \det \left(\frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}} \right) \right| \quad (10.41)$$

where $\frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}}$ is the Jacobian:

$$\frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}} = \begin{pmatrix} \frac{\partial z_1}{\partial \epsilon_1} & \dots & \frac{\partial z_1}{\partial \epsilon_k} \\ \vdots & \ddots & \vdots \\ \frac{\partial z_k}{\partial \epsilon_1} & \dots & \frac{\partial z_k}{\partial \epsilon_k} \end{pmatrix} \quad (10.42)$$

We design the transformation $\mathbf{z} = g(\boldsymbol{\epsilon})$ such that this Jacobian is tractable to compute. We give some examples below.

10.2.1.1 Gaussian with diagonal covariance (mean field)

Suppose we use a fully factorized Gaussian posterior. Then the reparameterization process becomes

$$\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (10.43)$$

$$\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon} \quad (10.44)$$

1 where the inference network generates the parameters of the transformation:
2

3 $(\boldsymbol{\mu}, \log \boldsymbol{\sigma}) = f_{\phi}^{\text{inf}}(\mathbf{x})$ (10.45)
4

5 Thus to sample from the posterior $q_{\phi}(\mathbf{z}|\mathbf{x})$, we sample $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, and then compute \mathbf{z} .
6

7 Given the sample, we need to evaluate the ELBO:

8 $f(\mathbf{z}) = \log p_{\theta}(\mathbf{x}|\mathbf{z}) + \log p_{\theta}(\mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})$ (10.46)
9

10 To evaluate the $p_{\theta}(\mathbf{x}|\mathbf{z})$ term, we can just plug \mathbf{z} into the likelihood. To evaluate the $\log q_{\phi}(\mathbf{z}|\mathbf{x})$
11 term, we need to use the change of variables formula from Equation (10.41). The Jacobian is given
12 by $\frac{\partial \mathbf{z}}{\partial \epsilon} = \text{diag}(\boldsymbol{\sigma})$. Hence
13

14

15 $\log q_{\phi}(\mathbf{z}|\mathbf{x}) = \sum_{k=1}^K [\log \mathcal{N}(\epsilon_k|0, 1) - \log \sigma_k] = - \sum_{k=1}^K \left[\frac{1}{2} \log(2\pi) + \frac{1}{2} \epsilon_k^2 + \log \sigma_k \right]$ (10.47)
16

17 Finally, to evaluate the $p(\mathbf{z})$ term, we can use the transformation $\mathbf{z} = \mathbf{0} + \mathbf{1} \odot \epsilon$, so the Jacobian is
18 the identity and we get

19

20 $\log p(\mathbf{z}) = \sum_{k=1}^K \left[\frac{1}{2} z_k^2 + \frac{1}{2} \log(2\pi) \right]$ (10.48)
21

22 An alternative is to use the objective
23

24 $f'(\mathbf{z}) = \log p_{\theta}(\mathbf{x}|\mathbf{z}) + D_{\text{KL}}(q_{\phi}(\mathbf{Z}|\mathbf{x}) \parallel p_{\theta}(\mathbf{Z}))$ (10.49)
25

26 In some cases, we evaluate the second term analytically, without needing Monte Carlo. For example,
27 if we assume a diagonal Gaussian prior, $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$, and diagonal gaussian posterior, $q(\mathbf{z}|\mathbf{x}) =$
28 $\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}))$, we can use Equation (5.78) to compute the KL in closed form:
29

30

31 $D_{\text{KL}}(q \parallel p) = -\frac{1}{2} \sum_{k=1}^K [\log \sigma_k^2 - \sigma_k^2 - \mu_k^2 + 1]$ (10.50)
32

33 The objective $f'(\mathbf{z})$ is often lower variance than $f(\mathbf{z})$, since it computes the KL analytically. However,
34 it is harder to generalize this objective to settings where the prior and/or posterior are not Gaussian.
35

36 10.2.1.2 Gaussian with full covariance

37 Now consider using a full covariance Gaussian posterior. We will compute a Cholesky decomposition
38 of the covariance, $\boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}^T$, where \mathbf{L} is a lower triangular matrix with non-zero entries on the
39 diagonal. Hence the reparameterization becomes
40

41 $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ (10.51)
42

43 $\mathbf{z} = \boldsymbol{\mu} + \mathbf{L}\epsilon$ (10.52)
44

45

The Jacobian of this affine transformation is $\frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}} = \mathbf{L}$. Since \mathbf{L} is a triangular matrix, its determinant is the product of its main diagonal, so

$$\log \left| \det \frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}} \right| = \sum_{k=1}^K \log |L_{kk}| \quad (10.53)$$

We can compute \mathbf{L} using

$$\mathbf{L} = \mathbf{M} \odot \mathbf{L}' + \text{diag}(\boldsymbol{\sigma}) \quad (10.54)$$

where \mathbf{M} is a masking matrix with 0s on and above the diagonal, and 1s below the diagonal, and where $(\boldsymbol{\mu}, \log \boldsymbol{\sigma}, \mathbf{L}')$ is predicted by the inference network. With this construction, the diagonal entries of \mathbf{L} are given by $\boldsymbol{\sigma}$, so

$$\log \left| \det \frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}} \right| = \sum_{k=1}^K \log |L_{kk}| = \sum_{k=1}^K \log \sigma_k \quad (10.55)$$

10.2.1.3 Gaussian with low-rank plus diagonal covariance

In high dimensions, an efficient alternative to using a Cholesky decomposition is the factor decomposition

$$\boldsymbol{\Sigma} = \mathbf{B}\mathbf{B}^\top + \mathbf{C}^2 \quad (10.56)$$

where \mathbf{B} is the factor loading matrix of size $d \times f$, where $f \ll d$ is the number of factors, d is the dimensionality of \mathbf{z} , and $\mathbf{C} = \text{diag}(c_1, \dots, c_d)$. This reduces the total number of variational parameters from $d + d(d+1)/2$ to $(f+2)d$. In [ONS18], they called this approach **VAF**C (short for variational approximation with factor covariance)

In the special case where $J = 1$, the covariance matrix becomes

$$\boldsymbol{\Sigma} = \mathbf{b}\mathbf{b}^\top + \text{diag}(\mathbf{c}^2) \quad (10.57)$$

In this case, it is possible to compute the natural gradient (Section 6.4) of the ELBO in closed form in $O(d)$ time, as shown in [Tra+20b; TND21], who call the approach **NAGVAC-1** (natural gradient Gaussian variational approximation). This can result in much faster convergence than following the normal gradient.

In Section 10.1.2, we show that this low rank approximation is much better than a diagonal approximation. See [Supplementary](#) Section 10.1 for more examples.

10.2.1.4 Other variational posteriors

Many other kinds of distribution can be written in a reparameterizable way, as described in [Moh+20]. This includes standard exponential family distributions, such as the gamma and Dirichlet, as well as more exotic forms, such as inverse autoregressive flows (see Section 10.4.3).

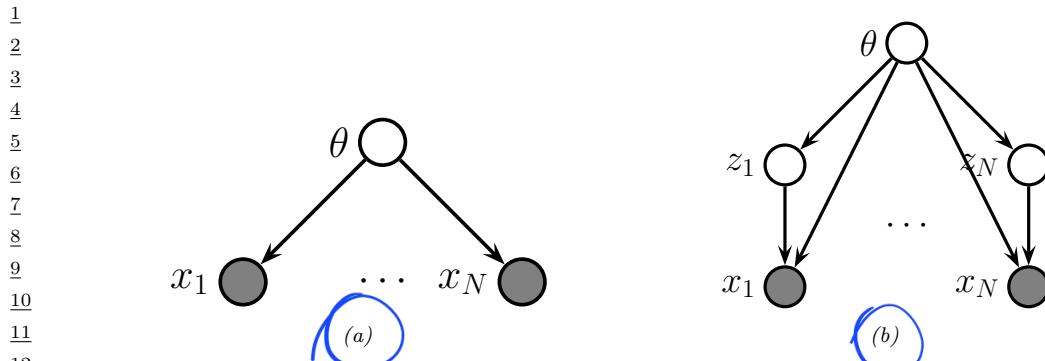


Figure 10.5: Graphical models with (a) Global latent parameter θ and observed variables $\mathbf{x}_{1:N}$. (b) Local latent variables $\mathbf{z}_{1:N}$, global latent parameter θ , and observed variables $\mathbf{x}_{1:N}$.

Explanation how to transform (a) into (b)

10.2.1.5 Example: Bayesian parameter inference

In this section, we use reparameterized SVI to infer the posterior for the parameters of a Gaussian mixture model (GMM). We will marginalize out the discrete latent variables, so just need to approximate the posterior over the global latent, $p(\boldsymbol{\theta}|\mathcal{D})$. This is sometimes called a “collapsed” model, since we have marginalized out all the local latent variables. That is, we have converted the model in Figure 10.5b to the one in Figure 10.5a. We choose a factored (mean field) variational posterior that is conjugate to the likelihood, but is also reparameterizable. This lets us fit the posterior with SGD.

For simplicity, we assume diagonal covariance matrices for each Gaussian mixture component. Thus the likelihood for one datapoint, $\mathbf{x} \in \mathbb{R}^D$, is

$$p(\mathbf{x}|\boldsymbol{\theta}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \text{diag}(\boldsymbol{\lambda}_k)^{-1}) \quad (10.58)$$

where $\boldsymbol{\mu}_k = (\mu_{k1}, \dots, \mu_{kD})$ are the means, $\boldsymbol{\lambda}_k = (\lambda_{k1}, \dots, \lambda_{kD})$ are the precisions, and $\boldsymbol{\pi} = (\pi_1, \dots, \pi_K)$ are the mixing weights. We use the following prior for these parameters:

$$p_{\xi}(\boldsymbol{\theta}) = \left[\prod_{k=1}^K \prod_{d=1}^D \mathcal{N}(\mu_{kd}|0, 1) \text{Ga}(\lambda_{kd}|5, 5) \right] \text{Dir}(\boldsymbol{\pi}|\mathbf{1}) \quad (10.59)$$

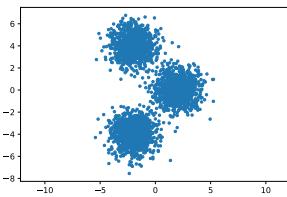
where ξ are the hyperparameters. We assume the following mean field posterior:

$$q_{\psi}(\boldsymbol{\theta}) = \left[\prod_{k=1}^K \prod_{d=1}^D \mathcal{N}(\mu_{kd}|m_{kd}, s_{kd}) \text{Ga}(\lambda_{kd}|\alpha_{kd}, \beta_{kd}) \right] \text{Dir}(\boldsymbol{\pi}|\mathbf{c}) \quad (10.60)$$

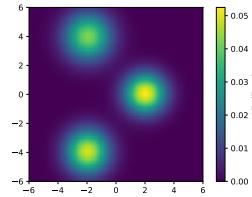
where $\psi = (\mathbf{m}_{1:K, 1:D}, \mathbf{s}_{1:K, 1:D}, \boldsymbol{\alpha}_{1:K, 1:D}, \boldsymbol{\beta}_{1:K, 1:D}, \mathbf{c})$ are the variational parameters for $\boldsymbol{\theta}$.

We can compute the ELBO using

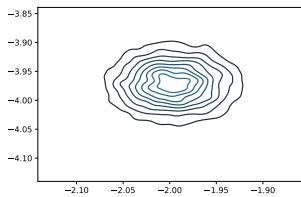
$$\mathcal{L}(\xi, \psi | \mathcal{D}) = \mathbb{E}_{q_{\psi}(\boldsymbol{\theta})} [\log p(\mathcal{D}|\boldsymbol{\theta}) + \log p_{\xi}(\boldsymbol{\theta}) - \log q_{\psi}(\boldsymbol{\theta})] \quad (10.61)$$



(a)



(b)



(c)

Figure 10.6: SVI for fitting a mixture of 3 Gaussians in 2d. (a) 3000 training points. (b) Fitted density, plugging in the posterior mean parameters. (c) Kernel density estimate fit to 10,000 samples from $q(\mu_1|\psi)$. Generated by [svi_gmm_demo_2d.ipynb](#).

Since the distributions are reparameterizable, we can and push gradients inside this expression. We can approximate the expectation by drawing a single posterior sample, and can approximate the log likelihood using minibatching. We can then update the variational parameters, (and optionally the hyperparameters of the prior, as we discussed in Section 10.1.3.2) using the pseudocode in Algorithm 10.2.

Algorithm 10.2: Reparameterized SVI for Bayesian parameter inference

```

1 Initialize  $\psi, \xi$ 
2 repeat
3   Sample minibatch  $\mathcal{B} = \{\mathbf{x}_b \sim \mathcal{D} : b = 1 : B\}$ 
4   Sample  $\epsilon \sim q_0$ 
5   Compute  $\tilde{\theta} = g(\psi, \epsilon)$ 
6   Compute  $\mathcal{L}(\psi|\mathcal{D}, \tilde{\theta}) = -\frac{N}{B} \sum_{\mathbf{x}_n \in \mathcal{B}} \log p(\mathbf{x}_n|\tilde{\theta}) - \log p_{\xi}(\tilde{\theta}) + \log q_{\psi}(\tilde{\theta})$ 
7   Update  $\xi := \xi - \eta \nabla_{\xi} \mathcal{L}(\xi, \psi|\mathcal{D}, \tilde{\theta})$ 
8   Update  $\psi := \psi - \eta \nabla_{\psi} \mathcal{L}(\xi, \psi|\mathcal{D}, \tilde{\theta})$ 
9 until converged

```

This example take a look and put into practice

Figure 10.6 gives an example of this in practice. We generate a dataset from a mixture of 3 Gaussians in 2d, using $\mu_1^* = [2, 0]$, $\mu_2^* = [-2, -4]$, $\mu_3^* = [-2, 4]$, precisions $\lambda_{dk}^* = 1$, and uniform mixing weights, $\pi^* = [1/3, 1/3, 1/3]$. Figure 10.6a shows the training set of 3000 points. We fit this using SVI, with a batch size of 500, for 1000 epochs, using the Adam optimizer. Figure 10.6b shows the predictions of the fitted model. More precisely, it shows $p(\mathbf{x}|\bar{\theta})$, where $\bar{\theta} = \mathbb{E}_{q(\theta|\psi)}[\theta]$. Figure 10.6c shows a kernel density estimate fit to 10,000 samples from $q(\mu_1|\psi)$. We see that the posterior mean is $\mathbb{E}[\mu_1] \approx [-2, -4]$. Due to label switching unidentifiability, we see this matches μ_2^* rather than μ_1^* .

1 **10.2.1.6 Example: MLE for LVMs**

3 In this section, we consider reparameterized SVI for computing the MLE for latent variable models
4 (LVMs) with continuous latents, such as variational autoencoders (Section 21.2). Unlike Sec-
5 tion 10.2.1.5, we cannot analytically marginalize out the local latents. Instead we will use amortized
6 inference, as in Section 10.1.5, which means we learn an inference network (with parameters ϕ) to
7 predict the local variational parameters ψ_n given input x_n . If we sample a single example x_n from
8 the dataset at each iteration, and a single latent variable z_n from the variational posterior, then we
9 get the pseudocode in Algorithm 10.3.

