

1 where  
 2

$$3 \quad \mathbf{m} = \int_{\mathbf{x}} \mathbf{g}(\mathbf{x}) p(\mathbf{x}) \quad (8.211a)$$

$$4 \quad \mathbf{S} = \int_{\mathbf{x}} (\mathbf{g}(\mathbf{x}) - \mathbf{m})(\mathbf{g}(\mathbf{x}) - \mathbf{m})^T p(\mathbf{x}) \quad (8.211b)$$

$$5 \quad \mathbf{C} = \int_{\mathbf{x}} (\mathbf{x} - \boldsymbol{\mu})(\mathbf{g}(\mathbf{x}) - \mathbf{m})^T p(\mathbf{x}) \quad (8.211c)$$

6 We can compute these expectations with a variety of methods, as we discuss in Section 8.8.3. See  
 7 Algorithm 14 for the pseudocode.

8 Once we have estimated  $\boldsymbol{\theta} = (\mathbf{A}, \mathbf{b}, \boldsymbol{\Omega})$ , the joint distribution follows from Equation (2.81):

$$9 \quad \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} \boldsymbol{\mu} \\ \mathbf{A}\boldsymbol{\mu} + \mathbf{b} \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma} & \boldsymbol{\Sigma}\mathbf{A}^T \\ \mathbf{A}\boldsymbol{\Sigma} & \boldsymbol{\Omega} + \boldsymbol{\Omega} + \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T \end{pmatrix} \right) = \mathcal{N} \left( \begin{pmatrix} \boldsymbol{\mu} \\ \mathbf{m} \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma} & \mathbf{C} \\ \mathbf{C}^T & \boldsymbol{\Omega} + \mathbf{S} \end{pmatrix} \right) \quad (8.212)$$

10 We see that this matches the moments of  $p(\mathbf{x}, \mathbf{y})$ , and hence minimizes the KL divergence.

---

11 **Algorithm 14:** Statistical linear regression.

---

```
12 1 def SLR( $\boldsymbol{\mu}$ ,  $\boldsymbol{\Sigma}$ ,  $\mathbf{m}$ ,  $\mathbf{S}$ ,  $\mathbf{C}$ ) :
 2   2    $\hat{\mathbf{A}} = \mathbf{C}^T \boldsymbol{\Sigma}^{-1}$ 
 3   3    $\hat{\mathbf{b}} = \mathbf{m} - \hat{\mathbf{A}}\boldsymbol{\mu}$ 
 4   4    $\hat{\boldsymbol{\Omega}} = \mathbf{S} - \hat{\mathbf{A}}\boldsymbol{\Sigma}\hat{\mathbf{A}}^T$ 
 5   5   Return ( $\hat{\mathbf{A}}$ ,  $\hat{\mathbf{b}}$ ,  $\hat{\boldsymbol{\Omega}}$ )
```

---

## 8.9.2 Prior linearization filter

We now discuss how to use SLR to perform online filtering. When computing the expectations, we use the prior predictive distribution. Hence this is called the **prior linearization filter** or **PrLF** [GFSS17]. We will use other distributions for the expectations in later sections.

In the prediction step, we approximate  $p(\mathbf{z}_{t-1}, \mathbf{z}_t | \mathbf{y}_{1:t-1})$  and then we compute the marginal distribution  $p(\mathbf{z}_t | \mathbf{y}_{1:t-1})$ . For notational simplicity, we use the UKF in the predict step:

$$34 \quad (\boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}, -) = \text{UnscentedPredict}(\boldsymbol{\mu}_{t-1|t-1}, \boldsymbol{\Sigma}_{t-1|t-1}, \mathbf{f}(\cdot, \mathbf{u}_t), \mathbf{Q}_t) \quad (8.213)$$

In the update step, we approximate  $p(\mathbf{z}_t, \mathbf{y}_t | \mathbf{y}_{1:t-1})$  by applying the SLR transform to the observation model, and then we perform Gaussian updating:

$$38 \quad (\mathbf{m}, \mathbf{S}, \mathbf{C}) = \text{UnscentedPredict}(\boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}, \mathbf{h}(\cdot, \mathbf{u}_t), \mathbf{0}) \quad (8.214)$$

$$39 \quad (\mathbf{H}_t, \mathbf{d}_t, \tilde{\mathbf{R}}_t) = \text{SLR}(\boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}, \mathbf{m}, \mathbf{S}, \mathbf{C}) \quad (8.215)$$

$$40 \quad (\boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t}, \ell_t) = \text{LinGaussUpdate}(\boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}, \mathbf{H}_t, \mathbf{d}_t, \tilde{\mathbf{R}}_t + \mathbf{R}_t, \mathbf{y}_t) \quad (8.216)$$

43 where  $\tilde{\mathbf{R}}_t$  is the extra covariance introduced by approximating the observation model.

44 A special case of the prior linearization filter arises if we set  $\tilde{\mathbf{R}}_t = \mathbf{0}$ , which ignores extra uncertainty  
 45 introduced by the linear approximation. This simpler case is known as the **statistically linearized**  
 46 **filter** [Gel74] (see also [Sar13, Sec 5.4]).

---

1 **8.9.3 Iterated posterior linearization filter**

3 In [GF+15] they propose a method called the **posterior linearization filter**, which uses SLR, but  
4 where the expectations for the measurement update step are taken with respect to the posterior  
5 instead of the prior. This ensures that the observation model is linearized near the location where  
6 the posterior has most of its mass.

7 Of course, this approach cannot be implemented as stated, because it requires knowledge of the  
8 posterior to approximate the posterior. However, we can approximate it by an iterative method, as  
9 shown in Algorithm 15, where  $J$  is the number of iterations. This is called the **iterated posterior**  
10 **linearization filter** or **IPLF**. (If we set  $J = 1$ , we recover the prior linearization filter.)  
11

---

12 **Algorithm 15:** Iterated posterior linearization filter.

```
14 1 def IPLF( $f, Q, h, R, y_{1:T}, J$ ) :
15 2 foreach  $t = 1 : T$  do
16 3   Predict step:
17 4      $(\mu_{t|t-1}, \Sigma_{t|t-1}, -) = \text{UnscentedPredict}(\mu_{t-1|t-1}, \Sigma_{t-1|t-1}, f, Q_t)$ 
18 5   Update step:
19 6      $\mu_{t|t} = \mu_{t|t-1}, \Sigma_{t|t} = \Sigma_{t|t-1}$ 
20 7     foreach  $j = 1 : J$  do
21 8        $(m, S, C) = \text{UnscentedPredict}(\mu_t, \Sigma_t, h(\cdot, u_t), 0)$ 
22 9        $(H_t, d_t, \tilde{R}_t) = \text{SLR}(\mu_t, \Sigma_t, m, S, C)$ 
23 10       $(\mu_t, \Sigma_t, \ell_t) = \text{LinGaussUpdate}(\mu_t, \Sigma_t, H_t, d_t, \tilde{R}_t + R_t, y_t)$ 
25 11  Return  $(\mu_{t|t}, \Sigma_{t|t})_{t=1}^T, \sum_{t=1}^T \ell_t$ 
26
```

---

27

28 The iterated EKF (Section 8.5.3) is similar to the IPLF, except it uses a first order Taylor series  
29 approximation of the measurement function to compute the moments (rather than a sigma point  
30 approximation), and it ignores the uncertainty due to linearization. This is much faster than IPLF.  
31 Unfortunately both IEKF and IPLF can diverge. A robust IEKF method, that uses line search to  
32 perform damped (partial) updates of  $\mu_{t|t}^j$ , is presented in [SHA15]. In [Rai+18b], a similar robust  
33 method is derived for IPLF.

34

35 **8.9.4 Iterated posterior linearization smoother**

37 In [GFSS17] they propose a method called the **iterated posterior linearization smoother** or  
38 **IPLS**, which extends the IPLF method of Section 8.9.3 to the offline setting. The basic idea is to  
39 compute expectations with respect to the previous smoothed posterior, and to iterate this process  
40 until convergence. See Algorithm 16 for the pseudocode.

41 The **iterated EKS** is similar to the IPLS, except it uses a first order Taylor series approximation  
42 of the measurement and dynamics functions at the posterior mean from the previous iteration, instead  
43 of using SLR. In [Bel94], it is shown to be equivalent to a Gauss-Newton method for computing the  
44 MAP estimate of the smoothing posterior.

45 Unfortunately both IEKS and IPLS can diverge. A robust IEKS method, that uses line search  
46 and Levenberg-Marquadt to update the parameters, is presented in [SS20a]. In [Lin+21c], a similar  
47

---

**Algorithm 16:** Iterated posterior linearization smoother

---

```

1 def IPLS( $f, \mathbf{Q}, \mathbf{h}, \mathbf{R}, \mathbf{y}_{1:T}, J$ ) :
2    $(\boldsymbol{\mu}_{t|T}, \boldsymbol{\Sigma}_{t|T})_{t=1}^T = \text{UnscentedSmoker}(\mathbf{f}, \mathbf{Q}, \mathbf{h}, \mathbf{R}, \mathbf{y}_{1:T})$ 
3   foreach  $j = 1 : J$  do
4     foreach  $t = 1 : T$  do
5        $(\mathbf{m}, \mathbf{S}, \mathbf{C}) = \text{UnscentedPredict}(\boldsymbol{\mu}_{t|T}, \boldsymbol{\Sigma}_{t|T}, \mathbf{f}(\cdot, \mathbf{u}_t), \mathbf{0})$ 
6        $(\mathbf{F}_t, \mathbf{b}_t, \tilde{\mathbf{Q}}_t) = \text{SLR}(\boldsymbol{\mu}_{t|T}, \boldsymbol{\Sigma}_{t|T}, \mathbf{m}, \mathbf{S}, \mathbf{C})$ 
7        $(\mathbf{m}, \mathbf{S}, \mathbf{C}) = \text{UnscentedPredict}(\boldsymbol{\mu}_{t|T}, \boldsymbol{\Sigma}_{t|T}, \mathbf{h}(\cdot, \mathbf{u}_t), \mathbf{0})$ 
8        $(\mathbf{H}_t, \mathbf{d}_t, \tilde{\mathbf{R}}_t) = \text{SLR}(\boldsymbol{\mu}_{t|T}, \boldsymbol{\Sigma}_{t|T}, \mathbf{m}, \mathbf{S}, \mathbf{C})$ 
9        $(\boldsymbol{\mu}_{t|T}, \boldsymbol{\Sigma}_{t|T})_{t=1}^T = \text{KalmanSmoker}(\{\mathbf{F}_t, \mathbf{b}_t, \mathbf{Q} + \tilde{\mathbf{Q}}_t, \mathbf{H}_t, \mathbf{d}_t, \mathbf{R} + \tilde{\mathbf{R}}_t\}, \mathbf{y}_{1:T})$ 
10
11
12
13
14 Return  $(\boldsymbol{\mu}_{t|T}, \boldsymbol{\Sigma}_{t|T})_{t=1}^T$ 

```

---

robust method is derived for IPLS.

In [HPR19] they extend IPLS to belief propagation in Forney factor graphs (Section 4.6.1.2), which enables the method to be applied to general graphical models with Gaussian potentials but nonlinear dependencies.

In [Kam+22], they present a method based on approximate expectation propagation (Section 10.7), that is very similar to IPLS, except that the distributions that are used to compute the SLR terms, needed to compute the Gaussian messages, are different. In particular, rather than using the smoothed posterior from the last iteration, it uses the “cavity” distributions, which is the current posterior minus the incoming message that was sent at the last iteration, similar to Section 8.3.4.4. The advantage of this is that the outgoing message does not double count the evidence. The disadvantage is that this may be numerically unstable.

### 8.9.5 Beyond additive Gaussian noise

In this section, we extend SLR to handle the case of non-Gaussian / non-additive noise, following [TGFS18].

#### 8.9.5.1 Generalized statistical linear regression

We assume  $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_X, \boldsymbol{\Sigma}_X)$ . We allow  $p(\mathbf{y}|\mathbf{x})$  to be any kind of distribution, provided it has finite variance. We only require that we can approximate the first and second **conditional moments**:

$$\mathbf{m}_Y(\mathbf{x}) = \mathbb{E}[\mathbf{y}|\mathbf{x}], \mathbf{V}_Y(\mathbf{x}) = \text{Cov}[\mathbf{y}|\mathbf{x}] \quad (8.217)$$

For example, if  $p(y|x) = \text{Poisson}(y|ce^x)$ , we have

$$m_Y(x) = V_Y(x) = ce^x \quad (8.218)$$

As usual we will approximate the conditional distribution by a linear Gaussian model of the form  $p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x} + \mathbf{b}, \boldsymbol{\Omega})$ . We can then easily compute the joint  $p(\mathbf{x}, \mathbf{y})$ , the marginal  $p(\mathbf{y})$ , and the posterior  $p(\mathbf{x}|\mathbf{y})$ .

The optimal parameters, in the sense of minimizing the mean squared norm of the residual,  $\|\mathbf{y} - \mathbf{Ax} - \mathbf{b}\|$ , and the corresponding error matrix, are given by the SLR equations in Algorithm 14, where the moments are computed using

$$\mathbf{m} = \mathbb{E}[\mathbf{Y}] = \mathbb{E}[\mathbb{E}[\mathbf{Y}|\mathbf{X}]] = \int_{\mathbf{x}} p(\mathbf{x}) \mathbf{m}_Y(\mathbf{x}) \quad (8.219)$$

$$\mathbf{S} = \mathbb{V}[\mathbf{Y}] = \mathbb{E}[\mathbb{V}[\mathbf{Y}|\mathbf{x}]] + \mathbb{V}[\mathbb{E}[\mathbf{Y}|\mathbf{X}]] \quad (8.220)$$

$$= \int_{\mathbf{x}} p(\mathbf{x}) \mathbf{V}_Y(\mathbf{x}) + \int_{\mathbf{x}} p(\mathbf{x}) (\mathbf{m}_Y(\mathbf{x}) - \mathbf{m})(\mathbf{m}_Y(\mathbf{x}) - \mathbf{m})^T \quad (8.221)$$

$$\mathbf{C} = \text{Cov}[\mathbf{X}, \mathbf{Y}] = \text{Cov}[\mathbf{X}, \mathbb{E}[\mathbf{Y}|\mathbf{X}]] = \int_{\mathbf{x}} p(\mathbf{x}) (\mathbf{x} - \boldsymbol{\mu}_X)(\mathbf{m}_Y(\mathbf{x}) - \mathbf{m}) \quad (8.222)$$

where expectations are taken wrt  $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_X, \boldsymbol{\Sigma}_X)$ . This is called **generalized statistical linear regression**.