10

11 **Algorithm 10.3:** Reparameterized amortized SVI for MLE of an LVM

12 1 Initialize θ, ϕ
13 2 repeat
14 3 Sample $x_n \sim p_{\mathcal{D}}$
15 4 Sample $\epsilon_n \sim q_0$
16 5 Compute $z_n = g(\phi, x_n, \epsilon_n)$
17 6 Compute $\mathcal{L}(\theta, \phi | x_n, z_n) = -\log p_{\theta}(x_n, z_n) + \log q_{\phi}(z_n | x_n)$
18 7 Update $\theta := \theta - \eta \nabla_{\theta} \mathcal{L}(\phi, \theta | x_n, z_n)$
19 8 Update $\phi := \phi - \eta \nabla_{\phi} \mathcal{L}(\phi, \theta | x_n, z_n)$
20 9 until converged

21

22

23

24 **10.2.2 Automatic differentiation VI**

25 To apply Gaussian VI, we need to transform constrained parameters (such as variance terms) to
26 unconstrained form, so they live in \mathbb{R}^D . This technique can be used for any distribution for which
27 we can define a bijection to \mathbb{R}^D . This approach is called **automatic differentiation variational**
28 **inference** or **ADVI** [Kuc+16]. We give the details below.

29

30 **10.2.2.1 Basic idea**

31 **Explanation of ADVI**

32 Our goal is to approximate the posterior $p(\theta | \mathcal{D}) \propto p(\theta)p(\mathcal{D} | \theta)$, where $\theta \in \Theta$ lives in some D -
33 dimensional constrained parameter space. Let $T : \Theta \rightarrow \mathbb{R}^D$ be a bijective mapping that maps from
34 the constrained space Θ to the unconstrained space \mathbb{R}^D , with inverse $T^{-1} : \mathbb{R}^D \rightarrow \Theta$. Let $u = T(\theta)$
35 be the unconstrained latent variables. We will use a Gaussian variational approximation to the
36 posterior for u , i.e.,: $q_{\psi}(u) = \mathcal{N}(u | \mu_d, \Sigma)$, where $\psi = (\mu, \Sigma)$.

37 By the change of variable formula Equation (2.257), we have

38
$$p(u) = p(T^{-1}(u)) |\det(\mathbf{J}_{T^{-1}}(u))| \quad (10.62)$$

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

40
$$\mathcal{L}(\psi) = E_{u \sim q_{\psi}(u)} [\log p(\mathcal{D} | T^{-1}(u)) + \log p(T^{-1}(u)) + \log |\det(\mathbf{J}_{T^{-1}}(u))|] + \mathbb{H}(\psi) \quad (10.63)$$

41 This is a tractable objective, assuming the Jacobian is tractable, since the final entropy term is
42 available in closed form, and we can use a Monte Carlo approximation of the expectation over u .

43

Since the objective is stochastic, and reparameterizable, we can use SGD to optimize it. However, [Ing20] propose **deterministic ADVI**, in which the samples $\epsilon_s \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ are held fixed during the optimization process. This is called the common random numbers trick (Section 11.6.1), and makes the objective a deterministic function; this allows for the use of more powerful second-order optimization methods, such as BFGS. (Of course, if the dataset is large, we might need to use minibatch subsampling, which reintroduces stochasticity.)

10.2.2.2 Example: ADVI for beta-binomial model

To illustrate ADVI, we consider the 1d beta-binomial model from Section 7.4.4. We want to approximate $p(\theta|\mathcal{D})$ using the prior $p(\theta) = \text{Beta}(\theta|a, b)$ and likelihood $p(\mathcal{D}|\theta) = \prod_i \text{Ber}(y_i|\theta)$, where the sufficient statistics are $N_1 = 10$, $N_0 = 1$, and the prior is uninformative, $a = b = 1$. We use the transformation $\theta = T^{-1}(z) = \sigma(z)$, and optimize the ELBO with SGD. The results of this method are shown in Figure 7.4 and show that the Gaussian fit is a good approximation, despite the skewed nature of the posterior.

10.2.2.3 Example: ADVI for GMMs

ADVI for Gaussian Mixture Models

In this section, we use ADVI to approximate the posterior of the parameters of a mixture of Gaussians. The difference from the VBEM algorithm of Section 10.3.6 is that we use ADVI combined with a Gaussian variational posterior, rather than using a mean field approximation defined by a product of conjugate distributions.

To apply ADVI, we marginalize out the discrete local discrete latents $m_n \in \{1, \dots, K\}$ analytically, so the likelihood has the form

$$p(\mathcal{D}|\boldsymbol{\theta}) = \prod_{n=1}^N \left[\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{y}_n | \boldsymbol{\mu}_k, \text{diag}(\boldsymbol{\Sigma}_k)) \right] \quad (10.64)$$

We use an uninformative Gaussian prior for the $\boldsymbol{\mu}_k$, a uniform LKJ prior for the \mathbf{L}_k , a log-normal prior for the $\boldsymbol{\sigma}_k$, and a uniform Dirichlet prior for the mixing weights $\boldsymbol{\pi}$. (See [Kuc+16, Fig 21] for a definition of the model in STAN syntax.) The posterior approximation for the unconstrained parameters is a block-diagonal gaussian. $q(\mathbf{u}) = \mathcal{N}(\mathbf{u} | \boldsymbol{\psi}_{\boldsymbol{\mu}}, \boldsymbol{\psi}_{\boldsymbol{\Sigma}})$, where the unconstrained parameters are computed using suitable bijections (see code for details).

We apply this method to the Old Faithful dataset from Figure 10.12, using $K = 10$ mixture components. The results are shown in Figure 10.7. In the top left, we show the special case where we constrain the posterior to be a MAP estimate, by setting $\boldsymbol{\psi}_{\boldsymbol{\Sigma}} = \mathbf{0}$. We see that there is no sparsity in the posterior, since there is no Bayesian “Occam factor” from marginalizing out the parameters. In panels c-d, we show 3 samples from the posterior. We see that the Bayesian method strongly prefers just 2 mixture components, although there is a small amount of support for some other Gaussian components (shown by the faint ellipses).

10.2.2.4 More complex posteriors

We can combine ADVI with any of the improved posterior approximations that we discuss in Section 10.4 — such as Gaussian mixtures [Mor+21b] or normalizing flows [ASD20] — to create a high-quality, automatic approximate inference scheme.

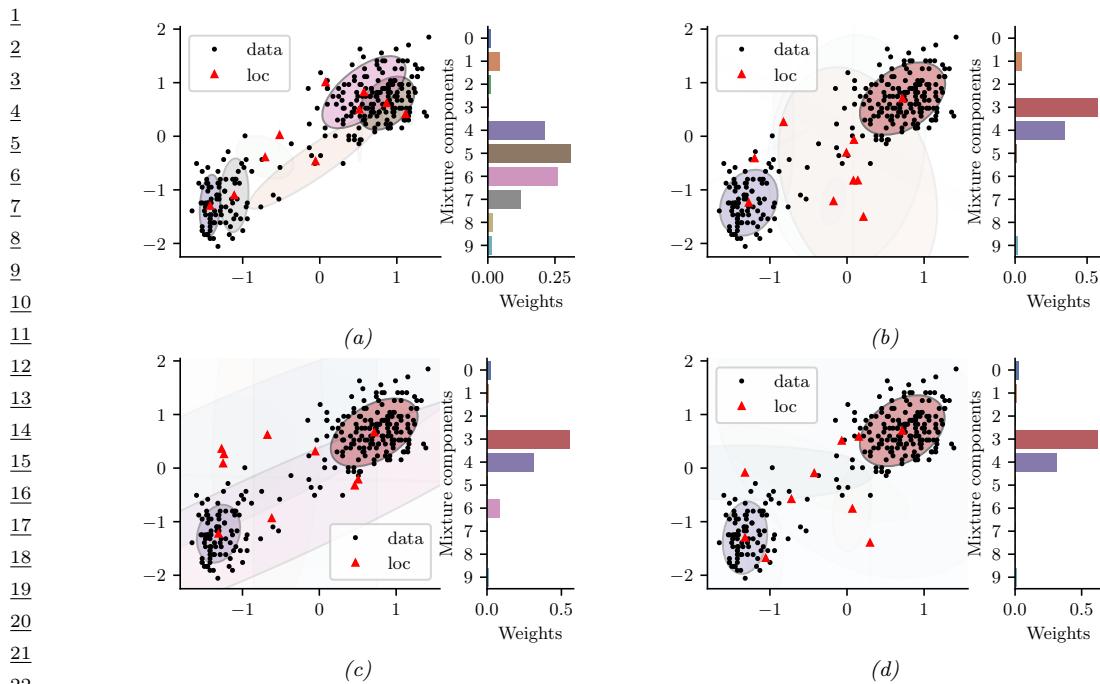


Figure 10.7: Posterior over the mixing weights (histogram) and the means and covariances of each Gaussian mixture component, using $K = 10$, when fitting the model to the Old Faithful dataset from Figure 10.12. (a) MAP approximation. (b-d) 3 samples from the Gaussian approximation. The intensity of the shading is proportional to the mixture weight. Generated by `gmm_advi_bijax.ipynb`.

10.2.3 Blackbox variational inference

In this section, we assume that we can evaluate $\tilde{\mathcal{L}}(\psi, \mathbf{z}) = \log p(\mathbf{z}, \mathbf{x}) - \log q_\psi(\mathbf{z})$ pointwise, but we do not assume we can take gradients of this function. (For example, \mathbf{z} may contain discrete variables.) We are thus treating the model as a “blackbox”. Hence this approach is called **blackbox variational inference** or **BBVI** [RGB14; ASD20].

Read more about this guy.

10.2.3.1 Estimating the gradient using REINFORCE

To estimate the gradient of the ELBO, we will use the **score function estimator**, also called the **REINFORCE estimator** (Section 6.3.4). In particular, suppose we write the ELBO as

$$\mathcal{L}(\psi) = \mathbb{E}_{q(\mathbf{z}|\psi)} [\tilde{\mathcal{L}}(\psi, \mathbf{z})] = \mathbb{E}_{q(\mathbf{z}|\psi)} [\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}|\psi)] \quad (10.65)$$

Then from Equation (6.58) we have

$$\nabla_\psi \mathcal{L}(\psi) = \mathbb{E}_{q(\mathbf{z}|\psi)} [\tilde{\mathcal{L}}(\psi, \mathbf{z}) \nabla_\psi \log q(\mathbf{z}|\psi)] \quad (10.66)$$

1 We can then compute a Monte Carlo approximation to this:

$$\widehat{\nabla_{\psi} \mathbb{L}(\psi_t)} = \frac{1}{S} \sum_{s=1}^S \tilde{\mathcal{L}}(\psi, \mathbf{z}_s) \nabla_{\psi} \log q_{\psi}(\mathbf{z}_s) |_{\psi=\psi_t} \quad (10.67)$$

6 We can pass this to any kind of gradient optimizer, such as SGD or Adam.

8 10.2.3.2 Reducing the variance using control variates

10 In practice, the variance of this estimator is quite large, so it is important to use methods such
11 as **control variates** or **CV** (Section 6.3.4.1). To see how this works, consider the naive gradient
12 estimator in Equation (10.67), which for the i 'th component we can write as

$$\widehat{\nabla_{\psi_i} \mathbb{L}(\psi_t)}^{\text{naive}} = \frac{1}{S} \sum_{s=1}^S \tilde{g}_i(\mathbf{z}_s) \quad (10.68)$$

$$\tilde{g}_i(\mathbf{z}_s) = g_i(\mathbf{z}_s) \times \tilde{\mathcal{L}}(\psi, \mathbf{z}_s) \quad (10.69)$$

$$g_i(\mathbf{z}_s) = \nabla_{\psi_i} \log q_{\psi}(\mathbf{z}_s) \quad (10.70)$$

19 The control variate version of this can be obtained by replacing $\tilde{g}_i(\mathbf{z}_s)$ with

$$\tilde{g}_i^{cv}(\mathbf{z}) = \tilde{g}_i(\mathbf{z}) + c_i(\mathbb{E}[b_i(\mathbf{z})] - b_i(\mathbf{z})) \quad (10.71)$$

22 where $b_i(\mathbf{z})$ is a baseline function and c_i is some constant, to be specified below. A convenient
23 baseline is the score function, $b_i(\mathbf{z}) = \nabla_{\psi_i} \log q_{\psi_i}(\mathbf{z}) = g_i(\mathbf{z})$, since this is correlated with $\tilde{g}_i(\mathbf{z})$, and
24 has the property that $\mathbb{E}[b_i(\mathbf{z})] = \mathbf{0}$, since the expected value of the score function is zero, as we
25 showed in Equation (3.44). Hence

$$\tilde{g}_i^{cv}(\mathbf{z}) = \tilde{g}_i(\mathbf{z}) - c_i g_i(\mathbf{z}) = g_i(\mathbf{z})(\tilde{\mathcal{L}}(\psi, \mathbf{z}) - c_i) \quad (10.72)$$

28 so the CV estimator is given by

$$\widehat{\nabla_{\psi_i} \mathbb{L}(\psi_t)}^{cv} = \frac{1}{S} \sum_{s=1}^S g_i(\mathbf{z}_s) \times (\tilde{\mathcal{L}}(\psi, \mathbf{z}_s) - c_i) \quad (10.73)$$

33 One can show that the optimal c_i that minimizes the variance of the CV estimator is

$$c_i = \frac{\text{Cov} [g_i(\mathbf{z}) \tilde{\mathcal{L}}(\psi, \mathbf{z}), g_i(\mathbf{z})]}{\mathbb{V}[g_i(\mathbf{z})]} \quad (10.74)$$

37 For more details, see e.g., [TND21].

40 10.3 Coordinate ascent VI

42 A common approximation in variational inference is to assume that all the latent variables are
43 independent, i.e.,

$$q_{\psi}(\mathbf{z}) = \prod_{j=1}^J q_j(z_j) \quad (10.75)$$

where J is the number of hidden variables, and $q_j(z_j)$ is shorthand for $q_{\psi_j}(z_j)$, where ψ_j are the variational parameters for the j 'th distribution. This is called the **mean field** approximation.

From Equation (10.11), the ELBO becomes

$$\mathcal{L}(\psi) = \int q_\psi(\mathbf{z}) \log p_\theta(\mathbf{x}, \mathbf{z}) d\mathbf{z} + \sum_{j=1}^J \mathbb{H}(q_j) \quad (10.76)$$

since the entropy of a product distribution is the sum of entropies of each component in the product. The first term also often decomposes according to the Markov properties of the graphical model. This allows us to use a coordinate ascent optimization scheme to estimate each ψ_j , as we explain in Section 10.3.1. This is called **coordinate ascent variational inference** or **CAVI**, and is an alternative to gradient-based VI.

14

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

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

where we define $\tilde{p}(\mathbf{z}) = p_\theta(\mathbf{z}, \mathbf{x})$ 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 :

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

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

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

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

1 Since $g_3(z_3) \propto \log f_3(z_3)$ we get
 2

$$3 \quad 4 \quad L_3(q_3) = \sum_{z_3} q_3(z_3) [\log f_3(z_3) - \log q_3(z_3)] + \text{const} = -D_{\text{KL}}(q_3 \parallel f_3) + \text{const} \quad (10.83)$$

5

6 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
 7 $q_3^*(z_3) = f_3(z_3)$.
 8

9 Now suppose that the joint distribution is defined by a Markov chain, where $z_1 \rightarrow z_2 \rightarrow z_3$, so
 10 $z_1 \perp z_3 | z_2$. 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
 11 of $q_3(z_3)$. Thus the ELBO simplifies to

$$12 \quad 13 \quad L_3(q_3) = \sum_{z_3} q_3(z_3) \left[\sum_{z_2} q_2(z_2) \log \tilde{p}(z_2, z_3) \right] + H(q_3) + \text{const} \quad (10.84)$$

14

$$15 \quad 16 \quad = \sum_{z_3} q_3(z_3) [\log f_3(z_3) - \log q_3(z_3)] + \text{const} \quad (10.85)$$

17

18 where

$$19 \quad 20 \quad f_3(z_3) \propto \exp \left[\sum_{z_2} q_2(z_2) \log \tilde{p}(z_2, z_3) \right] = \exp \left[\mathbb{E}_{z_{\text{mb}_3}} [\log \tilde{p}(z_2, z_3)] \right] \quad (10.86)$$

21

22 where $z_{\text{mb}_3} = z_2$ is the Markov blanket (Section 4.2.4.3) of z_3 . As before, the optimal variational
 23 distribution is given by $q_3(z_3) = f_3(z_3)$.