In general, we cannot compute these expectations exactly, so we need to use numerical integration (e.g., sigma point methods). Alternatively, we can make a Taylor series approximation to the conditional moments, and then we can compute the unconditional moments analytically. In more detail, we first compute

$$\mathbf{m}_Y(\mathbf{x}) \approx \underbrace{\mathbf{m}_Y(\boldsymbol{\mu})}_{\boldsymbol{\mu}_Y} + \underbrace{\nabla \mathbf{m}_Y(\boldsymbol{\mu})}_{\mathbf{J}} (\mathbf{x} - \boldsymbol{\mu}) \quad (8.223)$$

$$\mathbf{V}_Y(\mathbf{x}) \approx \mathbf{V}_Y(\boldsymbol{\mu}) \quad (8.224)$$

We can then compute the unconditional moments using Algorithm 17.

---

**Algorithm 17:** Linearized Gaussian approximation to a conditional distribution.

---

```

1 def CondMomentsPredict( $\boldsymbol{\mu}_X, \boldsymbol{\Sigma}_X, \mathbf{m}_Y(), \mathbf{V}_Y()$ ) :
2    $\mathbf{J} = \text{Jac}(\mathbf{m}_Y)(\boldsymbol{\mu}_X)$ 
3    $\mathbf{m} = \mathbf{m}_Y(\boldsymbol{\mu}_X)$ 
4    $\mathbf{S} = \mathbf{V}_Y(\boldsymbol{\mu}_X) + \mathbf{J} \boldsymbol{\Sigma}_X \mathbf{J}^T$ 
5    $\mathbf{C} = \boldsymbol{\Sigma}_X \mathbf{J}^T$ 
6   Return ( $\mathbf{m}, \mathbf{S}, \mathbf{C}$ )

```

---

### 8.9.5.2 Generalized IPLF

Once we have computed the moments, we can use SLR to compute a linear Gaussian approximation to  $p(\mathbf{y}|\mathbf{x})$  and hence we can compute  $p(\mathbf{x}, \mathbf{y})$ , from which we can extract  $p(\mathbf{y})$  for the predict step and  $p(\mathbf{x}|\mathbf{y})$  for the update step.

As in the IPLF method of Section 8.9.3, we can get improved accuracy by iterating the computation of the moments, using the most recent posterior estimate for the expectation. We call this the **generalized IPLF** or **GIPLF**. We give the pseudocode in Algorithm 18, where we assume (for notational simplicity) that the dynamics model has additive Gaussian noise, and where we use Taylor series linearization rather than sigma point methods. If the observation model is linear Gaussian, the resulting algorithm is equivalent to the IEKF in Section 8.5.3.

47

---

**Algorithm 18:** Generalized IPLF using Taylor series linearization.

---

```

1 def GIPLF( $f, Q, m_Y(), V_Y(), y_{1:T}, J$ ) :
2   foreach  $t = 1 : T$  do
3     Predict step:
4      $(\mu_{t|t-1}, \Sigma_{t|t-1}, -) = \text{LinearizedPredict}(\mu_{t-1|t-1}, \Sigma_{t-1|t-1}, f, Q_t)$ 
5     Update step:
6      $\mu_{t|t} = \mu_{t|t-1}, \Sigma_{t|t} = \Sigma_{t|t-1}$ 
7     foreach  $j = 1 : J$  do
8        $(m, S, C) = \text{CondMomentsPredict}(\mu_t, \Sigma_t, m_Y, V_Y)$ 
9        $(H_t, d_t, R_t) = \text{SLR}(\mu_t, \Sigma_t, m, S, C)$ 
10       $(\mu_t, \Sigma_t, \ell_t) = \text{LinGaussUpdate}(\mu_t, \Sigma_t, H_t, d_t, R_t, y_t)$ 
11
12
13
14 11 Return  $(\mu_{t|t}, \Sigma_{t|t})_{t=1}^T$ 
15
16
17
18

```

---

### 8.9.5.3 Extensions

The generalized IPLF algorithm can be extended in various ways. In [Hos+20b] they use it as a proposal distribution inside of a particle filtering algorithm (Section 13.2). In [TGFS18], they extend it to the smoothing (offline) setting. In [GFTS19], they apply it to fitting a Gaussian process with a Bernoulli likelihood.

In [WSS21], they propose a variety of “**Bayes-Newton**” methods for approximately computing Gaussian posteriors to probabilistic models with nonlinear and/or non-Gaussian likelihoods. This generalizes all of the above methods, and can be applied to SSMs and GPs.

## 8.10 Assumed density filtering

In this section, we discuss **assumed density filtering** or **ADF** [May79]. In this approach, we *assume* the posterior has a specific form (e.g., a Gaussian). At each step, we update the previous posterior with the new likelihood; the result will often not have the desired form (e.g., will no longer be Gaussian), so we project it to the closest approximating distribution of the required type.

In more detail, we assume (by induction) that our prior  $q_{t-1}(z_{t-1}) \approx p(z_{t-1}|y_{1:t-1})$  satisfies  $q_{t-1} \in \mathcal{Q}$ , where  $\mathcal{Q}$  is a family of tractable distributions. We can update the prior with the new measurement to get the approximate posterior as follows. First we compute the **one-step-ahead predictive distribution**

$$p_{t|t-1}(z_t|y_{1:t-1}) = \int p(z_t|z_{t-1})q_{t-1}(z_{t-1})dz_{t-1} \quad (8.225)$$

Then we update this prior with the likelihood for step  $t$  to get the posterior

$$p_t(z_t|y_{1:t}) = \frac{1}{Z_t} p(y_t|z_t)p_{t|t-1}(z_t) \quad (8.226)$$

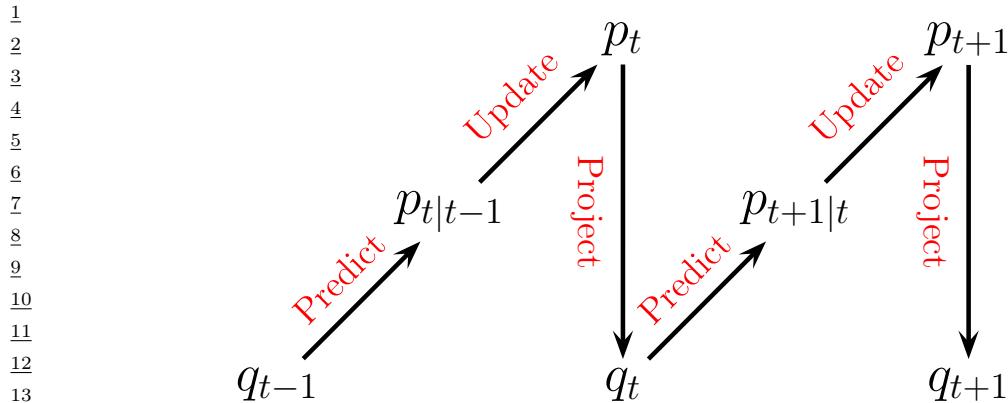


Figure 8.13: Illustration of the predict-update-project cycle of assumed density filtering.  $q_t \in \mathcal{Q}$  is a tractable distribution, whereas we may have  $p_{t|t-1} \notin \mathcal{Q}$  and  $p_t \notin \mathcal{Q}$ .

where

$$Z_t = \int p(\mathbf{y}_t | \mathbf{z}_t) p_{t|t-1}(\mathbf{z}_t) d\mathbf{z}_t \quad (8.227)$$

is the normalization constant. Unfortunately, we often find that the resulting posterior is no longer in our tractable family,  $p(\mathbf{z}_t) \notin \mathcal{Q}$ . So after Bayesian updating we seek the best tractable approximation by computing

$$q_t(\mathbf{z}_t | \mathbf{y}_{1:t}) = \underset{q \in \mathcal{Q}}{\operatorname{argmin}} D_{\text{KL}}(p_t(\mathbf{z}_t | \mathbf{y}_{1:t}) \| q(\mathbf{z}_t)) \quad (8.228)$$

This minimizes the Kullback-Leibler divergence from the approximation  $q(\mathbf{z}_t)$  to the “exact” posterior  $p_t(\mathbf{z}_t)$ , and can be thought of as **projecting**  $p$  onto the space of tractable distributions. Thus the overall algorithm consists of three steps — predict, update, and project — as sketched in Figure 8.13.

Computing  $\min_q D_{\text{KL}}(p \| q)$  is known as **moment projection**, since the optimal  $q$  should have the same moments as  $p$  (see Section 5.1.3.4). So in the Gaussian case, we just need to set the mean and covariance of  $q_t$  so they are the same as the mean and covariance of  $p_t$ . We will give some examples of this below. By contrast, computing  $\min_q D_{\text{KL}}(q \| p)$ , as in variational inference (Section 10.1), is known as **information projection**, and will result in mode seeking behavior (see Section 5.1.3.3), rather than trying to capture overall moments.

### 8.10.1 Connection with Gaussian filtering

When  $\mathcal{Q}$  is the set of Gaussian distributions, there is a close connection between ADF and Gaussian filtering, which we discussed in Section 8.8. GF corresponds to solving the following optimization problem

$$q_{t|t-1}(\mathbf{z}_t, \tilde{\mathbf{y}}_t) = \underset{q \in \mathcal{Q}}{\operatorname{argmin}} D_{\text{KL}}(p(\mathbf{z}_t, \tilde{\mathbf{y}}_t | \mathbf{y}_{1:t-1}) \| q(\mathbf{z}_t, \tilde{\mathbf{y}}_t | \mathbf{y}_{1:t-1})) \quad (8.229)$$

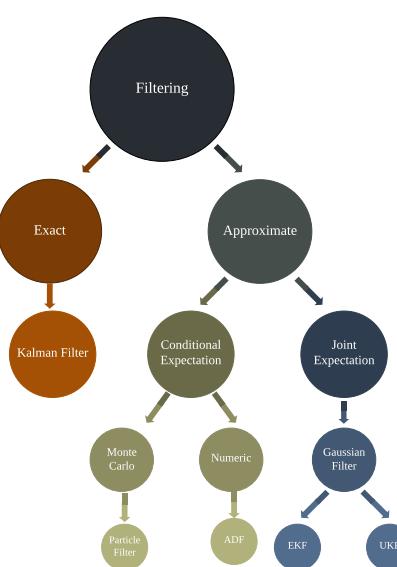


Figure 8.14: A taxonomy of filtering algorithms. Adapted from Figure 2 of [Wüt+16].

which can be solved by moment matching (see Section 8.8). We then condition this joint distribution on the event  $\tilde{\mathbf{y}}_t = \mathbf{y}_t$ , where  $\tilde{\mathbf{y}}_t$  is the unknown random variable and  $\mathbf{y}_t$  is its observed value. This gives  $p_t(\mathbf{z}_t|\mathbf{y}_{1:t})$ , which is easy to compute, due to the Gaussian assumption. By contrast, in Gaussian ADF, we first compute the (locally) exact posterior  $p_t(\mathbf{z}_t|\mathbf{y}_{1:t})$ , and then approximate it with  $q_t(\mathbf{z}_t|\mathbf{y}_{1:t})$  by projecting into  $\mathcal{Q}$ . Thus ADF approximates the conditional  $p_t(\mathbf{z}_t|\mathbf{y}_{1:t})$ , whereas GF approximates the joint  $q_{t|t-1}(\mathbf{z}_t, \tilde{\mathbf{y}}_t|\mathbf{y}_{1:t-1})$ , from which we derive  $p_t(\mathbf{z}_t|\mathbf{y}_{1:t})$  by conditioning.

ADF is more accurate than GF, since it directly approximates the posterior, but it is more computationally demanding, for reasons explained in [Wüt+16]. However, in [Kam+22] they propose an approximate form of expectation propagation (which is a generalization of ADF) in which the messages are computed using the same local joint Gaussian approximation as used in Gaussian filtering. See Figure 8.14 for a summary of how these different methods relate.

### 8.10.2 ADF for SLDS (Gaussian sum filter)

In this section, we apply ADF to inference in switching linear dynamical systems (SLDS, Section 29.9), which are a combination of HMM and LDS models. The resulting method is known as the **Gaussian sum filter** (see e.g., [Cro+11; Wil+17]).

A Gaussian sum filter approximates the belief state at each step by a mixture of  $K$  Gaussians. This can be implemented by running  $K$  Kalman filters in parallel. This is particularly well suited to switching SSMs. We now describe one version of this algorithm, known as the “second order generalized pseudo Bayes filter” (GPB2) [BSF88]. We assume that the prior belief state  $b_{t-1}$  is a mixture of  $K$  Gaussians, one per discrete state:

$$b_{t-1}^i \triangleq p(\mathbf{z}_{t-1}, m_{t-1} = i | \mathbf{y}_{1:t-1}) = \pi_{t-1|t-1}^i \mathcal{N}(\mathbf{z}_{t-1} | \boldsymbol{\mu}_{t-1|t-1}^i, \boldsymbol{\Sigma}_{t-1|t-1}^i) \quad (8.230)$$

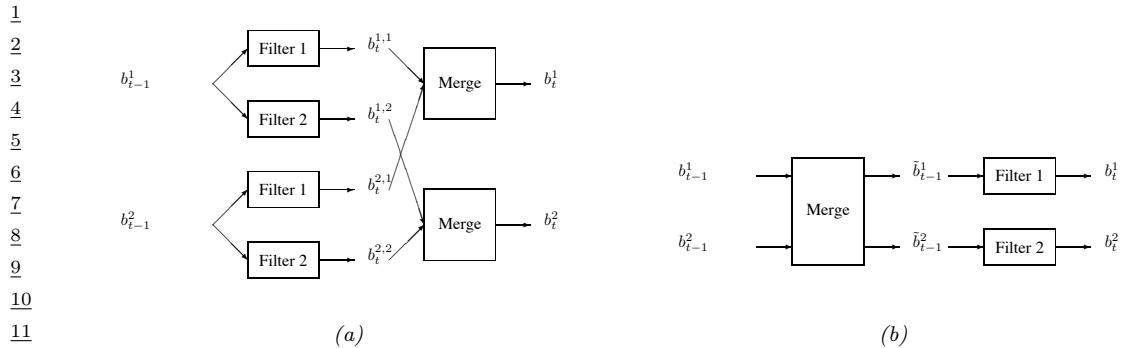


Figure 8.15: ADF for a switching linear dynamical system with 2 discrete states. (a) GPB2 method. (b) IMM method.

where  $i \in \{1, \dots, K\}$ . We then pass this through the  $K$  different linear models to get

$$b_t^{ij} \triangleq p(\mathbf{z}_t, m_{t-1} = i, m_t = j | \mathbf{y}_{1:t}) = \pi_{t|t}^{ij} \mathcal{N}(\mathbf{z}_t | \boldsymbol{\mu}_{t|t}^{ij}, \boldsymbol{\Sigma}_{t|t}^{ij}) \quad (8.231)$$

where  $\pi_{t|t}^{ij} = \pi_{t-1|t-1}^i A_{ij}$ , where  $A_{ij} = p(m_t = j | m_{t-1} = i)$ . Finally, for each value of  $j$ , we collapse the  $K$  Gaussian mixtures down to a single mixture to give

$$b_t^j \triangleq p(\mathbf{z}_t, m_t = j | \mathbf{y}_{1:t}) = \pi_{t|t}^j \mathcal{N}(\mathbf{z}_t | \boldsymbol{\mu}_{t|t}^j, \boldsymbol{\Sigma}_{t|t}^j) \quad (8.232)$$

See Figure 8.15a for a sketch.

The optimal way to approximate a mixture of Gaussians with a single Gaussian is given by  $q = \arg \min_q D_{\text{KL}}(q \| p)$ , where  $p(\mathbf{z}) = \sum_k \pi^k \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}^k, \boldsymbol{\Sigma}^k)$  and  $q(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$ . This can be solved by moment matching, that is,

$$\boldsymbol{\mu} = \mathbb{E}[\mathbf{z}] = \sum_k \pi^k \boldsymbol{\mu}^k \quad (8.233)$$

$$\boldsymbol{\Sigma} = \text{Cov}[\mathbf{z}] = \sum_k \pi^k \left( \boldsymbol{\Sigma}^k + (\boldsymbol{\mu}^k - \boldsymbol{\mu})(\boldsymbol{\mu}^k - \boldsymbol{\mu})^\top \right) \quad (8.234)$$

In the graphical model literature, this is called **weak marginalization** [Lau92], since it preserves the first two moments. Applying these equations to our model, we can go from  $b_t^{ij}$  to  $b_t^j$  as follows (where we drop the  $t$  subscript for brevity):

$$\pi^j = \sum_i \pi^{ij} \quad (8.235)$$

$$\pi^{j|i} = \frac{\pi^{ij}}{\sum_{j'} \pi^{ij'}} \quad (8.236)$$

$$\boldsymbol{\mu}^j = \sum_i \pi^{j|i} \boldsymbol{\mu}^{ij} \quad (8.237)$$

$$\boldsymbol{\Sigma}^j = \sum_i \pi^{j|i} (\boldsymbol{\Sigma}^{ij} + (\boldsymbol{\mu}^{ij} - \boldsymbol{\mu}^j)(\boldsymbol{\mu}^{ij} - \boldsymbol{\mu}^j)^\top) \quad (8.238)$$

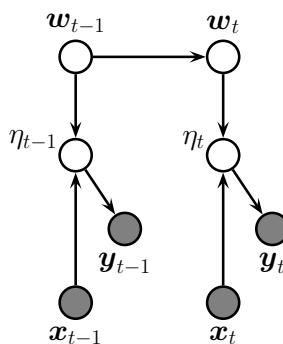


Figure 8.16: A dynamic logistic regression model.  $\mathbf{w}_t$  are the regression weights at time  $t$ , and  $\eta_t = \mathbf{w}_t^\top \mathbf{x}_t$ . Compare to Figure 29.24a.

This algorithm requires running  $K^2$  filters at each step. A cheaper alternative, known as **interactive multiple models** or **IMM** [BSF88], can be obtained by first collapsing the prior to a single Gaussian (by moment matching), and then updating it using  $K$  different Kalman filters, one per value of  $m_t$ . See Figure 8.15b for a sketch.

### 8.10.3 ADF for online logistic regression

In this section we discuss the application of ADF to online Bayesian parameter inference for a binary logistic regression model, based on [Zoe07]. The overall approach is similar to the online linear regression case (discussed in Section 29.7.2), but approximates the posterior after each update step, which is necessary since the likelihood is not conjugate to the prior.

We assume our model has the following form:

$$p(y_t | \mathbf{x}_t, \mathbf{w}_t) = \text{Ber}(y_t | \sigma(\mathbf{x}_t^\top \mathbf{w}_t)) \quad (8.239)$$

$$p(\mathbf{w}_t | \mathbf{w}_{t-1}) = \mathcal{N}(\mathbf{w}_t | \mathbf{w}_{t-1}, \mathbf{Q}) \quad (8.240)$$

where  $\mathbf{Q}$  is the covariance of the process noise, which allows the parameters to change slowly over time. We will assume  $\mathbf{Q} = \epsilon \mathbf{I}$ ; we can also set  $\epsilon = 0$ , as in the recursive least squares method (Section 29.7.2), if we believe the parameters will not change. See Figure 8.16 for an illustration of the model.

As our approximating family, we will use diagonal Gaussians, for computational efficiency. Thus the prior is the posterior from the previous timestep, and has the form

$$p(\mathbf{w}_{t-1} | \mathcal{D}_{1:t-1}) \approx p_{t-1}(\mathbf{w}_{t-1}) = \prod_j \mathcal{N}(w_{t-1}^j | \mu_{t-1|t-1}^j, \tau_{t-1|t-1}^j) \quad (8.241)$$

where  $\mu_{t-1|t-1}^j$  and  $\tau_{t-1|t-1}^j$  are the posterior mean and variance for parameter  $j$  given past data.

Now we discuss how to update this prior.

1 First we compute the one-step-ahead predictive density  $p_{t|t-1}(\mathbf{w}_t)$  using the standard linear-  
2 Gaussian update, i.e.,  $\boldsymbol{\mu}_{t|t-1} = \boldsymbol{\mu}_{t-1|t-1}$  and  $\boldsymbol{\tau}_{t|t-1} = \boldsymbol{\tau}_{t-1|t-1} + \mathbf{Q}$ , where we can set  $\mathbf{Q} = 0\mathbf{I}$  if there  
3 is no drift.  
4

5 Now we concentrate on the measurement update step. Define the scalar sum (corresponding to the  
6 logits, if we are using binary classification) as  $\eta_t = \mathbf{w}_t^\top \mathbf{x}_t$ . If  $p_{t|t-1}(\mathbf{w}_t) = \prod_j \mathcal{N}(w_t^j | \mu_{t|t-1}^j, \tau_{t|t-1}^j)$ ,  
7 then we can compute the 1d prior predictive distribution for  $\eta_t$  as follows:

$$\underline{8} \quad p(\eta_t | \mathcal{D}_{1:t-1}, \mathbf{x}_t) \approx p_{t|t-1}(\eta_t) = \mathcal{N}(\eta_t | m_{t|t-1}, v_{t|t-1}) \quad (8.242)$$

$$\underline{10} \quad m_{t|t-1} = \sum_j x_{t,j} \mu_{t|t-1}^j \quad (8.243)$$

$$\underline{12} \quad v_{t|t-1} = \sum_j x_{t,j}^2 \tau_{t|t-1}^j \quad (8.244)$$

15 The posterior for the 1d  $\eta_t$  is given by

$$\underline{16} \quad p(\eta_t | \mathcal{D}_{1:t}) \approx p_t(\eta_t) = \mathcal{N}(\eta_t | m_t, v_t) \quad (8.245)$$

$$\underline{18} \quad m_t = \int \eta_t \frac{1}{Z_t} p(y_t | \eta_t) p_{t|t-1}(\eta_t) d\eta_t \quad (8.246)$$

$$\underline{20} \quad v_t = \int \eta_t^2 \frac{1}{Z_t} p(y_t | \eta_t) p_{t|t-1}(\eta_t) d\eta_t - m_t^2 \quad (8.247)$$

$$\underline{22} \quad Z_t = \int p(y_t | \eta_t) p_{t|t-1}(\eta_t) d\eta_t \quad (8.248)$$

25 where  $p(y_t | \eta_t) = \text{Ber}(y_t | \eta_t)$ . These integrals are one dimensional, and so can be efficiently computed  
26 using Gaussian quadrature, as explained in [Zoe07; KB00].

27 Having inferred  $p_t(\eta_t)$ , we need to compute  $p_t(\mathbf{w} | \eta_t)$ . This can be done as follows. Define  $\delta_m$  as  
28 the change in the mean and  $\delta_v$  as the change in the variance:

$$\underline{29} \quad m_t = m_{t|t-1} + \delta_m, \quad v_t = v_{t|t-1} + \delta_v \quad (8.249)$$

31 Using the fact that  $p(\eta_t | \mathbf{w}) = \mathcal{N}(\eta_t | \mathbf{w}^\top \eta_t, 0)$  is a linear Gaussian system, with prior  $p(\mathbf{w}) =$   
32  $p(\mathbf{w} | \boldsymbol{\mu}_{t|t-1}, \boldsymbol{\tau}_{t|t-1})$  and “soft evidence”  $p(\eta_t) = \mathcal{N}(m_t, v_t)$ , we can derive the posterior for  $p(\mathbf{w} | \mathcal{D}_t)$  as  
33 follows:  
34

$$\underline{35} \quad p_t(w_t^i) = \mathcal{N}(w_t^i | \mu_{t|t}^i, \tau_{t|t}^i) \quad (8.250)$$

$$\underline{37} \quad \mu_{t|t}^i = \mu_{t|t-1}^i + a_i \delta_m \quad (8.251)$$

$$\underline{38} \quad \tau_{t|t}^i = \tau_{t|t-1}^i + a_i^2 \delta_v \quad (8.252)$$

$$\underline{40} \quad a_i \triangleq \frac{x_t^i \tau_{t|t-1}^i}{\sum_j (x_t^j)^2 + \tau_{t|t-1}^j} \quad (8.253)$$

43 Thus we see that the parameters which correspond to inputs  $i$  with larger magnitude (big  $|x_t^i|$ ) or  
44 larger uncertainty (big  $\tau_{t|t-1}^i$ ) get updated most, due to a large  $a_i$  factor, which makes intuitive sense.

45 As an example, we consider a 2d binary classification problem. We sequentially compute the  
46 posterior using the ADF, and compare to the offline estimate computed using a Laplace approximation.

47

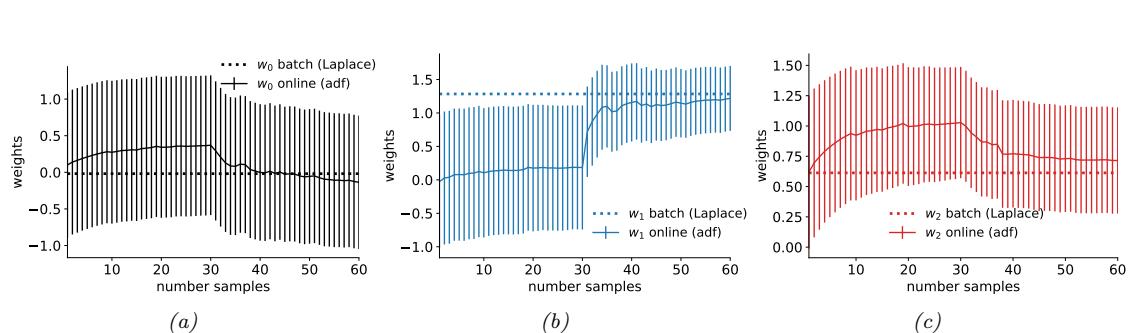


Figure 8.17: Bayesian inference applied to a 2d binary logistic regression problem,  $p(y = 1|\mathbf{x}) = \sigma(w_0 + w_1x_1 + w_2x_2)$ . We show the training data and the posterior predictive produced by different methods. (a) Offline MCMC approximation. (b) Offline Laplace approximation. (c) Online ADF approximation at the final step of inference. Generated by [adf\\_logistic\\_regression\\_demo.ipynb](#).

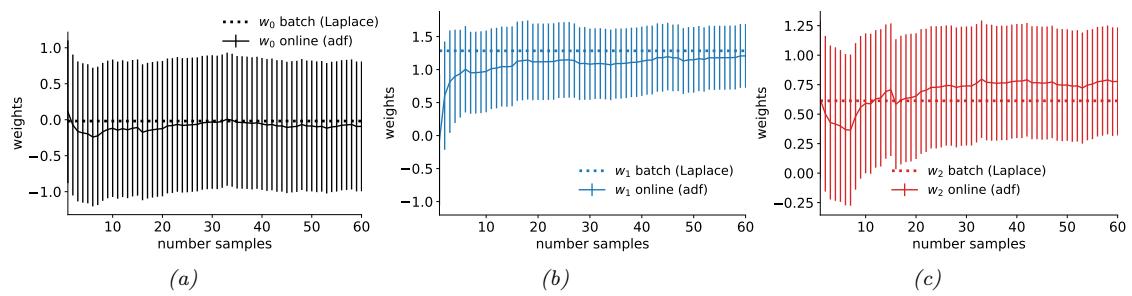


Figure 8.18: Same as Figure 8.17, except the order in which the data is visited is different. Generated by [adf\\_logistic\\_regression\\_demo.ipynb](#).

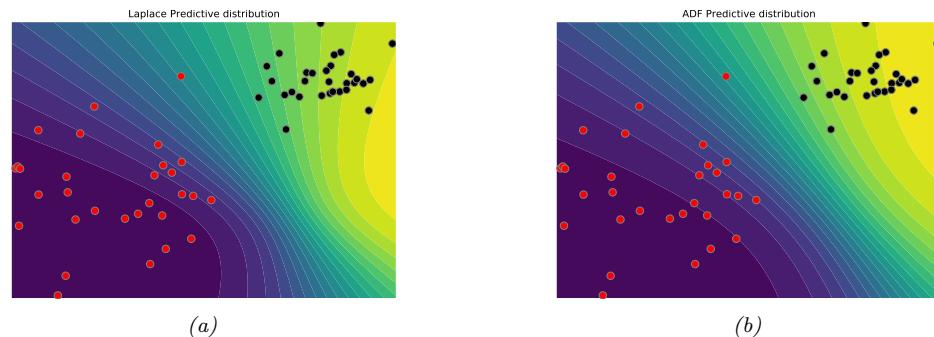


Figure 8.19: Predictive distribution for the binary logistic regression problem. (a) Result from Laplace approximation. (b) Result from ADF at the final step. Generated by [adf\\_logistic\\_regression\\_demo.ipynb](#).

1 In Figure 8.17 we plot the posterior marginals over the 3 parameters as a function of “time” (i.e., after  
2 conditioning on each training example one). We see that we converge to the offline MAP estimate. In  
3 Figure 8.18, we show the results of performing sequential Bayesian updating in a different ordering  
4 of the data. We still converge to approximate the same answer. In Figure 8.19, we see that the  
5 resulting posterior predictive distributions from the Laplace estimate and ADF estimate (at the end  
6 of training) are similar.  
7

8 Note that the whole algorithm only takes  $O(D)$  time and space per step, the same as SGD. However,  
9 unlike SGD, there are no step-size parameters, since the diagonal covariance implicitly specifies the  
10 size of the update for each dimension. Furthermore, we get a posterior approximation, not just a  
11 point estimate.