24 In general, when we have J groups of variables, the optimal variational distribution for the j 'th
 25 group is given by

$$26 \quad q_j(\mathbf{z}_j) \propto \exp \left[\mathbb{E}_{z_{\text{mb}_j}} [\log \tilde{p}(\mathbf{z}_j, z_{\text{mb}_j})] \right] \quad (10.87)$$

27

28 (Compare to the equation for Gibbs sampling in Equation (12.19).) The CAVI method simply
 29 computes q_j for each dimension j in turn, in an iterative fashion (see Algorithm 10.4). Convergence
 30 is guaranteed since the bound is convex wrt each of the factors q_i [Bis06, p. 466].
 31

32 **Algorithm 10.4:** Coordinate ascent variational inference (CAVI).

33
 34 1 Initialize $q_j(\mathbf{z}_j)$ for $j = 1 : J$
 35 2 **foreach** $t = 1 : T$ **do**
 36 3 **foreach** $j = 1 : J$ **do**
 37 4 Compute $g_j(\mathbf{z}_j) = \mathbb{E}_{z_{\text{mb}_i}} [\log \tilde{p}(\mathbf{z}_i, z_{\text{mb}_i})]$
 38 5 Compute $q_j(\mathbf{z}_j) \propto \exp(g_j(\mathbf{z}_j))$

41
 42 Note that the functional form of the q_i distributions does not need to be specified in advance, but
 43 will be determined by the form of the log joint. This is therefore called **free-form VI**, as opposed to
 44 fixed-form, where we explicitly choose a convenient distributional type for q (we discuss fixed-form
 45 VI in Section 10.2). We give some examples below that will make this clearer.
 46

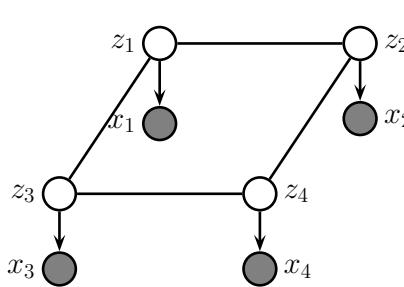


Figure 10.8: A grid-structured MRF with hidden nodes z_i and local evidence nodes x_i . The prior $p(\mathbf{z})$ is an undirected Ising model, and the likelihood $p(\mathbf{x}|\mathbf{z}) = \prod_i p(x_i|z_i)$ is a directed fully factored model.

10.3.2 Example: CAVI for the Ising model

In this section, we apply CAVI to perform mean field inference in an Ising model (Section 4.3.2.1), which is a kind of Markov random field defined on binary random variables, $z_i \in \{-1, +1\}$, arranged in a 2d grid.

Originally Ising models were developed as models of atomic spins for magnetic materials, although we will apply them to an image denoising problem. Specifically, let z_i be the hidden value of pixel i , and $x_i \in \mathbb{R}$ be the observed noisy value. See Figure 10.8 for the graphical model.

Let $L_i(z_i) \triangleq \log p(x_i|z_i)$ be the log likelihood for the i 'th pixel (aka the **local evidence** for node i in the graphical model). The overall likelihood has the form

$$p(\mathbf{x}|\mathbf{z}) = \prod_i p(x_i|z_i) = \exp\left(\sum_i L_i(z_i)\right) \quad (10.88)$$

Our goal is to approximate the posterior $p(\mathbf{z}|\mathbf{x})$. We will use an Ising model for the prior:

$$p(\mathbf{z}) = \frac{1}{Z_0} \exp(-\mathcal{E}_0(\mathbf{z})) \quad (10.89)$$

$$\mathcal{E}_0(\mathbf{z}) = -\sum_{i \sim j} W_{ij} z_i z_j \quad (10.90)$$

where we sum over each $i - j$ edge. Therefore the posterior has the form

$$p(\mathbf{z}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp(-\mathcal{E}(\mathbf{z})) \quad (10.91)$$

$$\mathcal{E}(\mathbf{z}) = \mathcal{E}_0(\mathbf{z}) - \sum_i L_i(z_i) \quad (10.92)$$

We will now make the following fully factored approximation:

$$q(\mathbf{z}) = \prod_i q_i(z_i) = \prod_i \text{Ber}(z_i|\mu_i) \quad (10.93)$$

where $\mu_i = \mathbb{E}_{q_i}[z_i]$ is the mean value of node i . To derive the update for the variational parameter μ_i , we first compute the unnormalized log joint, $\log \tilde{p}(\mathbf{z}) = -\mathcal{E}(\mathbf{z})$, dropping terms that do not involve

1
2 z_i :

3
4 $\log \tilde{p}(\mathbf{z}) = z_i \sum_{j \in \text{nbr}_i} W_{ij} z_j + L_i(z_i) + \text{const}$ (10.94)
5

6 This only depends on the states of the neighboring nodes. Hence
7

8
9 $q_i(z_i) \propto \exp(\mathbb{E}_{q_{-i}(\mathbf{z})} [\log \tilde{p}(\mathbf{z})]) = \exp \left(z_i \sum_{j \in \text{nbr}_i} W_{ij} \mu_j + L_i(z_i) \right)$ (10.95)
10

11 where $q_{-i}(\mathbf{z}) = \prod_{j \neq i} q(z_j)$. Thus we replace the states of the neighbors by their average values.
12 (Note that this replaces binary variables with continuous ones.)
13

14 We now simplify this expression. Let $m_i = \sum_{j \in \text{nbr}_i} W_{ij} \mu_j$ be the mean field influence on node i .
15 Also, let $L_i^+ \triangleq L_i(+1)$ and $L_i^- \triangleq L_i(-1)$. The approximate marginal posterior is given by
16

17 $q_i(z_i = 1) = \frac{e^{m_i + L_i^+}}{e^{m_i + L_i^+} + e^{-m_i + L_i^-}} = \frac{1}{1 + e^{-2m_i + L_i^- - L_i^+}} = \sigma(2a_i)$ (10.96)
18

19 $a_i \triangleq m_i + 0.5(L_i^+ - L_i^-)$ (10.97)
20

21 Similarly, we have $q_i(z_i = -1) = \sigma(-2a_i)$. From this we can compute the new mean for site i :
22

23 $\mu_i = \mathbb{E}_{q_i} [z_i] = q_i(z_i = +1) \cdot (+1) + q_i(z_i = -1) \cdot (-1)$ (10.98)
24

25 $= \frac{1}{1 + e^{-2a_i}} - \frac{1}{1 + e^{2a_i}} = \frac{e^{a_i}}{e^{a_i} + e^{-a_i}} - \frac{e^{-a_i}}{e^{-a_i} + e^{a_i}} = \tanh(a_i)$ (10.99)
26

27 We can turn the above equations into a fixed point algorithm by writing
28

29 $\mu_i^t = \tanh \left(\sum_{j \in \text{nbr}_i} W_{ij} \mu_j^{t-1} + 0.5(L_i^+ - L_i^-) \right)$ (10.100)
30

31 Following [MWJ99], we can use **damped updates** of the following form to improve convergence:
32

33
34 $\mu_i^t = (1 - \lambda) \mu_i^{t-1} + \lambda \tanh \left(\sum_{j \in \text{nbr}_i} W_{ij} \mu_j^{t-1} + 0.5(L_i^+ - L_i^-) \right)$ (10.101)
35

36 for $0 < \lambda < 1$. We can update all the nodes in parallel, or update them asynchronously.
37

38 Figure 10.9 shows the method in action, applied to a 2d Ising model with homogeneous attractive
39 potentials, $W_{ij} = 1$. We use parallel updates with a damping factor of $\lambda = 0.5$. (If we don't use
40 damping, we tend to get "checkerboard" artifacts.)
41

42 10.3.3 Variational Bayes 43

44 In Bayesian modeling, we treat the parameters $\boldsymbol{\theta}$ as latent variables. Thus our goal is to approximate
45 the parameter posterior $p(\boldsymbol{\theta}|\mathcal{D}) \propto p(\boldsymbol{\theta})p(\mathcal{D}|\boldsymbol{\theta})$. Applying VI to this problem is called **variational**
46 **Bayes** [Att00].
47

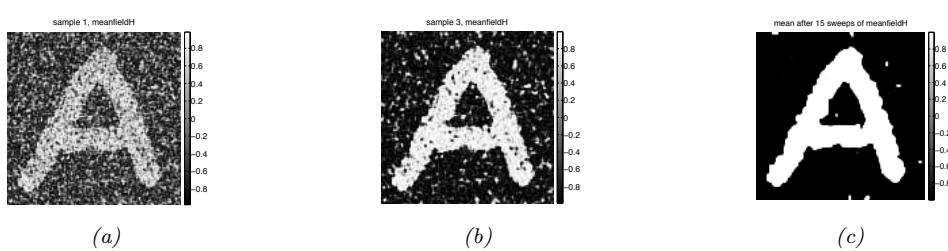


Figure 10.9: Example of image denoising using mean field (with parallel updates and a damping factor of 0.5). We use an Ising prior with $W_{ij} = 1$ and a Gaussian noise model with $\sigma = 2$. We show the results after 1, 3 and 15 iterations across the image. Compare to Figure 12.3, which shows the results of using Gibbs sampling. Generated by [ising_image_denoise_demo.ipynb](#).

In this section, we assume there are no latent variables except for the shared global parameters, so the model has the form

$$p(\boldsymbol{\theta}, \mathcal{D}) = p(\boldsymbol{\theta}) \prod_{n=1}^N p(\mathcal{D}_n | \boldsymbol{\theta}) \quad (10.102)$$

These conditional independencies are illustrated in Figure 10.5a.

We will fit the variational posterior by maximizing the ELBO

$$\mathcal{L}(\boldsymbol{\psi}_{\boldsymbol{\theta}} | \mathcal{D}) = \mathbb{E}_{q(\boldsymbol{\theta} | \boldsymbol{\psi}_{\boldsymbol{\theta}})} [\log p(\boldsymbol{\theta}, \mathcal{D})] + \mathbb{H}(q(\boldsymbol{\theta} | \boldsymbol{\psi}_{\boldsymbol{\theta}})) \quad (10.103)$$

We will assume the variational posterior factorizes over the parameters:

$$q(\boldsymbol{\theta} | \boldsymbol{\psi}_{\boldsymbol{\theta}}) = \prod_j q(\boldsymbol{\theta}_j | \boldsymbol{\psi}_{\theta_j}) \quad (10.104)$$

We can then update each $\boldsymbol{\psi}_{\theta_j}$ using CAVI (Section 10.3.1).

10.3.4 Example: VB for a univariate Gaussian

Consider inferring the parameters of a 1d Gaussian. The likelihood is given by $p(\mathcal{D} | \boldsymbol{\theta}) = \prod_{n=1}^N \mathcal{N}(x_n | \mu, \lambda^{-1})$, where μ is the mean and λ is the precision. Suppose we use a conjugate prior of the form

$$p(\mu, \lambda) = \mathcal{N}(\mu | \mu_0, (\kappa_0 \lambda)^{-1}) \text{Ga}(\lambda | a_0, b_0) \quad (10.105)$$

It is possible to derive the posterior $p(\mu, \lambda | \mathcal{D})$ for this model exactly, as shown in Section 3.4.3.3.

However, here we use the VB method with the following factored approximate posterior:

$$q(\mu, \lambda) = q(\mu | \boldsymbol{\psi}_{\mu}) q(\lambda | \boldsymbol{\psi}_{\lambda}) \quad (10.106)$$

We do not need to specify the forms for the distributions $q(\mu | \boldsymbol{\psi}_{\mu})$ and $q(\lambda | \boldsymbol{\psi}_{\lambda})$; the optimal forms will “fall out” automatically during the derivation (and conveniently, they turn out to be Gaussian and gamma respectively). Our presentation follows [Mac03, p429].

1 **10.3.4.1 Target distribution**

2 The unnormalized log posterior has the form

3 $\log \tilde{p}(\mu, \lambda) = \log p(\mu, \lambda, \mathcal{D}) = \log p(\mathcal{D}|\mu, \lambda) + \log p(\mu|\lambda) + \log p(\lambda)$ (10.107)

4 $= \frac{N}{2} \log \lambda - \frac{\lambda}{2} \sum_{n=1}^N (x_n - \mu)^2 - \frac{\kappa_0 \lambda}{2} (\mu - \mu_0)^2$

5 $+ \frac{1}{2} \log(\kappa_0 \lambda) + (a_0 - 1) \log \lambda - b_0 \lambda + \text{const}$ (10.108)

12 **10.3.4.2 Updating $q(\mu|\psi_\mu)$**

13 The optimal form for $q(\mu|\psi_\mu)$ is obtained by averaging over λ :

16 $\log q(\mu|\psi_\mu) = \mathbb{E}_{q(\lambda|\psi_\lambda)} [\log p(\mathcal{D}|\mu, \lambda) + \log p(\mu|\lambda)] + \text{const}$ (10.109)

18 $= -\frac{\mathbb{E}_{q(\lambda|\psi_\lambda)} [\lambda]}{2} \left\{ \kappa_0(\mu - \mu_0)^2 + \sum_{n=1}^N (x_n - \mu)^2 \right\} + \text{const}$ (10.110)

21 By completing the square one can show that $q(\mu|\psi_\mu) = \mathcal{N}(\mu|\mu_N, \kappa_N^{-1})$, where

24 $\mu_N = \frac{\kappa_0 \mu_0 + N \bar{x}}{\kappa_0 + N}, \quad \kappa_N = (\kappa_0 + N) \mathbb{E}_{q(\lambda|\psi_\lambda)} [\lambda]$ (10.111)

26 At this stage we don't know what $q(\lambda|\psi_\lambda)$ is, and hence we cannot compute $\mathbb{E}[\lambda]$, but we will derive
27 this below.