12 The overall approach is very similar to the generalized posterior linearization filter of Section 8.9.5.2,  
13 which uses quadrature (or the unscented transform) to compute a Gaussian approximation to the  
14 joint  $p(y_t, \mathbf{w}_t | \mathcal{D}_{1:t-1})$ , from which we can easily compute  $p(\mathbf{w}_t | \mathcal{D}_{1:t})$ . However, ADF approximates  
15 the posterior rather than the joint, as explained in Section 8.10.1.

16

#### 17 8.10.4 ADF for online DNNs

18  
19 In Section 17.5.3, we show how to use ADF to recursively approximate the posterior over the  
20 parameters of a deep neural network in an online fashion. This generalizes Section 8.10.3 to the case  
21 of nonlinear models.

22

### 23 8.11 Other inference methods for SSMs

24

25 There are a variety of other inference algorithms that can be applied to SSMs which we discuss  
26 elsewhere in this book. We give a very brief summary below, mostly focused on the case of offline  
27 inference (smoothing).

28

#### 29 8.11.1 Expectation propagation

30  
31 In Section 10.7 we discuss the expectation propagation (EP) algorithm, which can be viewed as an  
32 iterative version of ADF (Section 8.10). In particular, at each step we combine each exact local  
33 likelihood factor with approximate factors from both the past filtering distribution and the future  
34 smoothed posterior; these factors are combined to compute the locally exact posterior, which is then  
35 projected back to the tractable family (e.g., Gaussian), before moving to the next time step. This  
36 process can be iterated for increased accuracy. In many cases the local EP update is intractable,  
37 but we can make a local Gaussian approximation, similar to the one in general Gaussian filtering  
38 (Section 8.8), as explained in [Kam+22].  
39

#### 40 8.11.2 Variational inference

41

42 EP can be viewed as locally minimizing the inclusive KL,  $D_{\text{KL}}(p(\mathbf{z}_t | \mathbf{y}_{1:T}) \| q(\mathbf{z}_t | \mathbf{y}_{1:T}))$ , for each time  
43 step  $t$ . An alternative approach is to globally minimize the exclusive KL,  $D_{\text{KL}}(q(\mathbf{z}_{1:T} | \mathbf{y}_{1:T}) \| p(\mathbf{z}_{1:T} | \mathbf{y}_{1:T}))$ ;  
44 this is called variational inference, and is explained in Chapter 10. The difference between these two  
45 objectives is discussed in more detail in Section 5.1.3.3, but from a practical point of view, the main  
46 advantage of VI is that we can derive a tractable lower bound to the objective, and can then optimize  
47

it using stochastic optimization. This method is guaranteed to converge, unlike EP. For more details on VI applied to SSMs (both state estimation and parameter estimation), see e.g., [CWS21; Cou+20; Cou+21; BFY20; FLMM21; Cam+21].

### 8.11.3 MCMC

In Chapter 12 we discuss Markov chain Monte Carlo (MCMC) methods, which can be used to draw samples from intractable posteriors. In the case of SSMs, this includes both the distribution over states,  $p(\mathbf{z}_{1:T}|\mathbf{y}_{1:T})$ , and the distribution over parameters,  $p(\boldsymbol{\theta}|\mathbf{y}_{1:T})$ . In some cases, such as when using HMMs or linear-Gaussian SSMs, we can perform blocked Gibbs sampling, in which we use forwards filtering backwards sampling to sample an entire sequence from  $p(\mathbf{z}_{1:T}|\mathbf{y}_{1:T}, \boldsymbol{\theta})$ , followed by sampling the parameters,  $p(\boldsymbol{\theta}|\mathbf{z}_{1:T}, \mathbf{y}_{1:T})$  (see e.g., [CK96; Sco02; CMR05] for details.) Alternatively we can marginalize out the hidden states and just compute the parameter posterior  $p(\boldsymbol{\theta}|\mathbf{y}_{1:T})$ . When state inference is intractable, we can use gradient-based HMC methods (assuming the states are continuous), although this does not scale well to long sequences.

### 8.11.4 Particle filtering

In Section 13.2 we discuss particle filtering, which is a form of sequential Bayesian inference for SSMs which replaces the assumption that the posterior is (approximately) Gaussian with a more flexible representation, namely a set of weighted samples called “particles” (see e.g., [Aru+02; DJ11; NLS19]). Essentially the technique amounts to a form of importance sampling, combined with steps to prevent “particle impoverishment”, which refers to some samples receiving negligible weight because they are too improbable in the posterior (which grows with time). Particle filtering is widely used because it is very flexible, and has good theoretical properties. In practice it may require many samples to get a good approximation, but we can use heuristic methods, such as the extended or unscented Kalman filters, as proposal distributions, which can improve the efficiency significantly. In the offline setting, we can use particle smoothing (Section 13.5) or SMC (sequential Monte Carlo) samplers (Section 13.6).



# 9 Inference for graphical models

## 9.1 Introduction

In this chapter we consider posterior inference (i.e., computing marginals, modes, samples, etc) for probability distributions that can be represented by a PGM with some kind of sparse graph structure (i.e., it is not a fully connected graph). The algorithms we discuss will leverage the conditional independence properties encoded in the graph structure (discussed in Chapter 4) in order to perform efficient inference. In particular, we will use the principle of **dynamic programming** (DP), which finds an optimal solution by solving subproblems and then combining them.

DP can be implemented by computing local quantities for each node (or clique) in the graph, and then sending **messages** to neighboring nodes (or cliques) so that all nodes (cliques) can come to an overall consensus about the global solutions. Hence these are known as **message passing algorithms**. Each message can be interpreted as probability distribution about the value of a node given evidence from part of the graph. These distributions are often called **belief states**, so these algorithms are also called **belief propagation (BP)** algorithms.

The methods we discuss generalize the algorithms from Chapter 8, which only work with chain structured graphical models, to work with PGMs with arbitrary graph structure. This requires that we specify a **message passing schedule**. For chain structured models, a natural approach is to send messages forwards in time, and then backwards in time, as we discussed in Chapter 8. We can generalize this approach to work with trees, as we discuss in Section 9.2. For general graphs, there may be cycles or loops, as we discuss in Section 9.3. However, sending messages on loopy graphs may give incorrect answers. In such cases, we may wish to convert the graph to a tree, and then send messages on it, using the methods discussed in Section 9.4 and Section 9.5. We can also pose the inference problem as an optimization problem, as we discuss in Section 9.6.

## 9.2 Belief propagation on trees

The forwards-backwards algorithm for HMMs (Section 8.2.4) and the Kalman smoother algorithm for LDS (Section 8.3.3) can both be interpreted as message passing algorithms applied to chain structured graphical models. In this section, we generalize these algorithms to work with trees.

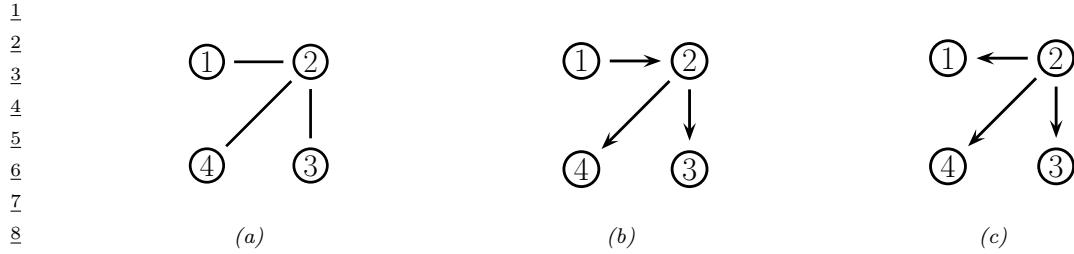


Figure 9.1: An undirected tree and two equivalent directed trees.

### 9.2.1 Directed vs undirected trees

Consider a pairwise *undirected* graphical model, which can be written as follows:

$$p^*(\mathbf{z}) \triangleq p(\mathbf{z}|\mathbf{y}) \propto \prod_{s \in \mathcal{V}} \psi_s(z_s|\mathbf{y}_s) \prod_{(s,t) \in \mathcal{E}} \psi_{s,t}(z_s, z_t) \quad (9.1)$$

where  $\psi_{s,t}(z_s, z_t)$  are the pairwise clique potential, one per edge,  $\psi_s(z_s|\mathbf{y}_s)$  are the local evidence potentials, one per node,  $\mathcal{V}$  is the set of nodes, and  $\mathcal{E}$  is the set of edges. (We will henceforth drop the conditioning on the observed values  $\mathbf{y}$  for brevity.)

Now suppose the corresponding graph structure is a tree, such as the one in Figure 9.1a. We can always convert this into a directed tree by picking an arbitrary node as the root, and then “picking the tree up by the root” and orienting all the edges away from the root. For example, if we pick node 1 as the root we get Figure 9.1b. This corresponds to the following directed graphical model:

$$p^*(\mathbf{z}) \propto p^*(z_1)p^*(z_2|z_1)p^*(z_3|z_2)p^*(z_4|z_2) \quad (9.2)$$

However, if we pick node 2 as the root, we get Figure 9.1c. This corresponds to the following directed graphical model:

$$p^*(\mathbf{z}) \propto p^*(z_2)p^*(z_1|z_2)p^*(z_3|z_2)p^*(z_4|z_2) \quad (9.3)$$

Since these graphs express the same conditional independence properties, they represent the same family of probability distributions, and hence we are free to use any of these parameterizations.

To make the model more symmetric, it is preferable to use an undirected tree. If we define the potentials as (possibly unnormalized) marginals (i.e.,  $\psi_s(z_s) \propto p^*(z_s)$  and  $\psi_{s,t}(z_s, z_t) = p^*(z_s, z_t)$ ), then we can write

$$p^*(\mathbf{z}) \propto \prod_{s \in \mathcal{V}} p^*(z_s) \prod_{(s,t) \in \mathcal{E}} \frac{p^*(z_s, z_t)}{p^*(z_s)p^*(z_t)} \quad (9.4)$$

For example, for Figure 9.1a we have

$$p^*(z_1, z_2, z_3, z_4) \propto p^*(z_1)p^*(z_2)p^*(z_3)p^*(z_4) \frac{p^*(z_1, z_2)p^*(z_2, z_3)p^*(z_2, z_4)}{p^*(z_1)p^*(z_2)p^*(z_2)p^*(z_3)p^*(z_2)p^*(z_4)} \quad (9.5)$$

```

1
2 // Collect to root
3 for each node  $s$  in post-order
4    $\text{bel}_s(z_s) \propto \psi_s(z_s) \prod_{t \in \text{ch}_s} m_{t \rightarrow s}(z_s)$ 
5    $t = \text{parent}(s)$ 
6    $m_{s \rightarrow t}(z_t) = \sum_{z_s} \psi_{st}(z_s, z_t) \text{bel}_s(z_s)$ 
7
8 // Distribute from root
9 for each node  $t$  in pre-order
10   $s = \text{parent}(t)$ 
11   $m_{s \rightarrow t}(z_t) = \frac{\text{bel}_s(z_s)}{m_{t \rightarrow s}(z_s)}$ 
12   $\text{bel}_t(z_t) \propto \text{bel}_t(z_t) m_{s \rightarrow t}(z_t)$ 
13
14
15
16
17
18 To see the equivalence with the directed representation, we can cancel terms to get
19
20  $p^*(z_1, z_2, z_3, z_4) \propto p^*(z_1, z_2) \frac{p^*(z_2, z_3)}{p^*(z_2)} \frac{p^*(z_2, z_4)}{p^*(z_2)}$  (9.6)
21
22  $= p^*(z_1) p^*(z_2|z_1) p^*(z_3|z_2) p^*(z_4|z_2)$  (9.7)
23  $= p^*(z_2) p^*(z_1|z_2) p^*(z_3|z_2) p^*(z_4|z_2)$  (9.8)
24
25 where  $p^*(z_t|z_s) = p^*(z_s, z_t)/p^*(z_s)$ .
26 Thus a tree can be represented as either an undirected or directed graph. Both representations
27 can be useful, as we will see.
28
29
30 

### 9.2.2 Sum-product algorithm


31 In this section, we assume that our model is an undirected tree, as in Equation (9.1). However, we
32 will pick an arbitrary node as a root, and orient all the edges downwards away from this root, so that
33 each node has a unique parent. For a directed, rooted tree, we can compute various node orderings.
34 In particular, in a pre-order, we traverse from the root to the left subtree and then to right subtree,
35 top to bottom. In a post-order, we traverse from the left subtree to the right subtree and then to
36 the root, bottom to top. We will use both of these below.
37 We now present the sum-product algorithm for trees. We first send messages from the leaves
38 to the root. This is the generalization of the forwards pass from Section 8.2.2. Let  $m_{s \rightarrow t}(z_t)$  denote
39 the message from node  $s$  to node  $t$ . This summarizes the belief state about  $z_t$  given all the evidence
40 in the tree below the  $s - t$  edge. Consider a node  $s$  in the ordering. We update its belief state by
41 combining the incoming messages from all its children with its own local evidence:
42
43  $\text{bel}_s(z_s) \propto \psi_s(z_s) \prod_{t \in \text{ch}_s} m_{t \rightarrow s}(z_s)$  (9.9)
44
45
46 To compute the outgoing message that  $s$  should send to its parent  $t$ , we pass the local belief through
47

```

Figure 9.2: Belief propagation on a pairwise, rooted tree.

To see the equivalence with the directed representation, we can cancel terms to get

$$p^*(z_1, z_2, z_3, z_4) \propto p^*(z_1, z_2) \frac{p^*(z_2, z_3)}{p^*(z_2)} \frac{p^*(z_2, z_4)}{p^*(z_2)} \quad (9.6)$$

$$= p^*(z_1) p^*(z_2|z_1) p^*(z_3|z_2) p^*(z_4|z_2) \quad (9.7)$$

$$= p^*(z_2) p^*(z_1|z_2) p^*(z_3|z_2) p^*(z_4|z_2) \quad (9.8)$$

where  $p^*(z_t|z_s) = p^*(z_s, z_t)/p^*(z_s)$ .

Thus a tree can be represented as either an undirected or directed graph. Both representations can be useful, as we will see.

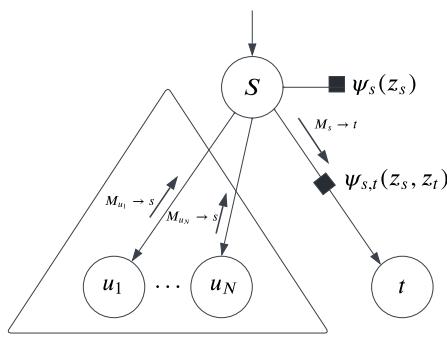
### 9.2.2 Sum-product algorithm

In this section, we assume that our model is an undirected tree, as in Equation (9.1). However, we will pick an arbitrary node as a root, and orient all the edges downwards away from this root, so that each node has a unique parent. For a directed, rooted tree, we can compute various node orderings. In particular, in a **pre-order**, we traverse from the root to the left subtree and then to right subtree, top to bottom. In a **post-order**, we traverse from the left subtree to the right subtree and then to the root, bottom to top. We will use both of these below.

We now present the **sum-product algorithm** for trees. We first send messages from the leaves to the root. This is the generalization of the forwards pass from Section 8.2.2. Let  $m_{s \rightarrow t}(z_t)$  denote the message from node  $s$  to node  $t$ . This summarizes the belief state about  $z_t$  given all the evidence in the tree below the  $s - t$  edge. Consider a node  $s$  in the ordering. We update its belief state by combining the incoming messages from all its children with its own local evidence:

$$\text{bel}_s(z_s) \propto \psi_s(z_s) \prod_{t \in \text{ch}_s} m_{t \rightarrow s}(z_s) \quad (9.9)$$

To compute the outgoing message that  $s$  should send to its parent  $t$ , we pass the local belief through



*Figure 9.3: Illustration of how the top-down message from  $s$  to  $t$  is computed during BP on a tree. The  $u_i$  nodes are the other children of  $s$ , besides  $t$ . Square nodes represent clique potentials.*

the pairwise potential linking  $s$  and  $t$ , and then marginalize out  $s$  to get

$$m_{s \rightarrow t}(z_t) = \sum_{z_s} \psi_{st}(z_s, z_t) \text{bel}_s(z_s) \quad (9.10)$$

At the root of the tree,  $\text{bel}_t(z_t) = p(z_t | \mathbf{y})$  will have seen all the evidence. It can then send messages back down to the leaves. The message that  $s$  sends to its child  $t$  should be the product of all the messages that  $s$  received from all its *other* children  $u$ , passed through the pairwise potential, and then marginalized:

$$m_{s \rightarrow t}(z_t) = \sum_{z_t} \left( \psi_s(z_s) \psi_{st}(z_s, z_t) \prod_{u \in \text{ch}_s \setminus t} m_{u \rightarrow s}(z_s) \right) \quad (9.11)$$

See Figure 9.3. Instead of multiplying all-but-one of the messages that  $s$  has received, we can multiply all of them and then divide out by the  $t \rightarrow s$  message from child  $t$ . The advantage of this is that the product of all the messages has already been computed in Equation (9.9), so we don't need to recompute that term. Thus we get

$$m_{s \rightarrow t}(z_t) = \sum_{z_s} \psi_{st}(z_s, z_t) \frac{\text{bel}_s(z_s)}{m_{t \rightarrow s}(z_s)} \quad (9.12)$$

We can think of  $\text{bel}_s(z_s)$  as the new updated posterior  $p(z_s | \mathbf{y})$  given all the evidence, and  $m_{t \rightarrow s}(z_s)$  as the prior predictive  $p(z_s | \mathbf{y}_t^-)$ , where  $\mathbf{y}_t^-$  is all the evidence in the subtree rooted at  $t$ . Thus the ratio contains the new evidence that  $t$  did not already know about from its own subtree. We use this to update the belief state at node  $t$  to get:

$$\text{bel}_t(z_t) \propto \text{bel}_t(z_t) m_{s \rightarrow t}(z_t) \quad (9.13)$$

(Note that Equation (9.9) is a special case of this where we don't divide out by  $m_{s \rightarrow t}$ , since in the upwards pass, there is no incoming message from the parent.) This is analogous to the backwards

1 smoothing equation in Equation (8.24), with  $\alpha_t(i)$  replaced by  $\text{bel}_t(z_t = i)$ ,  $A(i, j)$  replaced by  
2  $\psi_{st}(z_s = i, z_t = j)$ ,  $\gamma_{t+1}(j)$  replaced by  $\text{bel}_s(z_s = j)$ , and  $\alpha_{t+1|t}(j)$  replaced by  $m_{t \rightarrow s}(z_s = j)$ .

3 See Figure 9.2 for the overall pseudocode. This can be generalized to directed trees with multiple  
4 root nodes (known as **polytrees**) as described in [Supplementary](#) Section 9.1.1.

### 7 9.2.3 Max-product algorithm

8 In Section 9.2.2 we described the sum-product algorithm, that computes the posterior marginals:

$$\text{bel}_i(k) = \gamma_i(k) = p(z_i = k | \mathbf{y}) = \sum_{\mathbf{z}_{-i}} p(z_i = k, \mathbf{z}_{-i} | \mathbf{y}) \quad (9.14)$$

13 We can replace the sum operation with the max operation to get **max-product belief propagation**.  
14 The result of this computation are a set of **max marginals** for each node:

$$\zeta_i(k) = \max_{\mathbf{z}_{-i}} p(z_i = k, \mathbf{z}_{-i} | \mathbf{y}) \quad (9.15)$$

18 We can derive two different kinds of “MAP” estimates from these local quantities. The first is  
19  $\hat{\mathbf{z}}_i = \text{argmax}_k \gamma_i(k)$ ; this is known as the **maximizer of the posterior marginal** or **MPM** estimate  
20 (see e.g., [MMP87; SM12]); let  $\hat{\mathbf{z}} = [\hat{z}_1, \dots, \hat{z}_{N_z}]$  be the sequence of such estimates. The second is  
21  $\tilde{\mathbf{z}}_i = \text{argmax}_k \zeta_i(k)$ ; we call this the **maximizer of the max marginal** or **MMM** estimate; let  
22  $\tilde{\mathbf{z}} = [\tilde{z}_1, \dots, \tilde{z}_{N_z}]$ .

23 An interesting question is: what, if anything, do these estimates have to do with the “true” MAP  
24 estimate,  $\mathbf{z}^* = \text{argmax}_{\mathbf{z}} p(\mathbf{z} | \mathbf{y})$ ? We discuss this below.

#### 26 9.2.3.1 Connection between MMM and MAP

28 In [YW04], they showed that, if the max marginals are unique and computed exactly (e.g., if the  
29 graph is a tree), then  $\tilde{\mathbf{z}} = \mathbf{z}^*$ . This means we can recover the global MAP estimate by running max  
30 product BP and then setting each node to its local max (i.e., using the MMM estimate).

31 However, if there are ties in the max marginals (corresponding to the case where there is more  
32 than one globally optimal solution), this “local stitching” process may result in global inconsistencies.

33 If we have a tree-structured model, we can use a **traceback** procedure, analogous to the Viterbi  
34 algorithm (Section 8.2.7), in which we clamp nodes to their optimal values while working backwards  
35 from the root. For details, see e.g., [KF09a, p569].

36 Unfortunately, traceback does not work on general graphs. An alternative, iterative approach,  
37 proposed in [YW04], is follows. First we run max product BP, and clamp all nodes which have unique  
38 max marginals to their optimal values; we then clamp a single ambiguous node to an optimal value,  
39 and condition on all these clamped values as extra evidence, and perform more rounds of message  
40 passing, until all ties are broken. This may require many rounds of inference, although the number  
41 of non-clamped (hidden) variables get reduced at each round.

#### 43 9.2.3.2 Connection between MPM and MAP

45 In this section, we discuss the MPM estimate,  $\hat{\mathbf{z}}$ , which computes the maximum of the posterior  
46 marginals. In general, this does not correspond to the MAP estimate, even if the posterior marginals

1 are exact. To see why, note that MPM just looks at the belief state for each node given all the visible  
 2 evidence, but ignores any dependencies or constraints that might exist in the prior.  
 3

4 To illustrate why this could be a problem, consider the error correcting code example from  
 5 Section 5.5, where we defined  $p(\mathbf{z}, \mathbf{y}) = p(z_1)p(z_2)p(z_3|z_1, z_2)\prod_{i=1}^3 p(y_i|z_i)$ , where all variables are  
 6 binary. The priors  $p(z_1)$  and  $p(z_2)$  are uniform. The conditional term  $p(z_3|z_1, z_2)$  is deterministic,  
 7 and computes the parity of  $(z_1, z_2)$ . In particular, we have  $p(z_3 = 1|z_1, z_2) = \mathbb{I}(\text{odd}(z_1, z_2))$ , so that  
 8 the total number of 1s in the block  $\mathbf{z}_{1:3}$  is even. The likelihood terms  $p(y_i|z_i)$  represent a bit flipping  
 9 noisy channel model, with noise level  $\alpha = 0.2$ .

10 Suppose we observe  $\mathbf{y} = (1, 0, 0)$ . In this case, the exact posterior marginals are as follows:<sup>1</sup>  $\gamma_1 =$   
 11  $[0.3469, 0.6531]$ ,  $\gamma_2 = [0.6531, 0.3469]$ ,  $\gamma_3 = [0.6531, 0.3469]$ . The exact max marginals are all the same,  
 12 namely  $\zeta_i = [0.3265, 0.3265]$ . Finally, the 3 global MAP estimates are  $\mathbf{z}^* \in \{[0, 0, 0], [1, 1, 0], [1, 0, 1]\}$ ,  
 13 each of which corresponds to a single bit flip from the observed vector. The MAP estimates are all  
 14 valid code words (they have an even number of 1s), and hence are sensible hypotheses about the  
 15 value of  $\mathbf{z}$ . By contrast, the MPM estimate is  $\hat{\mathbf{z}} = [1, 0, 0]$ , which is not a legal codeword. (And in  
 16 this example, the MMM estimate is not well defined, since the max marginals are not unique.)

17 So, which method is better? This depends on our loss function, as we discuss in Section 34.1. If  
 18 we want to minimize the prediction error of each  $z_i$ , also called **bit error**, we should compute the  
 19 MPM. If we want to minimize the prediction error for the entire sequence  $\mathbf{z}$ , also called **word error**,  
 20 we should use MAP, since this can take global constraints into account.

21 For example, suppose we are performing speech recognition and someone says “recognize speech”.  
 22 MPM decoding may return “wreck a nice beach”, since locally it may be that “beach” is the most  
 23 probable interpretation of “speech” when viewed in isolation (see Figure 34.1). However, MAP  
 24 decoding would infer that “recognize speech” is the more likely overall interpretation, by taking into  
 25 account the language model prior,  $p(\mathbf{z})$ .

26 On the other hand, if we don’t have strong constraints, the MPM estimate can be more robust  
 27 [MMP87; SM12], since it marginalizes out the other nodes, rather than maxing them out. For  
 28 example, in the casino HMM example in Figure 8.4, we see that the MPM method makes 49 bit  
 29 errors (out of a total possible of  $T = 300$ ), and the MAP path makes 60 errors.

30

### 31 9.2.3.3 Connection between MPE and MAP

32 In the graphical models literature, computing the jointly most likely setting of all the latent variables,  
 33  $\mathbf{z}^* = \text{argmax}_{\mathbf{z}} p(\mathbf{z}|\mathbf{y})$ , is known as the **most probable explanation** or **MPE** [Pea88]. In that  
 34 literature, the term “MAP” is used to refer to the case where we maximize some of the hidden  
 35 variables, and marginalize (sum out) the rest. For example, if we maximize a single node,  $z_i$ , but  
 36 sum out all the others,  $\mathbf{z}_{-i}$ , we get the MPM  $\hat{z}_i = \text{argmax}_{z_i} \sum_{\mathbf{z}_{-i}} p(\mathbf{z}|\mathbf{y})$ .

37 We can generalize the MPM estimate to compute the best guess for a set of query variables  $Q$ ,  
 38 given evidence on a set of visible variables  $V$ , marginalizing out the remaining variables  $R$ , to get  
 39

$$40 \quad \mathbf{z}_Q^* = \arg \max_{\mathbf{z}_Q} \sum_{\mathbf{z}_R} p(\mathbf{z}_Q, \mathbf{z}_R | \mathbf{z}_V) \quad (9.16)$$

42 (Here  $\mathbf{z}_R$  are called **nuisance variables**, since they are not of interest, and are not observed.) In  
 43 [Pea88], this is called a MAP estimate, but we will call it an MPM estimate, to avoid confusion with  
 44 the ML usage of the term “MAP” (where we maximize everything jointly).

45 \_\_\_\_\_  
 46 1. See `error_correcting_code_demo.ipynb` for the code.

47

## 9.3 Loopy belief propagation

In this section, we extend belief propagation to work on graphs with cycles or loops; this is called **loopy belief propagation** or **LBP**. Unfortunately, this method may not converge, and even if it does, it is not clear if the resulting estimates are valid. Indeed, Judea Pearl, who invented belief propagation for trees, wrote the following about loopy BP in 1988:

When loops are present, the network is no longer singly connected and local propagation schemes will invariably run into trouble ... If we ignore the existence of loops and permit the nodes to continue communicating with each other as if the network were singly connected, messages may circulate indefinitely around the loops and the process may not converge to a stable equilibrium ... Such oscillations do not normally occur in probabilistic networks ... which tend to bring all messages to some stable equilibrium as time goes on. However, this asymptotic equilibrium is not coherent, in the sense that it does not represent the posterior probabilities of all nodes of the network — [Pea88, p.195]

Despite these reservations, Pearl advocated the use of belief propagation in loopy networks as an approximation scheme (J. Pearl, personal communication). [MWJ99] found empirically that it works on various graphical models, and it is now used in many real world applications, some of which we discuss below. In addition, there is now some theory justifying its use in certain cases, as we discuss below. (For more details, see e.g., [Yed11].)

### 9.3.1 Loopy BP for pairwise undirected graphs

In this section, we assume (for notational simplicity) that our model is an undirected pairwise PGM, as in Equation (9.1). However, unlike Section 9.2.2, we do not assume the graph is a tree. We can apply the same message passing equations as before. However, since there is no natural node ordering, we will do this in a parallel, asynchronous way. The basic idea is that all nodes receive messages from their neighbors in parallel, they then update their belief states, and finally they send new messages back out to their neighbors. This message passing process repeats until convergence. This kind of computing architecture is called a **systolic array**, due to its resemblance to a beating heart.

More precisely, we initialize all messages to the all 1's vector. Then, in parallel, each node absorbs messages from all its neighbors using

$$\text{bel}_s(z_s) \propto \psi_s(z_s) \prod_{t \in \text{nbr}_s} m_{t \rightarrow s}(z_s) \quad (9.17)$$

Then, in parallel, each node sends messages to each of its neighbors:

$$m_{s \rightarrow t}(z_t) = \sum_{z_s} \left( \psi_s(z_s) \psi_{st}(z_s, z_t) \prod_{u \in \text{nbr}_s \setminus t} m_{u \rightarrow s}(z_s) \right) \quad (9.18)$$

The  $m_{s \rightarrow t}$  message is computed by multiplying together all incoming messages, except the one sent by the recipient, and then passing through the  $\psi_{st}$  potential. We continue this process until convergence. If the graph is a tree, the method is guaranteed to converge after  $D(G)$  iterations, where  $D(G)$  is the **diameter** of the graph, that is, the largest distance between any two nodes.