29 **10.3.4.3 Updating $q(\lambda|\psi_\lambda)$**

31 The optimal form for $q(\lambda|\psi_\lambda)$ is given by

33 $\log q(\lambda|\psi_\lambda) = \mathbb{E}_{q(\mu|\psi_\mu)} [\log p(\mathcal{D}|\mu, \lambda) + \log p(\mu|\lambda) + \log p(\lambda)] + \text{const}$ (10.112)

35 $= (a_0 - 1) \log \lambda - b_0 \lambda + \frac{1}{2} \log \lambda + \frac{N}{2} \log \lambda$

37 $- \frac{\lambda}{2} \mathbb{E}_{q(\mu|\psi_\mu)} \left[\kappa_0(\mu - \mu_0)^2 + \sum_{n=1}^N (x_n - \mu)^2 \right] + \text{const}$ (10.113)

40 We recognize this as the log of a gamma distribution, hence $q(\lambda|\psi_\lambda) = \text{Ga}(\lambda|a_N, b_N)$, where

42 $a_N = a_0 + \frac{N+1}{2}$ (10.114)

45 $b_N = b_0 + \frac{1}{2} \mathbb{E}_{q(\mu|\psi_\mu)} \left[\kappa_0(\mu - \mu_0)^2 + \sum_{n=1}^N (x_n - \mu)^2 \right]$ (10.115)

1
2 **10.3.4.4 Computing the expectations**

3 To implement the updates, we have to specify how to compute the various expectations. Since
4 $q(\mu) = \mathcal{N}(\mu | \mu_N, \kappa_N^{-1})$, we have
5

6 $\mathbb{E}_{q(\mu)} [\mu] = \mu_N \tag{10.116}$

7 $\mathbb{E}_{q(\mu)} [\mu^2] = \frac{1}{\kappa_N} + \mu_N^2 \tag{10.117}$

8 Since $q(\lambda) = \text{Ga}(\lambda | a_N, b_N)$, we have

9
11 $\mathbb{E}_{q(\lambda)} [\lambda] = \frac{a_N}{b_N} \tag{10.118}$

12 We can now give explicit forms for the update equations. For $q(\mu)$ we have
15

16 $\mu_N = \frac{\kappa_0 \mu_0 + N \bar{x}}{\kappa_0 + N} \tag{10.119}$

17 $\kappa_N = (\kappa_0 + N) \frac{a_N}{b_N} \tag{10.120}$

18 and for $q(\lambda)$ we have

19
22 $a_N = a_0 + \frac{N + 1}{2} \tag{10.121}$

24
25 $b_N = b_0 + \frac{1}{2} \kappa_0 (\mathbb{E} [\mu^2] + \mu_0^2 - 2\mathbb{E} [\mu] \mu_0) + \frac{1}{2} \sum_{n=1}^N (x_n^2 + \mathbb{E} [\mu^2] - 2\mathbb{E} [\mu] x_n) \tag{10.122}$

27 We see that μ_N and a_N are in fact fixed constants, and only κ_N and b_N need to be updated
28 iteratively. (In fact, one can solve for the fixed points of κ_N and b_N analytically, but we don't do
29 this here in order to illustrate the iterative updating scheme.)
30

31
32 **10.3.4.5 Illustration**

33 Figure 10.10 gives an example of this method in action. The green contours represent the exact
34 posterior, which is Gaussian-gamma. The dotted red contours represent the variational approximation
35 over several iterations. We see that the final approximation is reasonably close to the exact solution.
36 However, it is more "compact" than the true distribution. It is often the case that mean field inference
37 underestimates the posterior uncertainty, for reasons explained in Section 5.1.4.1.
38

39
40 **10.3.4.6 Lower bound**

41 In VB, we maximize a lower bound on the log marginal likelihood:

42
43 $\mathcal{L}(\psi_\theta | \mathcal{D}) \leq \log p(\mathcal{D}) = \log \iint p(\mathcal{D} | \mu, \lambda) p(\mu, \lambda) d\mu d\lambda \tag{10.123}$

44 It is very useful to compute the lower bound itself, for three reasons. First, it can be used to assess
45 convergence of the algorithm. Second, it can be used to assess the correctness of one's code: as with
46
47

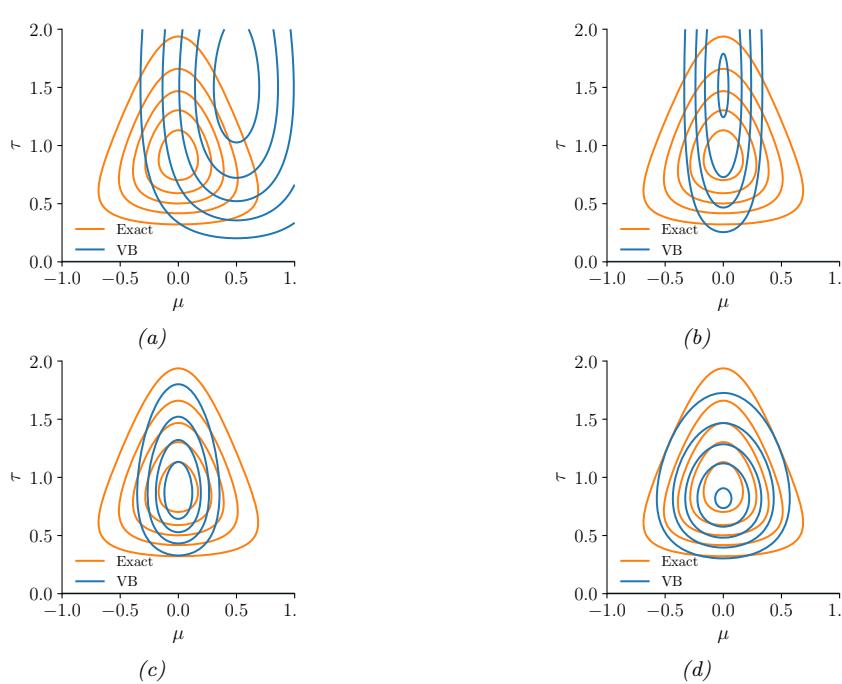


Figure 10.10: Factored variational approximation (orange) to the Gaussian-gamma distribution (blue). (a) Initial guess. (b) After updating $q(\mu|\psi_\mu)$. (c) After updating $q(\lambda|\psi_\lambda)$. (d) At convergence (after 5 iterations). Adapted from Fig. 10.4 of [Bis06]. Generated by `unigauss` `vb` `demo.ipynb`.

EM, if we use CAVI to optimize the objective, the bound should increase monotonically at each iteration, otherwise there must be a bug. Third, the bound can be used as an approximation to the marginal likelihood, which can be used for Bayesian model selection or empirical Bayes (see Section 10.1.3). In the case of the current model, one can show that the lower bound has the following form:

$$L = \text{const} + \frac{1}{2} \ln \frac{1}{\kappa_N} + \ln \Gamma(a_N) - a_N \ln b_N \quad (10.124)$$

10.3.5 Variational Bayes EM

In Bayesian latent variable models, we have two forms of hidden variables: local (or per example) hidden variables z_n , and global (shared) hidden variables θ , which represent the parameters of the model. See Figure 10.5b for an illustration. (Note that the parameters, which are fixed in number, are sometimes called **intrinsic variables**, whereas the local hidden variables are called **extrinsic variables**.) If $h = (\theta, z_{1:N})$ represents all the hidden variables, then the joint distribution is given by

$$p(\mathbf{h}, \mathcal{D}) = p(\boldsymbol{\theta}, \mathbf{z}_{1:N}, \mathcal{D}) = p(\boldsymbol{\theta}) \prod_{n=1}^N p(\mathbf{z}_n | \boldsymbol{\theta}) p(\mathbf{x}_n | \mathbf{z}_n, \boldsymbol{\theta}) \quad (10.125)$$

1 We will make the following mean field assumption:
2

3
4 $q(\boldsymbol{\theta}, \mathbf{z}_{1:N} | \psi_{1:N}, \boldsymbol{\psi}_{\boldsymbol{\theta}}) = q(\boldsymbol{\theta} | \boldsymbol{\psi}_{\boldsymbol{\theta}}) \prod_{n=1}^N q(\mathbf{z}_n | \psi_n)$ (10.126)
5

6 where $\boldsymbol{\psi} = (\psi_{1:N}, \boldsymbol{\psi}_{\boldsymbol{\theta}})$.
7

8 We will use VI to maximize the ELBO:

9 $L(\boldsymbol{\psi} | \mathcal{D}) = \mathbb{E}_{q(\boldsymbol{\theta}, \mathbf{z}_{1:N} | \psi_{1:N}, \boldsymbol{\psi}_{\boldsymbol{\theta}})} [\log p(\mathbf{z}_{1:N}, \boldsymbol{\theta}, \mathcal{D}) - \log q(\boldsymbol{\theta}, \mathbf{z}_{1:N})]$ (10.127)
10

11 If we use the mean field assumption, then we can apply the CAVI approach to optimize each set of
12 variational parameters. In particular, we can alternate between optimizing the $q_n(\mathbf{z}_n)$ in parallel,
13 independently of each other, with $q(\boldsymbol{\theta})$ held fixed, and then optimizing $q(\boldsymbol{\theta})$ with the q_n held fixed.
14 This is known as **variational Bayes EM** [BG06]. It is similar to regular EM, except in the E step,
15 we infer an approximate posterior for \mathbf{z}_n averaging out the parameters (instead of plugging in a point
16 estimate), and in the M step, we update the parameter posterior parameters using the expected
17 sufficient statistics.

18 Now suppose we approximate $q(\boldsymbol{\theta})$ by a delta function, $q(\boldsymbol{\theta}) = \delta(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})$. The Bayesian LVM ELBO
19 objective from Equation (10.127) simplifies to the “LVM ELBO”:

20
21 $L(\boldsymbol{\theta}, \boldsymbol{\psi}_{1:N} | \mathcal{D}) = \mathbb{E}_{q(\mathbf{z}_{1:N} | \psi_{1:N})} [\log p(\boldsymbol{\theta}, \mathcal{D}, \mathbf{z}_{1:N}) - \log q(\mathbf{z}_{1:N} | \psi_{1:N})]$ (10.128)

22 We can optimize this using the **variational EM** algorithm, which is a CAVI algorithm which updates
23 the $\boldsymbol{\psi}_n$ in parallel in the variational E step, and then updates $\boldsymbol{\theta}$ in the M step.
24

25 VEM is simpler than VBEM since in the variational E step, we compute $q(\mathbf{z}_n | \mathbf{x}_n, \hat{\boldsymbol{\theta}})$, instead of
26 $\mathbb{E}_{\boldsymbol{\theta}}[q(\mathbf{z}_n | \mathbf{x}_n, \boldsymbol{\theta})]$; that is, we plug in a point estimate of the model parameters, rather than averaging
27 over the parameters. For more details on VEM, see Section 10.1.3.

28 10.3.6 Example: VBEM for a GMM

30 Consider a standard Gaussian mixture model (GMM):
31

32 $p(\mathbf{z}, \mathbf{x} | \boldsymbol{\theta}) = \prod_n \prod_k \pi_k^{z_{nk}} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1})^{z_{nk}}$ (10.129)
33

34 where $z_{nk} = 1$ if datapoint n belongs to cluster k , and $z_{nk} = 0$ otherwise. Our goal is to approximate
35 the posterior $p(\mathbf{z}, \boldsymbol{\theta} | \mathbf{x})$ under the following conjugate prior
36

37 $p(\boldsymbol{\theta}) = \text{Dir}(\boldsymbol{\pi} | \boldsymbol{\alpha}) \prod_k \mathcal{N}(\boldsymbol{\mu}_k | \tilde{\boldsymbol{m}}, (\kappa \boldsymbol{\Lambda}_k)^{-1}) \text{Wi}(\boldsymbol{\Lambda}_k | \tilde{\mathbf{L}}, \tilde{\nu})$ (10.130)
38

39 where $\boldsymbol{\Lambda}_k$ is the precision matrix for cluster k . For the mixing weights, we usually use a symmetric
40 prior, $\boldsymbol{\alpha} = \alpha_0 \mathbf{1}$.

41 The exact posterior $p(\mathbf{z}, \boldsymbol{\theta} | \mathcal{D})$ is a mixture of K^N distributions, corresponding to all possible
42 labelings \mathbf{z} , which is intractable to compute. In this section, we derive a VBEM algorithm, which
43 will approximate the posterior around a local mode. We follow the presentation of [Bis06, Sec 10.2].
44 (See also Section 10.2.1.5 and Section 10.2.2.3, where we discuss variational approximations based on
45 stochastic gradient descent, which can scale better to large datasets compared to VBEM.)
46

47

1
2 **10.3.6.1 The variational posterior**

3 We will use the standard mean field approximation to the posterior: $q(\boldsymbol{\theta}, \mathbf{z}_{1:N}) = q(\boldsymbol{\theta}) \prod_n q_n(\mathbf{z}_n)$. At
4 this stage we have not specified the forms of the q functions; these will be determined by the form of
5 the likelihood and prior. Below we will show that the optimal forms are as follows:
6

7 $q_n(z_n) = \text{Cat}(\mathbf{z}_n | \mathbf{r}_n)$ (10.131)
8

9 $q(\boldsymbol{\theta}) = \text{Dir}(\boldsymbol{\pi} | \hat{\boldsymbol{\alpha}}) \prod_k \mathcal{N}(\boldsymbol{\mu}_k | \hat{\mathbf{m}}_k, (\hat{\kappa}_k \boldsymbol{\Lambda}_k)^{-1}) \text{Wi}(\boldsymbol{\Lambda}_k | \hat{\mathbf{L}}_k, \hat{\nu}_k)$ (10.132)
10

11 where \mathbf{r}_n are the posterior responsibilities, and the parameters with hats on them are the hyperpa-
12 rameters from the prior updated with data.
13

14
15 **10.3.6.2 Derivation of $q(\boldsymbol{\theta})$ (variational M step)**
16

17 Using the mean field recipe in Algorithm 10.4, we write down the log joint, and take expectations
18 over all variables except $\boldsymbol{\theta}$, so we average out the \mathbf{z}_n wrt $q(\mathbf{z}_n) = \text{Cat}(z_n | \mathbf{r}_n)$:

19
20
$$\log q(\boldsymbol{\theta}) = \log p(\boldsymbol{\pi}) + \underbrace{\sum_n \mathbb{E}_{q(z_n)} [\log p(\mathbf{z}_n | \boldsymbol{\pi})]}_{L_{\boldsymbol{\pi}}}$$

21
22
23
24
25 $+ \sum_k \left[\underbrace{\log p(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k) \sum_n \mathbb{E}_{q(z_n)} [z_{nk}] \log \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1})}_{L_{\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k}} \right] + \text{const}$ (10.133)
26
27
28

29 Since the expected log joint factorizes into a term involving $\boldsymbol{\pi}$ and terms involving $(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)$, we see
30 that the variational posterior also factorizes into the form
31

32
33 $q(\boldsymbol{\theta}) = q(\boldsymbol{\pi}) \prod_k q(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)$ (10.134)
34

35 For the $\boldsymbol{\pi}$ term, we have
36

37
38 $\log q(\boldsymbol{\pi}) = (\alpha_0 - 1) \sum_k \log \pi_k + \sum_k \sum_n r_{nk} \log \pi_k + \text{const}$ (10.135)
39

40 Exponentiating, we recognize this as a Dirichlet distribution:
41

42 $q(\boldsymbol{\pi}) = \text{Dir}(\boldsymbol{\pi} | \hat{\boldsymbol{\alpha}})$ (10.136)
43

44 $\hat{\alpha}_k = \alpha_0 + N_k$ (10.137)
45

46 $N_k = \sum_n r_{nk}$ (10.138)
47

For the $\boldsymbol{\mu}_k$ and $\boldsymbol{\Lambda}_k$ terms, we have

$$q(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k) = \mathcal{N}(\boldsymbol{\mu}_k | \widehat{\boldsymbol{m}}_k, (\widehat{\kappa}_k \boldsymbol{\Lambda}_k)^{-1}) \text{Wi}(\boldsymbol{\Lambda}_k | \widehat{\mathbf{L}}_k, \widehat{\nu}_k) \quad (10.139)$$

$$\widehat{\kappa}_k = \check{\kappa} + N_k \quad (10.140)$$

$$\widehat{\boldsymbol{m}}_k = (\check{\kappa} \check{\boldsymbol{m}} + N_k \bar{\boldsymbol{x}}_k) / \widehat{\kappa}_k \quad (10.141)$$

$$\widehat{\mathbf{L}}_k^{-1} = \check{\mathbf{L}}^{-1} + N_k \mathbf{S}_k + \frac{\check{\kappa} N_k}{\check{\kappa} + N_k} (\bar{\boldsymbol{x}}_k - \check{\boldsymbol{m}})(\bar{\boldsymbol{x}}_k - \check{\boldsymbol{m}})^\top \quad (10.142)$$

$$\widehat{\nu}_k = \check{\nu} + N_k \quad (10.143)$$

$$\bar{\boldsymbol{x}}_k = \frac{1}{N_k} \sum_n r_{nk} \boldsymbol{x}_n \quad (10.144)$$

$$\mathbf{S}_k = \frac{1}{N_k} \sum_n r_{nk} (\boldsymbol{x}_n - \bar{\boldsymbol{x}}_k)(\boldsymbol{x}_n - \bar{\boldsymbol{x}}_k)^\top \quad (10.145)$$

This is very similar to the M step for MAP estimation for GMMs, except here we are computing the parameters of the posterior for $\boldsymbol{\theta}$ rather than a point estimate $\hat{\boldsymbol{\theta}}$.

10.3.6.3 Derivation of $q(\boldsymbol{z})$ (variational E step)

The variational E step is more interesting, since it is quite different from the E step in regular EM, because we need to average over the parameters, rather than condition on them. In particular, we have

$$\begin{aligned} \log q(\boldsymbol{z}) &= \sum_n \sum_k z_{nk} \left(\mathbb{E}_{q(\boldsymbol{\pi})} [\log \pi_k] + \frac{1}{2} \mathbb{E}_{q(\boldsymbol{\Lambda}_k)} [\log |\boldsymbol{\Lambda}_k|] - \frac{D}{2} \log(2\pi) \right. \\ &\quad \left. - \frac{1}{2} \mathbb{E}_{q(\boldsymbol{\theta})} [(\boldsymbol{x}_n - \boldsymbol{\mu}_k)^\top \boldsymbol{\Lambda}_k (\boldsymbol{x}_n - \boldsymbol{\mu}_k)] \right) + \text{const} \end{aligned} \quad (10.146)$$

Using the fact that $q(\boldsymbol{\pi}) = \text{Dir}(\boldsymbol{\pi} | \boldsymbol{\alpha})$, one can show that

$$\exp(\mathbb{E}_{q(\boldsymbol{\pi})} [\log \pi_k]) = \frac{\exp(\psi(\widehat{\alpha}_k))}{\exp(\psi(\sum_{k'} \widehat{\alpha}_{k'}))} \triangleq \tilde{\pi}_k \quad (10.147)$$

where ψ is the **digamma function**:

$$\psi(x) = \frac{d}{dx} \log \Gamma(x) \quad (10.148)$$

This takes care of the first term.

For the second term, one can show

$$\mathbb{E}_{q(\boldsymbol{\Lambda}_k)} [\log |\boldsymbol{\Lambda}_k|] = \sum_{j=1}^D \psi\left(\frac{\widehat{\nu}_k + 1 - j}{2}\right) + D \log 2 + \log |\widehat{\mathbf{L}}_k| \quad (10.149)$$

Finally, for the expected value of the quadratic form, one can show

$$\mathbb{E}_{q(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)} [(\boldsymbol{x}_n - \boldsymbol{\mu}_k)^\top \boldsymbol{\Lambda}_k (\boldsymbol{x}_n - \boldsymbol{\mu}_k)] = D \widehat{\kappa}_k^{-1} + \widehat{\nu}_k (\boldsymbol{x}_n - \widehat{\boldsymbol{m}}_k)^\top \widehat{\mathbf{L}}_k (\boldsymbol{x}_n - \widehat{\boldsymbol{m}}_k) \triangleq \tilde{\Lambda}_k \quad (10.150)$$

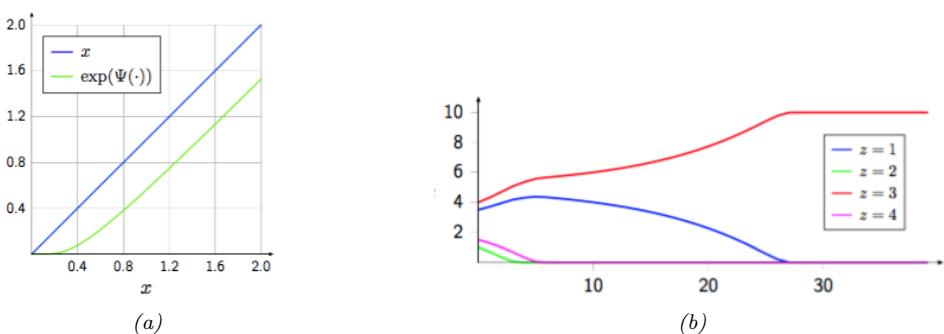


Figure 10.11: (a) We plot $\exp(\psi(x))$ vs x . We see that this function performs a form of shrinkage, so that small values get set to zero. (b) We plot N_k vs time for 4 different states (z values), starting from random initial values. We perform a series of VBEM updates, ignoring the likelihood term. We see that states that initially had higher counts get reinforced, and sparsely populated states get killed off. From [LK07]. Used with kind permission of Percy Liang.

Thus we get that the posterior responsibility of cluster k for datapoint n is

$$r_{nk} \propto \tilde{\pi}_k \tilde{\Lambda}_k^{\frac{1}{2}} \exp\left(-\frac{D}{2 \tilde{\kappa}_k} - \frac{\hat{\nu}_k}{2} (\mathbf{x}_n - \hat{\boldsymbol{\mu}}_k)^T \hat{\boldsymbol{\Lambda}}_k (\mathbf{x}_n - \hat{\boldsymbol{\mu}}_k)\right) \quad (10.151)$$

Compare this to the expression used in regular EM:

$$r_{nk}^{EM} \propto \hat{\pi}_k |\hat{\boldsymbol{\Lambda}}_k|^{\frac{1}{2}} \exp\left(-\frac{1}{2} (\mathbf{x}_n - \hat{\boldsymbol{\mu}}_k)^T \hat{\boldsymbol{\Lambda}}_k (\mathbf{x}_n - \hat{\boldsymbol{\mu}}_k)\right) \quad (10.152)$$

where $\hat{\pi}_k$ is the MAP estimate for π_k . The significance of this difference is discussed in Section 10.3.6.4.

10.3.6.4 Automatic sparsity inducing effects of VBEM

In regular EM, the E step has the form given in Equation (10.152), whereas in VBEM, the E step has the form given in Equation (10.151). Although they look similar, they differ in an important way. To understand this, let us ignore the likelihood term, and just focus on the prior. From Equation (10.147) we have

$$r_{nk}^{VB} = \tilde{\pi}_k = \frac{\exp(\psi(\hat{\alpha}_k))}{\exp(\psi(\sum_{k'} \hat{\alpha}_{k'}))} \quad (10.153)$$

And from the usual EM MAP estimation equations for GMM mixing weights (see e.g., [Mur22, Sec 8.7.3.4]) we have

$$r_{nk}^{EM} = \hat{\pi}_k = \frac{\hat{\alpha}_k - 1}{\sum_{k'} (\hat{\alpha}_{k'} - 1)} \quad (10.154)$$

where $\hat{\alpha}_k = \alpha_0 + N_k$, and $N_k = \sum_n r_{nk}$ is the expected number of assignments to cluster k .

We know from Figure 2.6 that using $\alpha_0 \ll 1$ causes π to be sparse, which will encourage \mathbf{r}_n to be sparse, which will “kill off” unnecessary mixture components (i.e., ones for which $N_k \ll N$, meaning very few datapoints are assigned to cluster k). To encourage this sparsity promoting effect, let us set $\alpha_0 = 0$. In this case, the updated parameters for the mixture weights are given by the following:

$$\tilde{\pi}_k = \frac{\exp(\psi(N_k))}{\exp(\psi(\sum_{k'} N_{k'}))} \quad (10.155)$$

$$\hat{\pi}_k = \frac{N_k - 1}{\sum_{k'} (N_{k'} - 1)} \quad (10.156)$$

Now consider a cluster which has no assigned data, so $N_k = 0$. In regular EM, $\hat{\pi}_k$ might end up negative, as pointed out in [FJ02]. (This will not occur if we use maximum likelihood training, which corresponds to $\alpha_0 = 1$, but this will not induce any sparsity, either.) This problem does not arise in VBEM, since we use the digamma function, which is always positive, as shown in Figure 10.11(a).

More interestingly, let us consider the effect of these updates on clusters that have unequal, but non-zero, number of assignments. Suppose we start with a random assignment of counts to 4 clusters, and iterate the VBEM algorithm, ignoring the contribution from the likelihood for simplicity. Figure 10.11(b) shows how the counts N_k evolve over time. We notice that clusters that started out with small counts end up with zero counts, and clusters that started out with large counts end up with even larger counts. In other words, the initially popular clusters get more and more members. This is called the **rich get richer** phenomenon; we will encounter it again in Supplementary Section 31.2, when we discuss Dirichlet process mixture models.

The reason for this effect is shown in Figure 10.11(a): we see that $\exp(\psi(N_k)) < N_k$, and is zero if N_k is sufficiently small, similar to the soft-thresholding behavior induced by ℓ_1 -regularization (see Section 15.2.6). Importantly, this effect of reducing N_k is greater on clusters with small counts.

We now demonstrate this automatic pruning method on a real example. We fit a mixture of 6 Gaussians to the Old Faithful dataset, using $\alpha_0 = 0.001$. Since the data only really “needs” 2 clusters, the remaining 4 get “killed off”, as shown in Figure 10.12. In Figure 10.13, we plot the initial and final values of α_k ; we see that $\bar{\alpha}_k = 0$ for all but two of the components k .

Thus we see that VBEM for GMMs with a sparse Dirichlet prior provides an efficient way to choose the number of clusters. Similar techniques can be used to choose the number of states in an HMM and other latent variable models. However, this **variational pruning effect** (also called **posterior collapse**), is not always desirable, since it can cause the model to “ignore” the latent variables \mathbf{z} if the likelihood function $p(\mathbf{x}|\mathbf{z})$ is sufficiently powerful. We discuss this more in Section 21.4.

37

38 10.3.6.5 Lower bound on the marginal likelihood 39

40 The VBEM algorithm is maximizing the following lower bound
41

$$42 \quad \mathcal{L} = \sum_{\mathbf{z}} \int d\boldsymbol{\theta} q(\mathbf{z}, \boldsymbol{\theta}) \log \frac{p(\mathbf{x}, \mathbf{z}, \boldsymbol{\theta})}{q(\mathbf{z}, \boldsymbol{\theta})} \leq \log p(\mathbf{x}) \quad (10.157)$$

45 This quantity increases monotonically with each iteration, as shown in Figure 10.14.
46

47

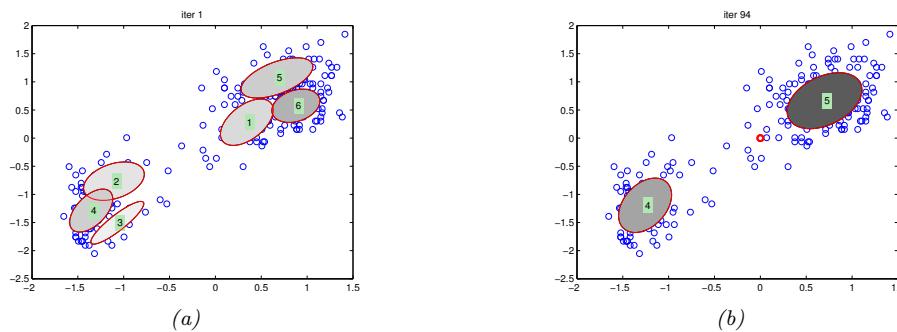


Figure 10.12: We visualize the posterior mean parameters at various stages of the VBEM algorithm applied to a mixture of Gaussians model on the Old Faithful data. Shading intensity is proportional to the mixing weight. We initialize with K-means and use $\alpha_0 = 0.001$ as the Dirichlet hyper-parameter. (The red dot on the right panel represents all the unused mixture components, which collapse to the prior at 0.) Adapted from Figure 10.6 of [Bis06]. Generated by [gmm_vb_em.ipynb](#).

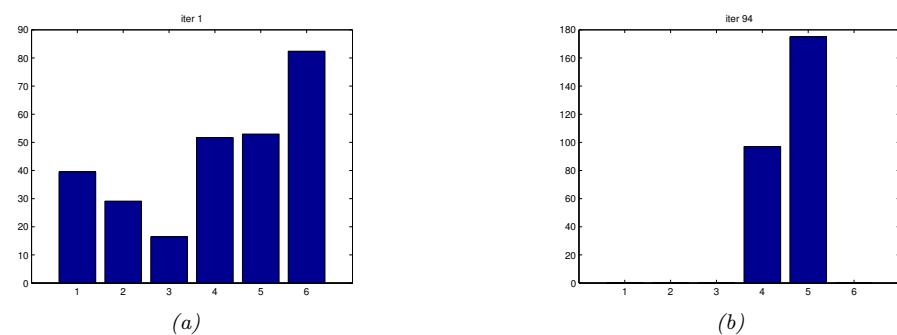


Figure 10.13: We visualize the posterior values of α_k for the model in Figure 10.12 after the first and last iteration of the algorithm. We see that unnecessary components get “killed off”. (Interestingly, the initially large cluster 6 gets “replaced” by cluster 5.) Generated by [gmm_vb_em.ipynb](#).

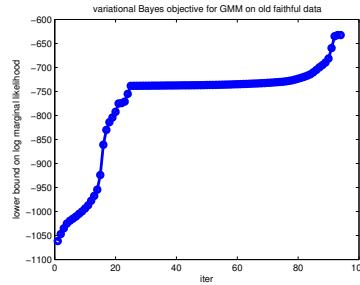


Figure 10.14: Lower bound vs iterations for the VB algorithm in Figure 10.12. The steep parts of the curve correspond to places where the algorithm figures out that it can increase the bound by “killing off” unnecessary mixture components, as described in Section 10.3.6.6. The plateaus correspond to slowly moving the clusters around. Generated by [gmm_vb_em.ipynb](#).

1 **10.3.6.6 Model selection using VBEM**

3 Section 10.3.6.4 discusses a way to choose K automatically, during model fitting, by “killing off”
4 unneeded clusters. An alternative approach is to fit several models, and then to use the variational
5 lower bound to the log marginal likelihood, $\mathcal{L}(K) \leq \log p(\mathcal{D}|K)$, to approximate $p(K|\mathcal{D})$. In particular,
6 if we have a uniform prior, we get the posterior
7

$$\underline{8} \quad p(K|\mathcal{D}) = \frac{p(\mathcal{D}|K)}{\sum_{K'} p(\mathcal{D}|K')} \approx \frac{e^{\mathcal{L}(K)}}{\sum_{K'} e^{\mathcal{L}(K')}} \quad (10.158)$$

11 It is shown in [BG06] that the VB approximation to the marginal likelihood is more accurate than
12 BIC [BG06]. However, the lower bound needs to be modified somewhat to take into account the
13 lack of identifiability of the parameters. In particular, although VB will approximate the volume
14 occupied by the parameter posterior, it will only do so around one of the local modes. With K
15 components, there are $K!$ equivalent modes, which differ merely by permuting the labels. Therefore a
16 more accurate approximation to the log marginal likelihood is to use $\log p(\mathcal{D}|K) \approx \mathcal{L}(K) + \log(K!)$.
17

18 **10.3.7 Variational message passing (VMP)**

19 In this section, we describe the CAVI algorithm for a generic model in which each complete conditional,
20 $p(\mathbf{z}_j|\mathbf{z}_{-j}, \mathbf{x})$, is in the exponential family, i.e.,

$$\underline{22} \quad p(\mathbf{z}_j|\mathbf{z}_{-j}, \mathbf{x}) = h(\mathbf{z}_j) \exp[\boldsymbol{\eta}_j(\mathbf{z}_{-j}, \mathbf{x})^\top \mathcal{T}(\mathbf{z}_j) - A_j(\boldsymbol{\eta}_j(\mathbf{z}_{-j}, \mathbf{x}))] \quad (10.159)$$

24 where $\mathcal{T}(\mathbf{z}_j)$ is the vector of sufficient statistics, $\boldsymbol{\eta}_j$ are the natural parameters, A_j is the log partition
25 function, and $h(\mathbf{z}_j)$ is the base distribution. This assumption holds if the prior $p(\mathbf{z}_j)$ is conjugate to
26 the likelihood, $p(\mathbf{z}_{-j}, \mathbf{x}|\mathbf{z}_j)$.