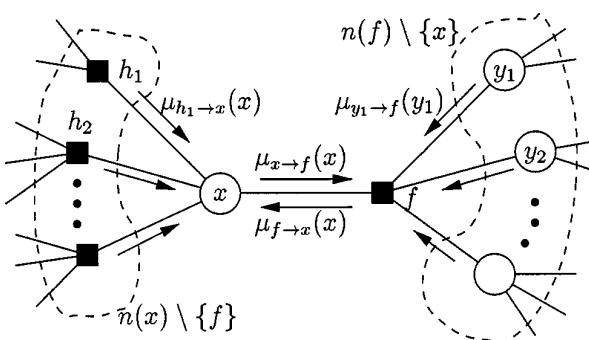


Figure 9.4: Message passing on a bipartite factor graph. Square nodes represent factors, and circles represent variables. The  $y_i$  nodes correspond to the neighbors  $x'_i$  of  $f$  other than  $x$ . From Figure 6 of [KFL01]. Used with kind permission of Brendan Frey.

### 9.3.2 Loopy BP for factor graphs

To implement loopy BP for general graphs, including those with higher-order clique potentials (beyond pairwise), it is useful to use a factor graph representation described in Section 4.6.1. In this section, we summarize the BP equations for the bipartite version of factor graphs, as derived in [KFL01].<sup>2</sup> For a version that works for Forney factor graphs, see [Loe+07].

In the case of bipartite factor graphs, we have two kinds of messages: variables to factors

$$m_{x \rightarrow f}(x) = \prod_{h \in \text{nbr}(x) \setminus \{f\}} m_{h \rightarrow x}(x) \quad (9.19)$$

and factors to variables

$$m_{f \rightarrow x}(x) = \sum_{x'} f(x, x') \prod_{x' \in \text{nbr}(f) \setminus \{x\}} m_{x' \rightarrow f}(x') \quad (9.20)$$

Here  $\text{nbr}(x)$  are all the factors that are connected to variable  $x$ , and  $\text{nbr}(f)$  are all the variables that are connected to factor  $f$ . These messages are illustrated in Figure 9.4. At convergence, we can compute the final beliefs as a product of incoming messages:

$$\text{bel}(x) \propto \prod_{f \in \text{nbr}(x)} m_{f \rightarrow x}(x) \quad (9.21)$$

The order in which the messages are sent can be determined using various heuristics, such as computing a spanning tree, and picking an arbitrary node as root. Alternatively, the update ordering can be chosen adaptively using **residual belief propagation** [EMK06]. Or fully parallel, asynchronous implementations can be used.

<sup>2</sup> For an efficient JAX implementation of these equations for discrete factor graphs, see <https://github.com/vicariousinc/PGMax>. For the Gaussian case, see <https://github.com/probml/pgm-jax>.

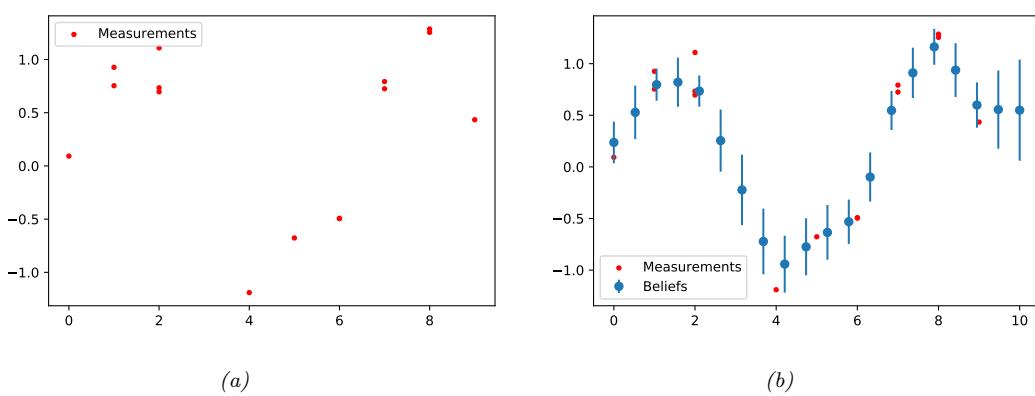


Figure 9.5: Interpolating noisy data using Gaussian belief propagation applied to a 1d MRF. Generated by [gauss-bp-1d-line.ipynb](#).

### 9.3.3 Gaussian belief propagation

It is possible to generalize (loopy) belief propagation to the Gaussian case, by using the ‘‘calculus for linear Gaussian models’’ in Section 2.2.7 to compute the messages and beliefs. Note that computing the posterior mean in a linear-Gaussian system is equivalent to solving a linear system, so these methods are also useful for linear algebra. See e.g., [PL03; Bic09; Du+18] for details.

As an example of Gaussian BP, consider the problem of interpolating noisy data in 1d, as discussed in [OED21]. In particular, let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be an unknown function for which we get  $N$  noisy measurements  $y_i$  at locations  $x_i$ . We want to estimate  $z_i = f(g_i)$  at  $G$  grid locations  $g_i$ . Let  $x_i$  be the closest location to  $g_i$ . Then we assume the measurement factor is as follows:

$$\psi_i(z_{i-1}, z_i) = \frac{1}{\sigma^2} (\hat{y}_i - y_i)^2 \quad (9.22)$$

$$\hat{y}_i = (1 - \gamma_i)z_{i-1} + \gamma_i z_i \quad (9.23)$$

$$\gamma_i = \frac{x_i - g_i}{g_i - g_{i-1}} \quad (9.24)$$

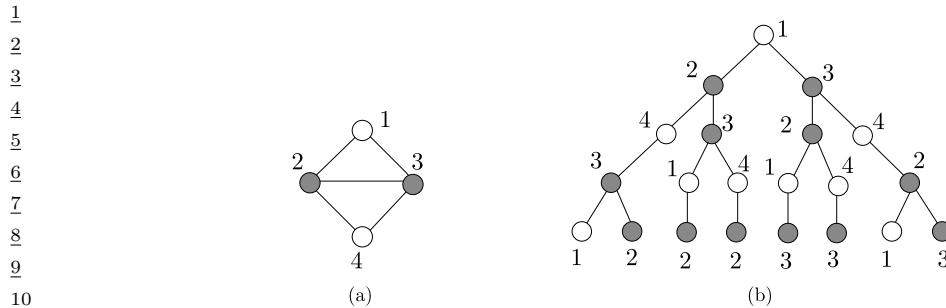
Here  $\hat{y}_i$  is the predicted measurement. The potential function makes the unknown function values  $z_{i-1}$  and  $z_i$  move closer to the observation, based on how far these grid points are from where the measurement was taken. In addition, we add a pairwise smoothness potential, that encodes the prior that  $z_i$  should be close to  $z_{i-1}$  and  $z_{i+1}$ :

$$\phi_i(z_{i-1}, z_i) = \frac{1}{\tau^2} \delta_i^2 \quad (9.25)$$

$$\delta_i = z_i - z_{i-1} \quad (9.26)$$

The overall model is

$$p(\mathbf{z}|\mathbf{x}, \mathbf{y}, \mathbf{g}, \sigma^2, \tau^2) \propto \prod_{i=1}^G \psi_i(z_{i-1}, z_i) \phi_i(z_{i-1}, z_i) \quad (9.27)$$



11 *Figure 9.6: (a) A simple loopy graph. (b) The computation tree, rooted at node 1, after 4 rounds of message*  
 12 *passing. Nodes 2 and 3 occur more often in the tree because they have higher degree than nodes 1 and 2.*  
 13 *From Figure 8.2 of [WJ08]. Used with kind permission of Martin Wainwright.*

14  
 15  
 16  
 17 Suppose the true underlying function is a sine wave. We show some sample data in Figure 9.5(a).  
 18 We then apply Gaussian BP. Since this model is a chain, and the model is linear-Gaussian, the  
 19 resulting posterior marginals, shown in Figure 9.5(b), are exact. We see that the method has inferred  
 20 the underlying sine shape just based on a smoothness prior.  
 21

22 To perform message passing in models with non-Gaussian potentials, we can extend the techniques  
 23 from Section 8.5.2 from chains to general graphs. For example, we can use local linearization, similar  
 24 to extended Kalman filter (see [PHR18]); or we can use **sigma point BP** [MHH14], similar to  
 25 unscented Kalman filter.

26 To extend to the non-Gaussian case, we can use **non-parametric BP** or **particle BP** (see e.g.,  
 27 [Sud+03; Isa03; Sud+10; Pac+14]), which uses ideas from particle filtering (Section 13.2).

28  
 29 **9.3.4 Convergence**

30 Loopy BP may not converge, or may only converge slowly. In this section, we discuss some techniques  
 31 that increase the chances of convergence, and the speed of convergence.  
 32

33  
 34 **9.3.4.1 When will LBP converge?**

35 The details of the analysis of when LBP will converge are beyond the scope of this chapter, but  
 36 we briefly sketch the basic idea. The key analysis tool is the **computation tree**, which visualizes  
 37 the messages that are passed as the algorithm proceeds. Figure 9.6 gives a simple example. In the  
 38 first iteration, node 1 receives messages from nodes 2 and 3. In the second iteration, it receives one  
 39 message from node 3 (via node 2), one from node 2 (via node 3), and two messages from node 4 (via  
 40 nodes 2 and 3). And so on.

41 The key insight is that  $T$  iterations of LBP is equivalent to exact computation in a computation  
 42 tree of height  $T + 1$ . If the strengths of the connections on the edges is sufficiently weak, then the  
 43 influence of the leaves on the root will diminish over time, and convergence will occur. See [MK05;  
 44 WJ08] and references therein for more information.

45  
 46  
 47

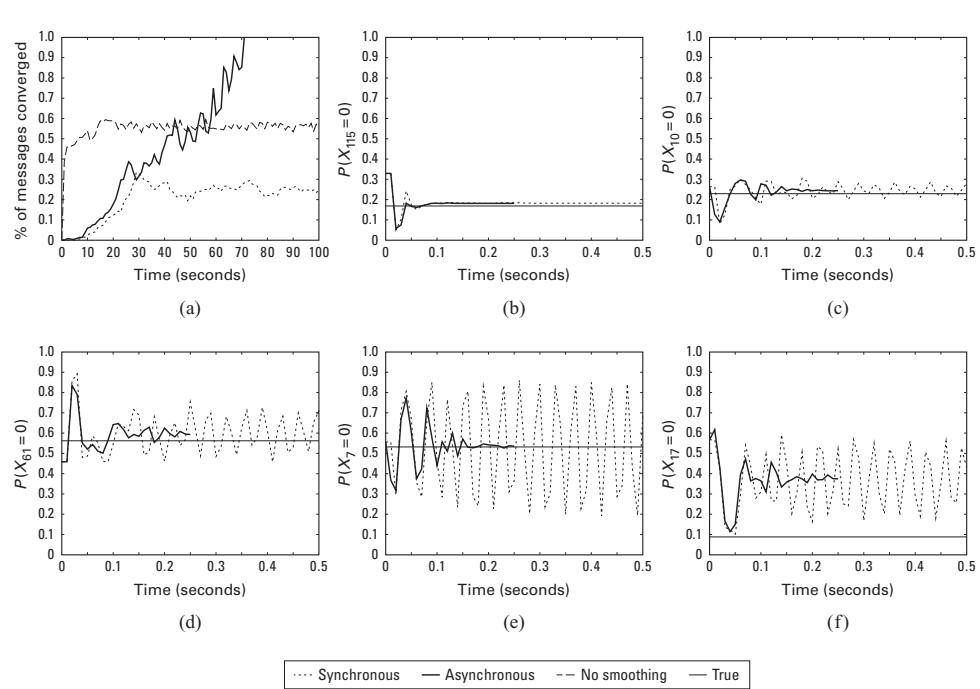


Figure 9.7: Illustration of the behavior of loopy belief propagation on an  $11 \times 11$  Ising grid with random potentials,  $w_{ij} \sim \text{Unif}(-C, C)$ , where  $C = 11$ . For larger  $C$ , inference becomes harder. (a) Percentage of messages that have converged vs time for 3 different update schedules: Dotted = damped synchronous (few nodes converge), dashed = undamped asynchronous (half the nodes converge), solid = damped asynchronous (all nodes converge). (b-f) Marginal beliefs of certain nodes vs time. Solid straight line = truth, dashed = synchronous, solid = damped asynchronous. From Figure 11.C.1 of [KF09a]. Used with kind permission of Daphne Koller.

### 9.3.4.2 Making LBP converge

Although the theoretical convergence analysis is very interesting, in practice, when faced with a model where LBP is not converging, what should we do?

One simple way to increase the chance of convergence is to use **damping**. That is, at iteration  $k$ , we use an update of the form

$$m_{t \rightarrow s}^k(x_s) = \lambda m_{t \rightarrow s}(x_s) + (1 - \lambda) m_{t \rightarrow s}^{k-1}(x_s) \quad (9.28)$$

where  $m_{t \rightarrow s}(x_s)$  is the standard undamped message, where  $0 \leq \lambda \leq 1$  is the damping factor. Clearly if  $\lambda = 1$  this reduces to the standard scheme, but for  $\lambda < 1$ , this partial updating scheme can help improve convergence. Using a value such as  $\lambda \sim 0.5$  is standard practice. The benefits of this approach are shown in Figure 9.7, where we see that damped updating results in convergence much more often than undamped updating (see [ZLG20] for some analysis of the benefits of damping).

It is possible to devise methods, known as **double loop algorithms**, which are guaranteed to converge to a local minimum of the same objective that LBP is minimizing [Yui01; WT01].