27 If Equation (10.159) holds, the mean field update node j becomes
28

$$\underline{29} \quad q_j(\mathbf{z}_j) \propto \exp [\mathbb{E} [\log p(\mathbf{z}_j|\mathbf{z}_{-j}, \mathbf{x})]] \quad (10.160)$$

$$\underline{30} \quad = \exp \left[\log h(\mathbf{z}_j) + \mathbb{E} [\boldsymbol{\eta}_j(\mathbf{z}_{-j}, \mathbf{x})]^\top \mathcal{T}(\mathbf{z}_j) - \mathbb{E} [A_j(\boldsymbol{\eta}_j(\mathbf{z}_{-j}, \mathbf{x}))] \right] \quad (10.161)$$

$$\underline{32} \quad \propto h(\mathbf{z}_j) \exp \left[\mathbb{E} [\boldsymbol{\eta}_j(\mathbf{z}_{-j}, \mathbf{x})]^\top \mathcal{T}(\mathbf{z}_j) \right] \quad (10.162)$$

34 Thus we update the local natural parameters using the expected values of the other nodes. These
35 become the new variational parameters:
36

$$\underline{37} \quad \boldsymbol{\psi}_j = \mathbb{E} [\boldsymbol{\eta}_j(\mathbf{z}_{-j}, \mathbf{x})] \quad (10.163)$$

39 We can generalize the above approach to work with any model where each full conditional is
40 conjugate. The resulting algorithm is known as **variational message passing** or **VMP** [WB05]
41 that works for any directed graphical model. VMP is similar to belief propagation (Section 9.3): at
42 each iteration, each node collects all the messages from its parents, and all the messages from its
43 children (which might require the children to get messages from their co-parents), and combines them
44 to compute the expected value of the node’s sufficient statistics. The messages that are sent are the
45 expected sufficient statistics of a node, rather than just a discrete or Gaussian distribution (as in BP).
46 Several software libraries have implemented this framework (see e.g., [Win; Min+18; Lut16; Wan17]).
47

VMP can be extended to the case where each full conditional is conditionally conjugate using the CVI framework in [Supplementary](#) Section 10.3.1. See also [ABV21], where they use local Laplace approximations to intractable factors inside of a message passing framework.

10.3.8 Autoconj

The VMP method requires the user to manually specify a graphical model; the corresponding node update equations are then computed for each node using a lookup table, for each possible combination of node types. It is possible to automatically derive these update equations for any conditionally conjugate directed graphical model using a technique called **autoconj** [HJT18]. This is analogous to the use of automatic differentiation (autodiff) to derive the gradient for any differentiable function. (Note that autoconj uses autodiff internally.) The resulting full conditionals can be used for CAVI, and also for Gibbs sampling (Section 12.3).

10.4 More accurate variational posteriors

In general, we can improve the tightness of the ELBO lower bound, and hence reduce the KL divergence of our posterior approximation, if we use more flexible posterior families (although optimizing within more flexible families may be slower, and can incur statistical error if the sample size is low [Bha+21]). In this section, we give several examples of more accurate variational posteriors, going beyond fully factored mean field approximations, or simple unimodal Gaussian approximations.

10.4.1 Structured mean field

The mean field assumption is quite strong, and can sometimes give poor results. Fortunately, sometimes we can exploit **tractable substructure** in our problem, so that we can efficiently handle some kinds of dependencies between the variables in the posterior in an analytic way, rather than assuming they are all independent. This is called the **structured mean field** approach [SJ95].

A common example arises when applying VI to time series models, such as HMMs, where the latent variables within each sequence are usually highly correlated across time. Rather than assuming a fully factorized posterior, we can treat each sequence $\mathbf{z}_{n,1:T}$ as a block, and just assume independence between blocks and the parameters: $q(\mathbf{z}_{1:N,1:T}, \boldsymbol{\theta}) = q(\boldsymbol{\theta}) \prod_{n=1}^N q(\mathbf{z}_{n,1:T})$, where $q(\mathbf{z}_{n,1:T}) = \prod_t q(\mathbf{z}_{n,t} | \mathbf{z}_{n,t-1})$. We can compute the joint distribution $q(\mathbf{z}_{n,1:T})$, taking into account the dependence between time steps, using the forwards-backwards algorithm. For details, see [JW14; Fot+14]. A similar approach was applied to the factorial HMM model, as we discuss in [Supplementary](#) Section 10.3.2.

An automatic way to derive a structured variational approximation to a probabilistic model, specified by a probabilistic programming language, is discussed in [AHG20].

10.4.2 Hierarchical (auxiliary variable) posteriors

Suppose $q_\phi(\mathbf{z}|\mathbf{x}) = \prod_k q_\phi(z_k|\mathbf{x})$ is a factorized distribution, such as a diagonal Gaussian. This does not capture dependencies between the latent variables (components of \mathbf{z}). We could of course use a full covariance matrix, but this might be too expensive.

An alternative approach is to use a hierarchical model, in which we add **auxiliary latent variables** \mathbf{a} , which are used to increase the flexibility of the variational posterior. In particular, we can still assume $q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{a})$ is conditionally factorized, but when we marginalize out \mathbf{a} , we induce dependencies between the elements of \mathbf{z} , i.e.,

$$q_\phi(\mathbf{z}|\mathbf{x}) = \int q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{a})q_\phi(\mathbf{a}|\mathbf{x})d\mathbf{a} \neq \prod_k q_\phi(z_k|\mathbf{x}) \quad (10.164)$$

This is called a **hierarchical variational model** [Ran16], or an **auxiliary variable deep generative model** [Maa+16].

In [TRB16], they model $q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{a})$ as a Gaussian process, which is a flexible nonparametric distribution (see Chapter 18), where \mathbf{a} are the inducing points. This combination is called a **variational GP**.

10.4.3 Normalizing flow posteriors

Normalizing flows are a class of probability models which work by passing a simple source distribution, such as a diagonal Gaussian, through a series of nonlinear, but invertible, mappings f to create a more complex distribution. This can be used to get more accurate posterior approximations than standard Gaussian VI, as we discuss in Section 23.1.2.2.

10.4.4 Implicit posteriors

In Chapter 26, we discuss implicit probability distributions, which are models which we can sample from, but which we cannot evaluate pointwise. For example, consider passing a Gaussian noise term, $\mathbf{z}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, through a nonlinear, *non-invertible* mapping f to create $\mathbf{z} = f(\mathbf{z}_0)$; it is easy to sample from $q(\mathbf{z})$, but it is intractable to evaluate the density $q(\mathbf{z})$ (unlike with flows). This makes it hard to evaluate the log density ratio $\log p_\theta(\mathbf{z})/q_\psi(\mathbf{z}|\mathbf{x})$, which is needed to compute the ELBO. However, we can use the same method as is used in GANs (generative adversarial networks, Chapter 26), in which we train a classifier that discriminates prior samples from samples from the variational posterior by evaluating $T(\mathbf{x}, \mathbf{z}) = \log q_\psi(\mathbf{z}|\mathbf{x}) - \log p_\theta(\mathbf{z})$. See e.g., [TR19] for details.

10.4.5 Combining VI with MCMC inference

There are various ways to combine variational inference with MCMC to get an improved approximate posterior. In [SKW15], they propose **Hamiltonian variational inference**, in which they train an inference network to initialize an HMC sampler (Section 12.5). The gradient of the log posterior (wrt the latents), which is needed by HMC, is given by

$$\nabla_{\mathbf{z}} \log p_\theta(\mathbf{z}|\mathbf{x}) = \nabla_{\mathbf{z}} \log [p_\theta(\mathbf{x}, \mathbf{z}) - \log p_\theta(\mathbf{x})] = \nabla_{\mathbf{z}} \log p_\theta(\mathbf{x}, \mathbf{z}) \quad (10.165)$$

This is easy to compute. They use the final sample to approximate the posterior $q_\phi(\mathbf{z}|\mathbf{x})$. To compute the entropy of this distribution, they also learn an auxiliary inverse inference network to reverse the HMC Markov chain.

A simpler approach is proposed in [Hof17]. Here they train an inference network to initialize an HMC sampler, using the standard ELBO for ϕ , but they optimize the generative parameters θ using

a stochastic approximation to the log marginal likelihood, given by $\log p_{\theta}(\mathbf{z}, \mathbf{x})$ where \mathbf{z} is a sample from the HMC chain. This does not require learning a reverse inference network, and avoids problems with variational pruning, since it does not use the ELBO for training the generative model.

10.5 Tighter bounds

Another way to improve the quality of the posterior approximation is to optimize q wrt a bound that is a tighter approximation to the log marginal likelihood compared to the standard ELBO. We give some examples below.

10.5.1 Multi-sample ELBO (IWAE bound)

In this section, we discuss a method known as the **importance weighted autoencoder** or **IWAE** [BGS16], which is a way to tighten the variational lower bound by using self-normalized importance sampling (Section 11.5.2). (It can also be interpreted as standard ELBO maximization in an expanded model, where we add extra auxiliary variables [CMD17; DS18; Tuc+19].)

Let the inference network $q_{\phi}(\mathbf{z}|\mathbf{x})$ be viewed as a proposal distribution for the target posterior $p_{\theta}(\mathbf{z}|\mathbf{x})$. Define $w_s^* = \frac{p_{\theta}(\mathbf{x}, \mathbf{z}_s)}{q_{\phi}(\mathbf{z}_s|\mathbf{x})}$ as the unnormalized importance weight for a sample, and $w_s = w_s^*/(\sum_{s'=1}^S w_{s'}^*)$ as the normalized importance weights. From Equation (11.43) we can compute an estimate of the marginal likelihood $p(\mathbf{x})$ using

$$\hat{p}_S(\mathbf{x}|\mathbf{z}_{1:S}) \triangleq \frac{1}{S} \sum_{k=1}^S \frac{p_{\theta}(\mathbf{x}, \mathbf{z}_s)}{q_{\phi}(\mathbf{z}_s|\mathbf{x})} = \frac{1}{S} \sum_{k=1}^S w_s \quad (10.166)$$

This is unbiased, i.e., $\mathbb{E}_{q_{\phi}(\mathbf{z}_{1:S}|\mathbf{x})} [\hat{p}_S(\mathbf{x}|\mathbf{z}_{1:S})] = p(\mathbf{x})$, where $q_{\phi}(\mathbf{z}_{1:S}|\mathbf{x}) = \prod_{s=1}^S q_{\phi}(\mathbf{z}_s|\mathbf{x})$. In addition, since the estimator is always positive, we can take logarithms, and thus obtain a stochastic lower bound on the log likelihood:

$$\mathbb{L}_S(\phi, \theta|\mathbf{x}) \triangleq \mathbb{E}_{q_{\phi}(\mathbf{z}_{1:S}|\mathbf{x})} \left[\log \left(\frac{1}{S} \sum_{s=1}^S w_s \right) \right] = \mathbb{E}_{q_{\phi}(\mathbf{z}_{1:S}|\mathbf{x})} [\log \hat{p}_S(\mathbf{z}_{1:S})] \quad (10.167)$$

$$\leq \log \mathbb{E}_{q_{\phi}(\mathbf{z}_{1:S}|\mathbf{x})} [\hat{p}_S(\mathbf{z}_{1:S})] = \log p(\mathbf{x}) \quad (10.168)$$

where we used Jensen's inequality in the penultimate line, and the unbiased property in the last line. This is called the **multi-sample ELBO** or **IWAE bound** [BGS16]. The gradients of this expression wrt θ and ϕ are given in Equation (10.179). If $S = 1$, \mathbb{L}_S reduces to the standard ELBO:

$$\mathbb{L}_1(\phi, \theta|\mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log w] = \int q_{\phi}(\mathbf{z}|\mathbf{x}) \log \frac{p_{\theta}(\mathbf{z}, \mathbf{x})}{q_{\phi}(\mathbf{z}|\mathbf{x})} d\mathbf{z} \quad (10.169)$$

One can show [BGS16] that increasing the number of samples S is guaranteed to make the bound tighter, thus making it a better proxy for the log likelihood. Intuitively, averaging the S samples inside the log removes the need for every sample \mathbf{z}_s to explain the data \mathbf{x} . This encourages the proposal distribution q to be less concentrated than the single-sample variational posterior.

1 2 **10.5.1.1 Pathologies of optimizing the IWAE bound**

3 Unfortunately, increasing the number of samples in the IWAE bound can decrease the signal to noise
4 ratio, resulting in learning a worse model [Rai+18a]. Intuitively, the reason this happens is that
5 increasing S reduces the dependence of the bound on the quality of the inference network, which
6 makes the gradient of the ELBO wrt ϕ less informative (higher variance).

7 One solution to this is to use the **doubly reparameterized gradient estimator** [TL18b].
8 Another approach is to use alternative estimation methods that avoid ELBO maximization, such
9 as using the thermodynamic variational objective (see Section 10.5.2) or the reweighted wake-sleep
10 algorithm (see Section 10.6).
11

12 13 **10.5.2 The thermodynamic variational objective (TVO)**

14 In [MLW19; Bre+20b], they present the **thermodynamic variational objective** or **TVO**. This
15 is an alternative to IWAE for creating tighter variational bounds, which has certain advantages,
16 particularly for posteriors that are not reparameterizable (e.g., discrete latent variables). The
17 framework also has close connections with the reweighted wake-sleep algorithm from Section 10.6, as
18 we will see in Section 10.5.3.

19 The TVO technique uses **thermodynamic integration**, also called **path sampling**, which is
20 a technique used in physics and phylogenetics to approximate intractable normalization constants
21 of high dimensional distributions (see e.g., [GM98; LP06; FP08]). This is based on the insight
22 that it is easier to calculate the ratio of two unknown constants than to calculate the constants
23 themselves. This is similar to the idea behind annealed importance sampling (Section 11.5.4), but TI
24 is deterministic. For details, see [MLW19; Bre+20b].
25

26 27 **10.5.3 Minimizing the evidence upper bound**

28 Recall that the evidence lower bound or ELBO is given by
29

$$\underline{30} \quad L(\boldsymbol{\theta}, \phi | \mathbf{x}) = \log p_{\boldsymbol{\theta}}(\mathbf{x}) - D_{\text{KL}}(q_{\phi}(\mathbf{z} | \mathbf{x}) \| p_{\boldsymbol{\theta}}(\mathbf{z} | \mathbf{x})) \leq \log p_{\boldsymbol{\theta}}(\mathbf{x}) \quad (10.170)$$

31 By analogy, we can define the **evidence upper bound** or **EUBO** as follows:

$$\underline{34} \quad \text{EUBO}(\boldsymbol{\theta}, \phi | \mathbf{x}) = \log p_{\boldsymbol{\theta}}(\mathbf{x}) + D_{\text{KL}}(p_{\boldsymbol{\theta}}(\mathbf{z} | \mathbf{x}) \| q_{\phi}(\mathbf{z} | \mathbf{x})) \geq \log p_{\boldsymbol{\theta}}(\mathbf{x}) \quad (10.171)$$

35 Minimizing this wrt the variational parameters ϕ , as an alternative to maxmimizing the ELBO, was
36 proposed in [MLW19], where they showed that it can sometimes converge to the true $\log p_{\boldsymbol{\theta}}(\mathbf{x})$ faster.