1 Unfortunately, these methods are rather slow and complicated, and the accuracy of the resulting  
 2 marginals is usually not much greater than with standard LBP. (Indeed, oscillating marginals is  
 3 sometimes a sign that the LBP approximation itself is a poor one.) Consequently, these techniques  
 4 are not very widely used.  
 5

6

### 7 9.3.4.3 Increasing the convergence rate with adaptive scheduling

8 The standard approach when implementing LBP is to perform **synchronous updates**, where all  
 9 nodes absorb messages in parallel, and then send out messages in parallel. That is, the new messages  
 10 at iteration  $k + 1$  are computed in parallel using  
 11

$$12 \quad \mathbf{m}_{1:E}^{k+1} = (f_1(\mathbf{m}^k), \dots, f_E(\mathbf{m}^k)) \quad (9.29)$$

13 where  $E$  is the number of edges, and  $f_i(\mathbf{m})$  is the function that computes the message for edge  $i$   
 14 given all the old messages. This is analogous to the Jacobi method for solving linear systems of  
 15 equations.  
 16

17 It is well known [Ber97] that the Gauss-Seidel method, which performs **asynchronous updates**  
 18 in a fixed round-robin fashion, converges faster when solving linear systems of equations. We can  
 19 apply the same idea to LBP, using updates of the form  
 20

$$21 \quad \mathbf{m}_i^{k+1} = f_i(\{\mathbf{m}_j^{k+1} : j < i\}, \{\mathbf{m}_j^k : j > i\}) \quad (9.30)$$

22 where the message for edge  $i$  is computed using new messages (iteration  $k + 1$ ) from edges earlier in  
 23 the ordering, and using old messages (iteration  $k$ ) from edges later in the ordering.  
 24

25 This raises the question of what order to update the messages in. One simple idea is to use a fixed  
 26 or random order. The benefits of this approach are shown in Figure 9.7, where we see that (damped)  
 27 asynchronous updating results in convergence much more often than synchronous updating.  
 28

29 However, we can do even better by using an adaptive ordering. The intuition is that we should  
 30 focus our computational efforts on those variables that are most uncertain. [EMK06] proposed  
 31 a technique known as **residual belief propagation**, in which messages are scheduled to be sent  
 32 according to the norm of the difference from their previous value. That is, we define the residual of  
 33 new message  $m_{s \rightarrow t}$  at iteration  $k$  to be  
 34

$$35 \quad r(s, t, k) = \|\log m_{s \rightarrow t} - \log m_{s \rightarrow t}^k\|_\infty = \max_j |\log \frac{m_{s \rightarrow t}(j)}{m_{s \rightarrow t}^k(j)}| \quad (9.31)$$

36 We can store messages in a priority queue, and always send the one with highest residual. When a  
 37 message is sent from  $s$  to  $t$ , all of the other messages that depend on  $m_{s \rightarrow t}$  (i.e., messages of the form  
 38  $m_{t \rightarrow u}$  where  $u \in \text{nbr}(t) \setminus s$ ) need to be recomputed; their residual is recomputed, and they are added  
 39 back to the queue. In [EMK06], they showed (experimentally) that this method converges more  
 40 often, and much faster, than using synchronous updating, asynchronous updating with a fixed order.  
 41

42 A refinement of residual BP was presented in [SM07]. In this paper, they use an upper bound on  
 43 the residual of a message instead of the actual residual. This means that messages are only computed  
 44 if they are going to be sent; they are not just computed for the purposes of evaluating the residual.  
 45 This was observed to be about five times faster than residual BP, although the quality of the final  
 46 results is similar.  
 47

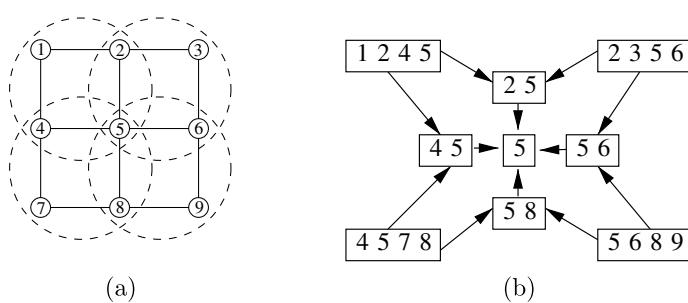


Figure 9.8: (a) Clusters superimposed on a  $3 \times 3$  lattice graph. (b) Corresponding hyper-graph. Nodes represent clusters, and edges represent set containment. From Figure 4.5 of [WJ08]. Used with kind permission of Martin Wainwright.

### 9.3.5 Accuracy

For a graph with a single loop, one can show that the max-product version of LBP will find the correct MAP estimate, if it converges [Wei00]. For more general graphs, one can bound the error in the approximate marginals computed by LBP, as shown in [WJW03; IFW05; Vin+10b].

Much stronger results are available in the case of Gaussian models. In particular, it can be shown that, if the method converges, the means are exact, although the variances are not (typically the beliefs are over confident). See e.g., [WF01a; JMW06; Bic09; Du+18] for details.

### 9.3.6 Generalized belief propagation

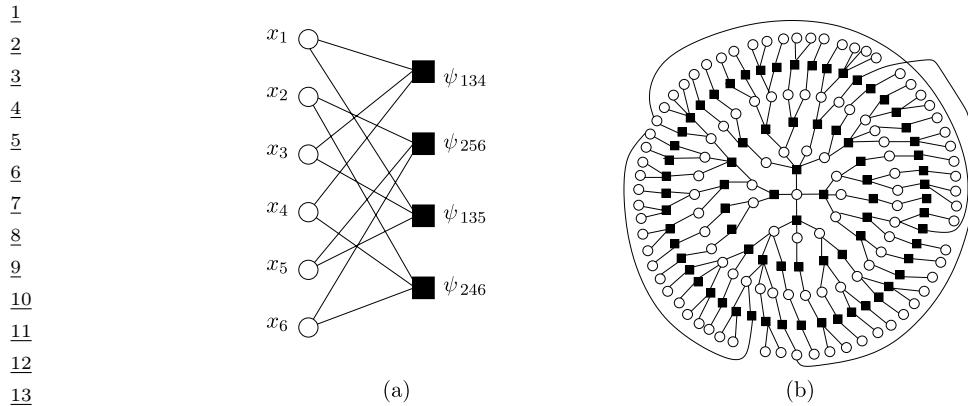
We can improve the accuracy of loopy BP by clustering together nodes that form a tight loop. This is known as the **cluster variational method**, or **generalized belief propagation** [YFW00].

The result of clustering is a hyper-graph, which is a graph where there are hyper-edges between sets of vertices instead of between single vertices. Note that a junction tree (Section 9.5) is a kind of hyper-graph. We can represent a hyper-graph using a poset (partially ordered set) diagram, where each node represents a hyper-edge, and there is an arrow  $e_1 \rightarrow e_2$  if  $e_2 \subset e_1$ . See Figure 9.8 for an example.

If we allow the size of the largest hyper-edge in the hyper-graph to be as large as the treewidth of the graph, then we can represent the hyper-graph as a tree, and the method will be exact, just as LBP is exact on regular trees (with treewidth 1). In this way, we can define a continuum of approximations, from LBP all the way to exact inference. See Supplementary Section 10.4.3.3 for more information.

### 9.3.7 Convex BP

In Supplementary Section 10.4.3 we analyse LBP from a variational perspective, and show that the resulting optimization problem, for both standard and generalized BP, is non-convex. However it is possible to create a version of **convex BP**, as we explain in Supplementary Section 10.4.4, which has the advantage that it will always converge.



14 *Figure 9.9: (a) A simple factor graph representation of a (2,3) low-density parity check code. Each message*  
 15 *bit (hollow round circle) is connected to two parity factors (solid black squares), and each parity factor is*  
 16 *connected to three bits. Each parity factor has the form  $\psi_{stu}(x_s, x_t, x_u) = \mathbb{I}(x_s \otimes x_t \otimes x_u = 1)$ , where  $\otimes$  is*  
 17 *the xor operator. The local evidence factors for each hidden node are not shown. (b) A larger example of a*  
 18 *random LDPC code. We see that this graph is “locally tree-like”, meaning there are no short cycles; rather,*  
 19 *each cycle has length  $\sim \log m$ , where  $m$  is the number of nodes. This gives us a hint as to why loopy BP*  
 20 *works so well on such graphs. (Note, however, that some error correcting code graphs have short loops, so this*  
 21 *is not the full explanation.) From Figure 2.9 from [WJ08]. Used with kind permission of Martin Wainwright.*

### 24 9.3.8 Application: error correcting codes

25 LBP was first proposed by Judea Pearl in his 1988 book [Pea88]. He recognized that applying BP to  
 26 loopy graphs might not work, but recommended it as a heuristic.

27 However, the main impetus behind the interest in LBP arose when McEliece, MacKay, and Cheng  
 28 [MMC98] showed that a popular algorithm for error correcting codes, known as **turbocodes** [BGT93],  
 29 could be viewed as an instance of LBP applied to a certain kind of graph.

30 We introduced error correcting codes in Section 5.5. Recall that the basic idea is to send the  
 31 source message  $\mathbf{x} \in \{0, 1\}^m$  over a noisy channel, and for the receiver to try to infer it given noisy  
 32 measurements  $\mathbf{y} \in \{0, 1\}^m$  or  $\mathbf{y} \in \mathbb{R}^m$ . That is, the receiver needs to compute  $\mathbf{x}^* = \operatorname{argmax}_{\mathbf{x}} p(\mathbf{x}|\mathbf{y}) =$   
 33  $\operatorname{argmax}_{\mathbf{x}} \tilde{p}(\mathbf{x})$ .

34 It is standard to represent  $\tilde{p}(\mathbf{x})$  as a factor graph (Section 4.6.1), which can easily represent any  
 35 deterministic relationships (parity constraints) between the bits. A factor graph is a bipartite graph  
 36 with  $x_i$  nodes on one side, and factors on the other. A graph in which each node is connected to  
 37  $n$  factors, and in which each factor is connected to  $k$  nodes, is called an  $(n, k)$  code. Figure 9.9(a)  
 38 shows a simple example of a (2,3) code, where each bit (hollow round circle) is connected to two  
 39 parity factors (solid black squares), and each parity factor is connected to three bits. Each parity  
 40 factor has the form

$$42 \quad \psi_{stu}(x_s, x_t, x_u) \triangleq \begin{cases} 1 & \text{if } x_s \otimes x_t \otimes x_u = 1 \\ 0 & \text{otherwise} \end{cases} \quad (9.32)$$

45 If the degrees of the parity checks and variable nodes remain bounded as the blocklength  $m$  increases,  
 46 this is called a **low-density parity check code**, or **LDPC code**. (Turbocodes are constructed in  
 47

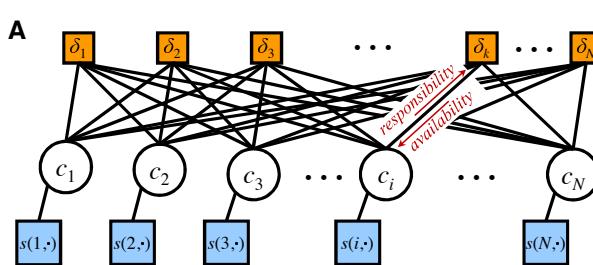


Figure 9.10: Factor graphs for affinity propagation. Circles are variables, squares are factors. Each  $c_i$  node has  $N$  possible states. From Figure S2 of [FD07a]. Used with kind permission of Brendan Frey.

a similar way.)

Figure 9.9(b) shows an example of a randomly constructed LDPC code. This graph is “locally tree-like”, meaning there are no short cycles; rather, each cycle has length  $\sim \log m$ . This fact is important to the success of LBP, which is only guaranteed to work on tree-structured graphs. Using methods such as these, people have been able to approach the lower bound in Shannon’s channel coding theorem, meaning they have produced codes with very little redundancy for a given amount of noise in the channel. See e.g., [MMC98; Mac03] for more details. Such codes are widely used, e.g., in modern cellphones.

### 9.3.9 Application: Affinity propagation

In this section, we discuss **affinity propagation** [FD07a], which can be seen as an improvement to K-medoids clustering, which takes as input a pairwise similarity matrix. The idea is that each data point must choose another data point as its exemplar or centroid; some data points will choose themselves as centroids, and this will automatically determine the number of clusters. More precisely, let  $c_i \in \{1, \dots, N\}$  represent the centroid for datapoint  $i$ . The goal is to maximize the following function

$$J(\mathbf{c}) = \sum_{i=1}^N S(i, c_i) + \sum_{k=1}^N \delta_k(\mathbf{c}) \quad (9.33)$$

where  $S(i, c_i)$  is the similarity between data point  $i$  and its centroid  $c_i$ . The second term is a penalty term that is  $-\infty$  if some data point  $i$  has chosen  $k$  as its exemplar (i.e.,  $c_i = k$ ), but  $k$  has not chosen itself as an exemplar (i.e., we do not have  $c_k = k$ ). More formally,

$$\delta_k(\mathbf{c}) = \begin{cases} -\infty & \text{if } c_k \neq k \text{ but } \exists i : c_i = k \\ 0 & \text{otherwise} \end{cases} \quad (9.34)$$

This encourages “representative” samples to vote for themselves as centroids, thus encouraging clustering behavior.

The objective function can be represented as a factor graph. We can either use  $N$  nodes, each with  $N$  possible values, as shown in Figure 9.10, or we can use  $N^2$  binary nodes (see [GF09] for the details). We will assume the former representation.

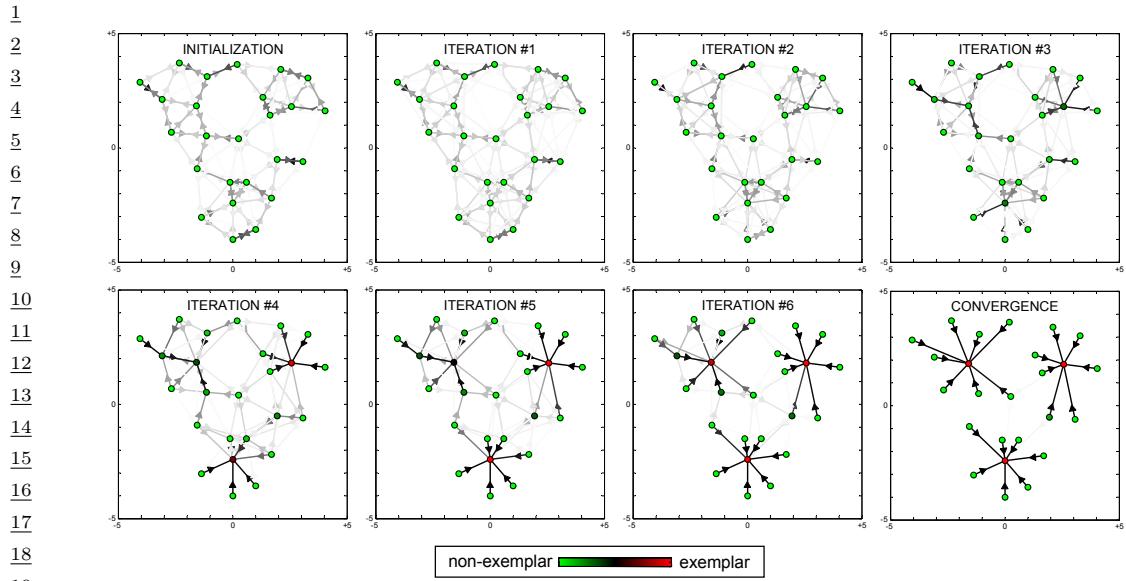


Figure 9.11: Example of affinity propagation. Each point is colored coded by how much it wants to be an exemplar (red is the most, green is the least). This can be computed by summing up all the incoming availability messages and the self-similarity term. The darkness of the  $i \rightarrow k$  arrow reflects how much point  $i$  wants to belong to exemplar  $k$ . From Figure 1 of [FD07a]. Used with kind permission of Brendan Frey.