38 The above bound is for a specific input \mathbf{x} . If we sample \mathbf{x} from the generative model, and
39 minimize $\mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{x})} [\text{EUBO}(\boldsymbol{\theta}, \phi | \mathbf{x})]$ wrt ϕ , we recover the sleep phase of the wake-sleep algorithm (see
40 Section 10.6.2).

41 Now suppose we sample \mathbf{x} from the empirical distribution, and minimize $\mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [\text{EUBO}(\boldsymbol{\theta}, \phi | \mathbf{x})]$
42 wrt ϕ . To approximate the expectation, we can use self-normalized importance sampling, as in
43 Equation (10.188), to get

$$\underline{45} \quad \nabla_{\phi} \text{EUBO}(\boldsymbol{\theta}, \phi | \mathbf{x}) = \sum_{s=1}^S \bar{w}_s \nabla_{\phi} \log q_{\phi}(\mathbf{z}^s | \mathbf{x}) \quad (10.172)$$

1 where $\bar{w}_s = w^{(s)} / (\sum_{s'} w^{(s')})$, and $w^{(s)} = \frac{p(\mathbf{x}, \mathbf{z}^s)}{q(\mathbf{z}^s | \phi_t)}$. This is equivalent to the “daydream” update (aka
 2 “wake-phase ϕ update”) of the wake-sleep algorithm (see Section 10.6.3).
 3

4 5 10.6 Wake-sleep algorithm 6

7 So far in this chapter we have focused on fitting latent variable models by maximizing the ELBO.
 8 This has two main drawbacks. First, it does not work well when we have discrete latent variables,
 9 because in such cases we cannot use the reparameterization trick; thus we have to use higher variance
 10 estimators, such as REINFORCE (see Section 10.2.3). Second, even in the case where we can use
 11 the reparameterization trick, the lower bound may not be very tight. We can improve the tightness
 12 by using the IWAE multi-sample bound (Section 10.5.1), but paradoxically this may not result in
 13 learning a better model, for reasons discussed in Section 10.5.1.1.

14 In this section, we discuss a different way to jointly train generative and inference models, which
 15 avoids some of the problems with ELBO maximization. The method is known as the **wake-sleep**
 16 **algorithm** [Hin+95; BB15b; Le+19; FT19]. because it alternates between two steps: in the wake
 17 phase, we optimize the generative model parameters θ to maximize the marginal likelihood of the
 18 observed data (we approximate $\log p_\theta(\mathbf{x})$ by drawing importance samples from the inference network),
 19 and in the sleep phase, we optimize the inference model parameters ϕ to learn to invert the generative
 20 model by training the inference network on labeled (\mathbf{x}, \mathbf{z}) pairs, where \mathbf{x} are samples generated by
 21 the current model parameters. This can be viewed as a form of **adaptive importance sampling**,
 22 which iteratively improves its proposal, while simultaneously optimizing the model. We give further
 23 details below.
 24

25 10.6.1 Wake phase 26

27 In the **wake phase**, we minimize the KL divergence from the empirical distribution to the model’s
 28 distribution:
 29

$$30 \quad \mathcal{L}(\theta) = D_{\text{KL}}(p_{\mathcal{D}}(\mathbf{x}) \| p_{\theta}(\mathbf{x})) = \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})}[-\log p_{\theta}(\mathbf{x})] + \text{const} \quad (10.173)$$

31 where $p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{z}) p_{\theta}(\mathbf{x} | \mathbf{z}) d\mathbf{z}$. This is equivalent to maximizing the likelihood of the observed
 32 data:
 33

$$34 \quad \ell(\theta) = \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})}[\log p_{\theta}(\mathbf{x})] \quad (10.174)$$

36 Since the log marginal likelihood $\log p_{\theta}(\mathbf{x})$ cannot be computed exactly, we will approximate it.
 37 In the original wake-sleep paper, they proposed to use the ELBO lower bound. In the **reweighted**
 38 **wake-sleep** (RWS) algorithm of [BB15b; Le+19], they propose to use the IWAE bound from
 39 Section 10.5.1 instead. In particular, if we draw S samples from the inference network, $\mathbf{z}_s \sim q_{\phi}(\mathbf{z} | \mathbf{x})$,
 40 we get the following estimator:
 41

$$42 \quad \ell(\theta | \phi, \mathbf{x}) = \log \left(\frac{1}{S} \sum_{s=1}^S w_s \right) \quad (10.175)$$

45 where $w_s = \frac{p_{\theta}(\mathbf{x}, \mathbf{z}_s)}{q_{\phi}(\mathbf{z}_s | \mathbf{x})}$. Note that this is the same as the IWAE bound in Equation (10.168).
 46

1 We now discuss how to compute the gradient of this objective wrt θ or ϕ . Using the log-derivative
2 trick, we have that
3

4

$$\nabla \log w_s = \frac{1}{w_s} \nabla w_s \quad (10.176)$$

5

6 Hence

7

$$\nabla \ell(\theta|\phi, \mathbf{x}) = \frac{1}{\sum_{s=1}^S w_s} \left(\frac{1}{S} \sum_{s=1}^S \nabla w_s \right) \quad (10.177)$$

8

9

$$= \frac{1}{\sum_{s=1}^S w_s} \left(\sum_{s=1}^S w_s \nabla \log w_s \right) \quad (10.178)$$

10

11

$$= \sum_{s=1}^S \bar{w}_s \nabla \log w_s \quad (10.179)$$

12

13 where $\bar{w}_s = w_s / (\sum_{s'=1}^S w_{s'})$.

14 In the case of the derivatives wrt θ , we have

15

$$\nabla_\theta \log w_s = \frac{1}{w_s} \nabla_\theta w_s = \frac{q_\phi(z_s|\mathbf{x})}{p_\theta(\mathbf{x}, z_s)} \nabla_\theta \frac{p_\theta(\mathbf{x}, z_s)}{q_\phi(z_s|\mathbf{x})} = \frac{1}{p_\theta(\mathbf{x}, z_s)} \nabla_\theta p_\theta(\mathbf{x}, z_s) = \nabla_\theta \log p_\theta(\mathbf{x}, z_s) \quad (10.180)$$

16

17 and hence we get

18

$$\nabla_\theta \ell(\theta|\phi, \mathbf{x}) \sum_{s=1}^S \bar{w}_s \nabla \log p_\theta(\mathbf{x}, z_s) \quad (10.181)$$

19

20 10.6.2 Sleep phase

21 In the **sleep phase**, we try to minimize the KL divergence between the true posterior (under the
22 current model) and the inference network's approximation to that posterior:
23

24

$$\mathcal{L}(\phi) = \mathbb{E}_{p_\theta(\mathbf{x})} [D_{\text{KL}}(p_\theta(z|\mathbf{x}) \| q_\phi(z|\mathbf{x}))] = \mathbb{E}_{p_\theta(\mathbf{z}, \mathbf{x})} [-\log q_\phi(z|\mathbf{x})] + \text{const} \quad (10.182)$$

25

26 Equivalently, we can maximize the following log likelihood objective:

27

$$\ell(\phi|\theta) = \mathbb{E}_{(\mathbf{z}, \mathbf{x}) \sim p_\theta(\mathbf{z}, \mathbf{x})} [\log q_\phi(z|\mathbf{x})] \quad (10.183)$$

28

29 where $p_\theta(\mathbf{z}, \mathbf{x}) = p_\theta(\mathbf{z})p_\theta(\mathbf{x}|\mathbf{z})$. We see that the sleep phase amounts to maximum likelihood training
30 of the inference network based on samples from the generative model. These “fantasy samples”,
31 created while the network “dreams”, can be easily generated using ancestral sampling (Section 4.2.5).
32 If we use S such samples, the objective becomes

33

$$\ell(\phi|\theta) = \frac{1}{S} \sum_{s=1}^S \log q_\phi(z'_s|\mathbf{x}'_s) \quad (10.184)$$

34

1 where $(\mathbf{z}'_s, \mathbf{x}'_s) \sim p_{\theta}(\mathbf{z}, \mathbf{x})$. The gradient of this is given by
 2

$$3 \quad 4 \quad \nabla_{\phi} \ell(\phi | \theta) = \frac{1}{S} \sum_{s=1}^S \nabla_{\phi} \log q_{\phi}(\mathbf{z}'_s | \mathbf{x}'_s) \quad (10.185)$$

5

6 We do not require $q_{\phi}(\mathbf{z}' | \mathbf{x})$ to be reparameterizable, since the samples are drawn from a distribution
 7 that is independent of ϕ . This means it is easy to apply this method to models with discrete latent
 8 variables.
 9

10.6.3 Daydream phase

13 The disadvantage of the sleep phase is that the inference network, $q_{\phi}(\mathbf{z} | \mathbf{x})$, is trying to follow a
 14 moving target, $p_{\theta}(\mathbf{z} | \mathbf{x})$. Furthermore, it is only being trained on synthetic data from the model,
 15 not on real data. The reweighted wake-sleep algorithm of [BB15b] proposed to learn the inference
 16 network by using real data from the empirical distribution, in addition to fantasy data. They call
 17 the case where you use real data the “**wake-phase q update**”, but we will call it the “**daydream**
 18 **phase**”, since, unlike sleeping, the system uses real data \mathbf{x} to update the inference model, instead of
 19 fantasies.¹ [Le+19] went further, and proposed to only use the wake and daydream phases, and to
 20 skip the sleep phase entirely.

21 In more detail, the new objective which we want to minimize becomes

$$22 \quad 23 \quad \mathcal{L}(\phi | \theta) = \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [D_{\text{KL}}(p_{\theta}(\mathbf{z} | \mathbf{x}) \| q_{\phi}(\mathbf{z} | \mathbf{x}))] \quad (10.186)$$

24 We can compute a single sample approximation to the negative of the above expression as follows:
 25

$$26 \quad 27 \quad \ell(\phi | \theta, \mathbf{x}) = \mathbb{E}_{p_{\theta}(\mathbf{z} | \mathbf{x})} [\log q_{\phi}(\mathbf{z} | \mathbf{x})] \quad (10.187)$$

28 where $\mathbf{x} \sim p_{\mathcal{D}}$. We can approximate this expectation using importance sampling, with q_{ϕ} as the
 29 proposal. This results in the following estimator of the gradient for each datapoint:
 30

$$31 \quad 32 \quad \nabla_{\phi} \ell(\phi | \theta, \mathbf{x}) = \int p_{\theta}(\mathbf{z} | \mathbf{x}) \nabla_{\phi} \log q_{\phi}(\mathbf{z} | \mathbf{x}) d\mathbf{z} \approx \sum_{s=1}^S \bar{w}_s \nabla_{\phi} \log q_{\phi}(\mathbf{z}_s | \mathbf{x}) \quad (10.188)$$

33

34 where $\mathbf{z}_s \sim q_{\phi}(\mathbf{z}_s | \mathbf{x})$ and \bar{w}_s are the normalized weights.
 35

36 We see that Equation (10.188) is very similar to Equation (10.185). The key difference is that in
 37 the daydream phase, we sample from $(\mathbf{x}, \mathbf{z}_s) \sim p_{\mathcal{D}}(\mathbf{x})q_{\phi}(\mathbf{z} | \mathbf{x})$, where \mathbf{x} is a real datapoint, whereas
 38 in the sleep phase, we sample from $(\mathbf{x}'_s, \mathbf{z}'_s) \sim p_{\theta}(\mathbf{z}, \mathbf{x})$, where \mathbf{x}'_s is generated datapoint.
 39

10.6.4 Summary of algorithm

41 We summarize the RWS algorithm in Algorithm 10.5. The disadvantage of the RWS algorithm is
 42 that it does not optimize a single well-defined objective, so it is not clear if the method will converge,
 43 in contrast to ELBO maximization. On the other hand, the method is fairly simple, since it consists
 44 of two alternating weighted maximum likelihood problems. It can also be shown to “sandwich” a
 45

46 1. We thank Rif A. Saurous for suggesting this term.
 47

Algorithm 10.5: One SGD update using wake-sleep algorithm.

- 1 Sample \mathbf{x}_n from dataset
- 2 Draw S samples from inference network: $\mathbf{z}_s \sim q(\mathbf{z}|\mathbf{x}_n)$
- 3 Compute unnormalized weights: $w_s = \frac{p(\mathbf{x}_n, \mathbf{z}_s)}{q(\mathbf{z}_s|\mathbf{x}_n)}$
- 4 Compute normalized weights: $\bar{w}_s = \frac{w_s}{\sum_{s'=1}^S w_{s'}}$
- 5 Optional: Compute estimate of log likelihood: $\log p(\mathbf{x}_n) = \log(\frac{1}{S} \sum_{s=1}^S w_s)$
- 6 Wake phase: Update $\boldsymbol{\theta}$ using $\sum_{s=1}^S \bar{w}_s \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{z}_s, \mathbf{x}_n)$
- 7 Daydream phase: Update $\boldsymbol{\phi}$ using $\sum_{s=1}^S \bar{w}_s \nabla_{\boldsymbol{\phi}} \log q_{\boldsymbol{\phi}}(\mathbf{z}_s|\mathbf{x}_n)$
- 8 Optional sleep phase: Draw S samples from model, $(\mathbf{x}'_s, \mathbf{z}'_s) \sim p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})$ and update $\boldsymbol{\phi}$ using $\frac{1}{S} \sum_{s=1}^S \nabla_{\boldsymbol{\phi}} \log q_{\boldsymbol{\phi}}(\mathbf{z}'_s|\mathbf{x}'_s)$
- 9 b

lower and upper bound of the log marginal likelihood. We can think of this in terms of the two joint distributions $p_{\theta}(x, z) = p_{\theta}(z)p_{\theta}(x|z)$ and $q_{D,\phi}(x, z) = p_D(x)q_{\phi}(z|x)$:

$$\text{wake phase } \min_{\theta} D_{\text{KL}}(q_{\mathcal{D}, \phi}(x, z) \| p_{\theta}(x, z)) \quad (10.189)$$

$$\text{daydream phase } \min_{\phi} D_{\text{KL}}(p_{\theta}(x, z) \| q_{\mathcal{D}, \phi}(x, z)) \quad (10.190)$$

25 10.7 Expectation propagation (EP)

One problem with lower bound maximization (i.e., standard VI) is that we are minimizing $D_{\text{KL}}(q \parallel p)$, which induces **zero-forcing** behavior, as we discussed in Section 5.1.4.1. This means that $q(\mathbf{z}|\mathbf{x})$ tends to be too compact (over-confident), to avoid the situation in which $q(\mathbf{z}|\mathbf{x}) > 0$ but $p(\mathbf{z}|\mathbf{x}) = 0$, which would incur infinite KL penalty.

Although zero-forcing can be desirable behavior for some multi-modal posteriors (e.g., mixture models), it is not so reasonable for many unimodal posteriors (e.g., Bayesian logistic regression, or GPs with log-concave likelihoods). One way to avoid this problem is to minimize $D_{\text{KL}}(p \parallel q)$, which is zero-avoiding, as we discussed in Section 5.1.4.1. This tends to result in broad posteriors, which avoids overconfidence. In this section, we discuss **expectation propagation** or EP [Min01b], which can be seen as a local approximation to $D_{\text{KL}}(p \parallel q)$.

38 10.7.1 Algorithm

We assume the exact posterior can be written as follows:

$$p(\boldsymbol{\theta}|\mathcal{D}) = \frac{1}{Z_p} \hat{p}(\boldsymbol{\theta}), \quad \hat{p}(\boldsymbol{\theta}) = p_0(\boldsymbol{\theta}) \prod_{k=1}^K f_k(\boldsymbol{\theta}) \quad (10.191)$$

44
45 where $\hat{p}(\theta)$ is the unnormalized posterior, p_0 is the prior, f_k corresponds to the k 'th likelihood term
46 or **local factor** (also called a **site potential**). Here $Z_p = p(\mathcal{D})Z_0$ is the normalization constant for
47

the posterior, where Z_0 is the normalization constant for the prior. To simplify notation, we let $f_0(\boldsymbol{\theta}) = p_0(\boldsymbol{\theta})$ be the prior.

We will approximate the posterior as follows:

$$q(\boldsymbol{\theta}) = \frac{1}{Z_q} \hat{q}(\boldsymbol{\theta}), \quad \hat{q}(\boldsymbol{\theta}) = p_0(\boldsymbol{\theta}) \prod_{k=1}^K \tilde{f}_k(\boldsymbol{\theta}) \quad (10.192)$$

where $\tilde{f}_k \in \mathcal{Q}$ is the approximate local factor, and \mathcal{Q} is some tractable family in the exponential family, usually a Gaussian [Gel+14b].

We will optimize each \tilde{f}_i in turn, keeping the others fixed. We initialize each \tilde{f}_i using an uninformative distribution from the family \mathcal{Q} , so $q(\boldsymbol{\theta}) = p_0(\boldsymbol{\theta})$.

To compute the new local factor \tilde{f}_i^{new} , we proceed as follows. First we compute the **cavity distribution** by deleting the \tilde{f}_i from the approximate posterior by dividing it out:

$$q_{-i}^{\text{cavity}}(\boldsymbol{\theta}) = \frac{q(\boldsymbol{\theta})}{\tilde{f}_i(\boldsymbol{\theta})} \propto \prod_{k \neq i} \tilde{f}_k(\boldsymbol{\theta}) \quad (10.193)$$

This division operation can be implemented by subtracting the natural parameters, as explained in Section 2.3.3.2. The cavity distribution represents the effect of all the factors except for f_i (which is approximated by \tilde{f}_i).

Next we (conceptually) compute the **tilted distribution** by multiplying the exact factor f_i onto the cavity distribution:

$$q_i^{\text{tilted}}(\boldsymbol{\theta}) = \frac{1}{Z_i} f_i(\boldsymbol{\theta}) q_{-i}^{\text{cavity}}(\boldsymbol{\theta}) \quad (10.194)$$

where $Z_i = \int q_{-i}^{\text{cavity}}(\boldsymbol{\theta}) f_i(\boldsymbol{\theta}) d\boldsymbol{\theta}$ is the normalization constant for the tilted distribution. This is the result of combining the current approximation, excluding factor i , with the exact f_i term.

Unfortunately, the resulting tilted distribution may be outside of our model family (e.g., if we combine a Gaussian prior with a non-Gaussian likelihood). So we will approximate the tilted distribution as follows:

$$q_i^{\text{proj}}(\boldsymbol{\theta}) = \text{proj}(q_i^{\text{tilted}}) \triangleq \underset{\tilde{q} \in \mathcal{Q}}{\operatorname{argmin}} D(q_i^{\text{tilted}} || \tilde{q}) \quad (10.195)$$

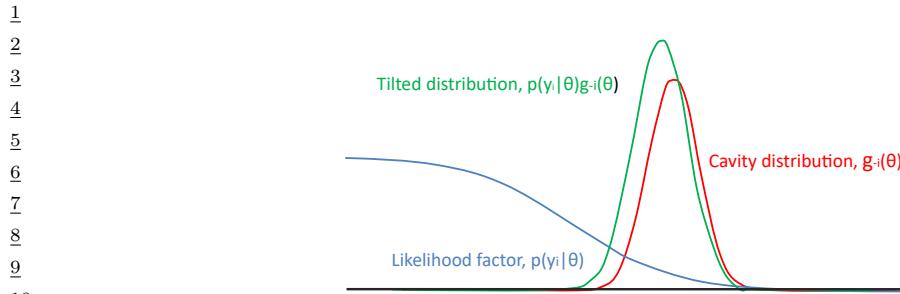
This can be thought of as projecting the tilted distribution into the approximation family. If $D(q_i^{\text{tilted}} || q) = D_{\text{KL}}(q_i^{\text{tilted}} || q)$, this can be done by moment matching, as shown in Section 5.1.4.2.

For example, suppose the cavity distribution is Gaussian, $q_{-i}^{\text{cavity}}(\boldsymbol{\theta}) = \mathcal{N}_c(\boldsymbol{\theta} | \mathbf{r}_{-i}, \mathbf{Q}_{-i})$, using the canonical parameterization. Then the log of the tilted distribution is given by

$$\log q_i^{\text{tilted}}(\boldsymbol{\theta}) = \alpha \log f_i(\boldsymbol{\theta}) - \frac{1}{2} \boldsymbol{\theta}^\top \mathbf{Q}_{-i} \boldsymbol{\theta} + \mathbf{r}_{-i}^\top \boldsymbol{\theta} + \text{const} \quad (10.196)$$

Let $\hat{\boldsymbol{\theta}}$ be a local maximum of this objective. If \mathcal{Q} is the set of Gaussians, we can compute the projected tilted distribution as a Gaussian with the following parameters:

$$\mathbf{Q}_{\setminus i} = -\nabla_{\boldsymbol{\theta}}^2 \log q_i^{\text{tilted}}(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}}, \quad \mathbf{r}_{\setminus i} = \mathbf{Q}_{\setminus i} \hat{\boldsymbol{\theta}} \quad (10.197)$$



11 *Figure 10.15: Combining a logistic likelihood factor $f_i = p(y_i | \theta)$ with the cavity prior, $q_{-i}^{cavity} = g_{-i}(\theta)$, to get*
 12 *the tilted distribution, $q_i^{tilted} = p(y_i | \theta)g_{-i}(\theta)$. Adapted from Figure 2 of [Gel+14b].*

16 This is called **Laplace propagation** [SVE04]. For more general distributions, we can use Monte
 17 Carlo approximations; this is known as **blackbox EP** [HL+16a; Li+18c].

18 Finally, we compute a local factor that, if combined with the cavity distribution, would give the
 19 same results as this projected distribution:

$$21 \quad \tilde{f}_i^{\text{new}}(\theta) = \frac{q_i^{\text{proj}}(\theta)}{q_{-i}^{\text{cavity}}(\theta)} \quad (10.198)$$

24 We see that $q_{-i}^{\text{cavity}}(\theta)\tilde{f}_i^{\text{new}}(\theta) = q_i^{\text{proj}}(\theta)$, so combining this approximate factor with the cavity
 25 distribution results in a distribution which is the best possible approximation (within \mathcal{Q}) to the
 26 results of using the exact factor.

27

28 10.7.2 Example

30 Figure 10.15 illustrates the process of combining a very non-Gaussian likelihood f_i with a Gaussian
 31 cavity prior q_{-i}^{cavity} to yield a nearly Gaussian tilted distribution q_i^{tilted} , which can then be approximated
 32 by a Gaussian using projection.

33 Thus instead of trying to “Gaussianize” each likelihood term f_i in isolation (as is done, e.g., in
 34 EKF), we try to find the best local factor \tilde{f}_i (within some family) that achieves approximately the
 35 same effect, when combined with all the other terms (represented by the cavity distribution, q_{-i}), as
 36 using the exact factor f_i . That is, we choose a local factor that works well in the context of all the
 37 other factors.

38

39 10.7.3 EP as generalized ADF

40 We can view EP as a generalization of the ADF algorithm discussed in Section 8.6. ADF is a
 41 form of sequential Bayesian inference. At each step, it maintains a tractable approximation to
 42 the posterior, $q_t(z) \in \mathcal{Q}$, updates it with the likelihood from the next observation, $\hat{p}_{t+1}(z) \propto$
 43 $q_t(z)p(x_t | z)$, and then projects the resulting updated posterior back to the tractable family using
 44 $q_{t+1} = \operatorname{argmin}_{q \in \mathcal{Q}} D_{\text{KL}}(\hat{p}_{t+1} \| q)$. ADF minimizes KL in the desired direction. However, it is a
 45 sequential algorithm, designed for the online setting. In the batch setting, the method can given
 46

47

different results depending on the order in which the updates are performed. In addition, if we perform multiple passes over the data, we will include the same likelihood terms multiple times, resulting in an overconfident posterior. EP overcomes this problem.

10.7.4 Optimization issues

In practice, EP can be numerically unstable. For example, if we use Gaussians as our local factors, we might end up with negative variance when we subtract the natural parameters. To reduce the chance of this, it is common to use damping, in which we perform a partial update of each factor with a step size of δ . More precisely, we change the final step to be the following:

$$\tilde{f}_i^{\text{new}}(\boldsymbol{\theta}) = \left(\tilde{f}_i(\boldsymbol{\theta}) \right)^{1-\delta} \left(\frac{q_i^{\text{proj}}(\boldsymbol{\theta})}{q_{-i}^{\text{cavity}}} \right)^\delta \quad (10.199)$$

This can be implemented by scaling the natural parameters by δ . [ML02] suggest $\delta = 1/K$ as a safe strategy (where K is the number of factors), but this results in very slow convergence. [Gel+14b] suggest starting with $\delta = 0.5$, and then reducing to $\delta = 1/K$ over K iterations.

In addition to numerical stability, there is no guarantee that EP will converge in its vanilla form, although empirically it can work well, especially with log-concave factors f_i (e.g., as in GP classifiers).

10.7.5 Power EP and α -divergence

We also have a choice about what divergence measure $D(q_i^{\text{tilted}} || q)$ to use when we approximate the tilted distribution. If we use $D_{\text{KL}}(q_i^{\text{tilted}} || q)$, we recover classic EP, as described above. If we use $D_{\text{KL}}(q || q_i^{\text{tilted}})$, we recover the reverse KL used in standard variational inference. We can generalize the above results by using α -divergences (Section 2.7.1.2), which allow us to interpolate between mode seeking and mode covering behavior, as shown in Figure 2.20. We can optimize the α -divergence by using the **power EP** method of [Min04].

Algorithmically, this is a fairly small modification to regular EP. In particular, we first compute the cavity distribution, $q_{-i}^{\text{cavity}} \propto \frac{q}{f_i^\alpha}$; we then approximate the tilted distribution, $q_i^{\text{proj}} = \text{proj}(q_{-i}^{\text{cavity}} f_i^\alpha)$;

and finally we compute the new factor $\tilde{f}_i^{\text{new}} \propto \left(\frac{q_i^{\text{proj}}}{q_{-i}^{\text{cavity}}} \right)^{1/\alpha}$.

10.7.6 Stochastic EP

The main disadvantage of EP in the big data setting is that we need to store the $\tilde{f}_n(\boldsymbol{\theta})$ terms for each datapoint n , so we can compute the cavity distribution. If $\boldsymbol{\theta}$ has D dimensions, and we use full covariance Gaussians, this requires $O(ND^2)$ memory.

The idea behind **stochastic EP** [LHLT15] is to approximate the local factors with a shared factor that acts like an aggregated likelihood, i.e.,

$$\prod_{n=1}^N f_n(\boldsymbol{\theta}) \approx \tilde{f}(\boldsymbol{\theta})^N \quad (10.200)$$

1 where typically $f_n(\boldsymbol{\theta}) = p(\mathbf{x}_n | \boldsymbol{\theta})$. This exploits the fact that the posterior only cares about approximating the product of the likelihoods, rather than each likelihood separately. Hence it suffices for $\tilde{f}(\boldsymbol{\theta})$ to approximate the average likelihood.

2 We can modify EP to this setting as follows. First, when computing the cavity distribution, we
3 use

4

$$q_{-1}(\boldsymbol{\theta}) \propto q(\boldsymbol{\theta}) / \tilde{f}(\boldsymbol{\theta}) \quad (10.201)$$

5 We then compute the tilted distribution
6

7

$$q_{\setminus n}(\boldsymbol{\theta}) \propto f_n(\boldsymbol{\theta}) q_{-1}(\boldsymbol{\theta}) \quad (10.202)$$

8 Next we derive the new local factor for this datapoint using moment matching:

9

$$\tilde{f}_n(\boldsymbol{\theta}) = \text{proj}(q_{\setminus n}(\boldsymbol{\theta})) / q_{-1}(\boldsymbol{\theta}) \quad (10.203)$$

10 Finally, we perform a damped update of the average likelihood $\tilde{f}(\boldsymbol{\theta})$ using this new local factor:

11

$$\tilde{f}_{\text{new}}(\boldsymbol{\theta}) = \tilde{f}_{\text{old}}(\boldsymbol{\theta})^{1-1/N} \tilde{f}_n(\boldsymbol{\theta})^{1/N} \quad (10.204)$$

12 The ADF algorithm is similar to SEP, in that we compute the tilted distribution $q_{\setminus t} \propto f_t q_{t-1}$ and
13 then project it, without needing to keep the f_t factors. The difference is that instead of using the
14 cavity distribution $q_{-1}(\boldsymbol{\theta})$ as a prior, it uses the posterior from the previous time step, q_{t-1} . This
15 avoids the need to compute and store \tilde{f} , but results in overconfidence in the batch setting.

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

11 Monte Carlo methods

11.1 Introduction

In this chapter, we discuss **Monte Carlo methods**, which are a stochastic approach to solving numerical integration problems. The name refers to the “Monte Carlo” casino in Monaco; this was used as a codename by von Neumann and Ulam, who invented the technique while working on the atomic bomb during WWII. Since then, the technique has become widely adopted in physics, statistics, machine learning, and many areas of science and engineering.

In this chapter, we give a brief introduction to some key concepts. In Chapter 12, we discuss MCMC, which is the most widely used MC method for high-dimensional problems. In Chapter 13, we discuss SMC, which is widely used for MC inference in state space models, but can also be applied more generally. For more details on MC methods, see e.g., [Liu01; RC04; KTB11; BZ20].

11.2 Monte Carlo integration

We often want to compute the expected value of some function of a random variable, $\mathbb{E}[f(\mathbf{X})]$. This requires computing the following integral:

$$\mathbb{E}[f(\mathbf{x})] = \int f(\mathbf{x})p(\mathbf{x})d\mathbf{x} \tag{11.1}$$

where $\mathbf{x} \in \mathbb{R}^n$, $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$, and $p(\mathbf{x})$ is the target distribution of \mathbf{X} .¹ In low dimensions (up to, say, 3), we can compute the above integral efficiently using **numerical integration**, which (adaptively) computes a grid, and then evaluates the function at each point on the grid.² But this does not scale to higher dimensions.

An alternative approach is to draw multiple random samples, $\mathbf{x}_n \sim p(\mathbf{x})$, and then to compute

$$\mathbb{E}[f(\mathbf{x})] \approx \frac{1}{N_s} \sum_{n=1}^{N_s} f(\mathbf{x}_n) \tag{11.2}$$

This is called **Monte Carlo integration**. It has the advantage over numerical integration that the function is only evaluated in places where there is non-negligible probability, so it does not

1. In many cases, the target distribution may be the posterior $p(\mathbf{x}|\mathbf{y})$, which can be hard to compute; in such problems, we often work with the unnormalized distribution, $\tilde{p}(\mathbf{x}) = p(\mathbf{x}, \mathbf{y})$, instead, and then normalize the results using $Z = \int p(\mathbf{x}, \mathbf{y})d\mathbf{x} = p(\mathbf{y})$.

2. In 1d, numerical integration is called **quadrature**; in higher dimensions, it is called **cubature** [Sar13].