We can find a strong local maximum of the objective by using max-product loopy belief propagation (Section 9.3). Referring to the model in Figure 9.10, each variable node  $c_i$  sends a message to each factor node  $\delta_k$ . It turns out that this vector of  $N$  numbers can be reduced to a scalar message, denoted  $r_{i \rightarrow k}$ , known as the responsibility. This is a measure of how much  $i$  thinks  $k$  would make a good exemplar, compared to all the other exemplars  $i$  has looked at. In addition, each factor node  $\delta_k$  sends a message to each variable node  $c_i$ . Again this can be reduced to a scalar message,  $a_{i \leftarrow k}$ , known as the availability. This is a measure of how strongly  $k$  believes it should be an exemplar for  $i$ , based on all the other data points  $k$  has looked at.

As usual with loopy BP, the method might oscillate, and convergence is not guaranteed. However, by using damping, the method is very reliable in practice. If the graph is densely connected, message passing takes  $O(N^2)$  time, but with sparse similarity matrices, it only takes  $O(E)$  time, where  $E$  is the number of edges or non-zero entries in  $S$ .

The number of clusters can be controlled by scaling the diagonal terms  $S(i, i)$ , which reflect how much each data point wants to be an exemplar. Figure 9.11 gives a simple example of some 2d data, where the negative Euclidean distance was used to measured similarity. The  $S(i, i)$  values were set to be the median of all the pairwise similarities. The result is 3 clusters. Many other results are reported in [FD07a], who show that the method significantly outperforms K-medoids.

---

### 9.3.10 Emulating BP with graph neural nets

There is a close connection between message passing in PGMs and message passing in graph neural networks (GNNs), which we discuss in Section 16.3.6. However, for PGMs, the message computations are computing using (non-learned) update equations that work for any model; all that is needed is the graph structure  $G$ , model parameters  $\theta$ , and evidence  $v$ . By contrast, GNNs are trained to emulate specific functions using labeled input-output pairs.

It is natural to wonder what happens if we train a GNN on the exact posterior marginals derived from a small PGM, and then apply that trained GNN to a different test PGM. In [Yoo+18; Zha+19d], they show this method can work quite well if the test PGM is similar in structure to the one used for training.

An alternative approach is to start with a known PGM, and then “unroll” the BP message passing algorithm to produce a layered feedforward model, whose connectivity is derived from the graph. The resulting network can then be trained discriminatively for some end-task (not necessarily computing posterior marginals). Thus the BP procedure applied to the PGM just provides a way to design the neural network structure. This method is called **deep unfolding** (see e.g., [HLRW14]), and can often give very good results. (See also [SW20] for a more recent version of this approach, called **“neural enhanced BP”**.)

These neural methods are useful if the PGM is fixed, and we want to repeatedly perform inference or prediction with it, using different values of the evidence, but where the set of nodes which are observed is always the same. This is an example of amortized inference, where we train a model to emulate the results of running an iterative optimization scheme (see Section 10.3.6 for more discussion).

## 9.4 The variable elimination (VE) algorithm

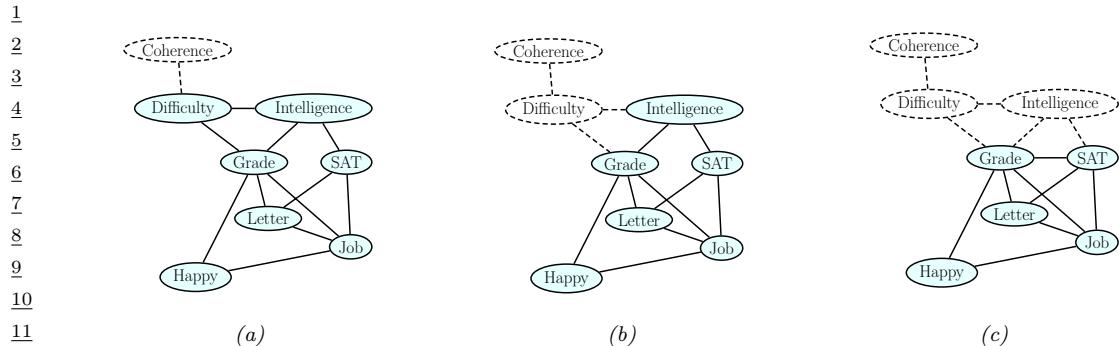
In this section, we discuss an algorithm to compute a posterior marginal  $p(\mathbf{z}_Q|\mathbf{y})$  for any query set  $Q$ , assuming  $p$  is defined by a graphical model. Unlike loopy BP, it is guaranteed to give the correct answers even if the graph has cycles. We assume all the hidden nodes are discrete, although a version of the algorithm can be created for the Gaussian case by using the rules for sum and product defined in Section 2.2.7.

### 9.4.1 Derivation of the algorithm

We will explain the algorithm by applying it to an example. Specifically, we consider the student network from Section 4.2.2.2. Suppose we want to compute  $p(J=1)$ , the marginal probability that a person will get a job. Since we have 8 binary variables, we could simply enumerate over all possible assignments to all the variables (except for  $J$ ), adding up the probability of each joint instantiation:

$$p(J) = \sum_L \sum_S \sum_G \sum_H \sum_I \sum_D \sum_C p(C, D, I, G, S, L, J, H) \quad (9.35)$$

However, this would take  $O(2^7)$  time. We can be smarter by **pushing sums inside products**. This is the key idea behind the **variable elimination** algorithm [ZP96], also called **bucket elimination** [Dec96], or, in the context of genetic pedigree trees, the **peeling algorithm** [CTS78].



*Figure 9.12: Example of the elimination process, in the order C, D, I, H, G, S, L. When we eliminate I (figure c), we add a fill-in edge between G and S, since they are not connected. Adapted from Figure 9.10 of [KF09a].*

In our example, we get

$$\begin{aligned}
p(J) &= \sum_{L,S,G,H,I,D,C} p(C, D, I, G, S, L, J, H) \\
&= \sum_{L,S,G,H,I,D,C} \psi_C(C)\psi_D(D,C)\psi_I(I)\psi_G(G,I,D)\psi_S(S,I)\psi_L(L,G) \\
&\quad \times \psi_J(J,L,S)\psi_H(H,G,J) \\
&= \sum_{L,S} \psi_J(J,L,S) \sum_G \psi_L(L,G) \sum_H \psi_H(H,G,J) \sum_I \psi_S(S,I)\psi_I(I) \\
&\quad \times \sum_D \psi_G(G,I,D) \sum_C \psi_C(C)\psi_D(D,C)
\end{aligned}$$

We now evaluate this expression, working right to left as shown in Table 9.1. First we multiply together all the terms in the scope of the  $\sum_{\gamma}$  operator to create the temporary factor

$$\tau'_*(C, D) = \psi_C(C)\psi_D(D, C) \quad (9.36)$$

<sup>34</sup> Then we marginalize out  $C$  to get the new factor.

$$\frac{36}{37} \quad \tau_1(D) = \sum_C \tau'_1(C, D) \quad (9.37)$$

<sup>39</sup> Next we multiply together all the terms in the scope of the  $\sum_D$  operator and then marginalize out <sup>40</sup> to create

$$\frac{41}{42} \cdot ((G, L, E)) = ((G, L, E) - (D)) \quad (2.22)$$

$$\begin{aligned} \frac{43}{44} \quad \tau_2(G, I, D) &= \psi_G(G, I, D) \cdot \tau_1(D) \\ \frac{43}{44} \quad \tau_2(G, I) &= \sum \tau'_2(G, I, D) \end{aligned} \quad (9.39)$$

45

---


$$\begin{aligned}
& \sum_L \sum_S \psi_J(J, L, S) \sum_G \psi_L(L, G) \sum_H \psi_H(H, G, J) \sum_I \psi_S(S, I) \psi_I(I) \sum_D \psi_G(G, I, D) \underbrace{\sum_C \psi_C(C) \psi_D(D, C)}_{\tau_1(D)} \\
& \sum_L \sum_S \psi_J(J, L, S) \sum_G \psi_L(L, G) \sum_H \psi_H(H, G, J) \sum_I \psi_S(S, I) \psi_I(I) \underbrace{\sum_D \psi_G(G, I, D) \tau_1(D)}_{\tau_2(G, I)} \\
& \sum_L \sum_S \psi_J(J, L, S) \sum_G \psi_L(L, G) \sum_H \psi_H(H, G, J) \underbrace{\sum_I \psi_S(S, I) \psi_I(I) \tau_2(G, I)}_{\tau_3(G, S)} \\
& \sum_L \sum_S \psi_J(J, L, S) \sum_G \psi_L(L, G) \underbrace{\sum_H \psi_H(H, G, J) \tau_3(G, S)}_{\tau_4(G, J)} \\
& \sum_L \sum_S \psi_J(J, L, S) \underbrace{\sum_G \psi_L(L, G) \tau_4(G, J) \tau_3(G, S)}_{\tau_5(J, L, S)} \\
& \sum_L \underbrace{\sum_S \psi_J(J, L, S) \tau_5(J, L, S)}_{\tau_6(J, L)} \\
& \sum_L \underbrace{\tau_6(J, L)}_{\tau_7(J)}
\end{aligned}$$

---

Table 9.1: Eliminating variables from Figure 4.38 in the order  $C, D, I, H, G, S, L$  to compute  $P(J)$ .

The above technique can be used to compute any marginal of interest, such as  $p(J)$  or  $p(J, H)$ . To compute a conditional, we can take a ratio of two marginals, where the visible variables have been clamped to their known values (and hence don't need to be summed over). For example,

$$p(J = j | I = 1, H = 0) = \frac{p(J = j, I = 1, H = 0)}{\sum_{j'} p(J = j', I = 1, H = 0)} \quad (9.40)$$

### 9.4.2 Computational complexity of VE

The running time of VE is clearly exponential in the size of the largest factor, since we have to sum over all of the corresponding variables. Some of the factors come from the original model (and are thus unavoidable), but new factors may also be created in the process of summing out. For example, in Table 9.1, we created a factor involving  $G, I$  and  $S$ ; but these nodes were not originally present together in any factor.

The order in which we perform the summation is known as the **elimination order**. This can have a large impact on the size of the intermediate factors that are created. For example, consider the ordering in Table 9.1: the largest created factor (beyond the original ones in the model) has size 3, corresponding to  $\tau_5(J, L, S)$ . Now consider the ordering in Table 9.2: now the largest factors are  $\tau_1(I, D, L, J, H)$  and  $\tau_2(D, L, S, J, H)$ , which are much bigger.

---

1  
 2      $\sum_D \sum_C \psi_D(D, C) \sum_H \sum_L \sum_S \psi_J(J, L, S) \sum_I \underbrace{\psi_I(I) \psi_S(S, I) \sum_G \psi_G(G, I, D) \psi_L(L, G) \psi_H(H, G, J)}_{\tau_1(I, D, L, J, H)}$   
 3  
 4  
 5      $\sum_D \sum_C \psi_D(D, C) \sum_H \sum_L \sum_S \psi_J(J, L, S) \underbrace{\sum_I \psi_I(I) \psi_S(S, I) \tau_1(I, D, L, J, H)}_{\tau_2(D, L, S, J, H)}$   
 6  
 7  
 8      $\sum_D \sum_C \psi_D(D, C) \sum_H \sum_L \underbrace{\sum_S \psi_J(J, L, S) \tau_2(D, L, S, J, H)}_{\tau_3(D, L, J, H)}$   
 9  
 10  
 11      $\sum_D \sum_C \psi_D(D, C) \sum_H \underbrace{\sum_L \tau_3(D, L, J, H)}_{\tau_4(D, J, H)}$   
 12  
 13  
 14      $\sum_D \sum_C \psi_D(D, C) \underbrace{\sum_H \tau_4(D, J, H)}_{\tau_5(D, J)}$   
 15  
 16  
 17      $\sum_D \sum_C \underbrace{\psi_D(D, C) \tau_5(D, J)}_{\tau_6(D, J)}$   
 18  
 19  
 20      $\underbrace{\sum_D \tau_6(D, J)}_{\tau_7(J)}$   
 21  
 22  
 23

---

Table 9.2: Eliminating variables from Figure 4.38 in the order  $G, I, S, L, H, C, D$ .

24  
 25  
 26  
 27  
 28     We can determine the size of the largest factor graphically, without worrying about the actual  
 29     numerical values of the factors, by running the VE algorithm “symbolically”. When we eliminate a  
 30     variable  $z_t$ , we connect together all variables that share a factor with  $z_t$  (to reflect the new temporary  
 31     factor  $\tau'_t$ ). The edges created by this process are called **fill-in edges**. For example, Figure 9.12  
 32     shows the fill-in edges introduced when we eliminate in the  $C, D, I, \dots$  order. The first two steps do  
 33     not introduce any fill-ins, but when we eliminate  $I$ , we connect  $G$  and  $S$ , to capture the temporary  
 34     factor

35  
 36      $\tau'_3(G, S, I) = \psi_S(S, I) \psi_I(I) \tau_2(G, I)$  (9.41)

37  
 38     Let  $\mathcal{G}_\prec$  be the (undirected) graph induced by applying variable elimination to  $\mathcal{G}$  using elimination  
 39     ordering  $\prec$ . The temporary factors generated by VE correspond to maximal **cliques** in the graph  
 40      $\mathcal{G}_\prec$ . For example, with ordering  $(C, D, I, H, G, S, L)$ , the maximal cliques are as follows:

41      $\{C, D\}, \{D, I, G\}, \{G, L, S, J\}, \{G, J, H\}, \{G, I, S\}$  (9.42)

43     It is clear that the time complexity of VE is

44  
 45      $\sum_{c \in \mathcal{C}(G_\prec)} K^{|c|}$  (9.43)

where  $\mathcal{C}(\mathcal{G})$  are the (maximal) cliques in graph  $\mathcal{G}$ ,  $|c|$  is the size of the clique  $c$ , and we assume for notational simplicity that all the variables have  $K$  states each.

Let us define the **induced width** of a graph given elimination ordering  $\prec$ , denoted  $w_\prec$ , as the size of the largest factor (i.e., the largest clique in the induced graph) minus 1. Then it is easy to see that the complexity of VE with ordering  $\prec$  is  $O(K^{w_\prec+1})$ . The smallest possible induced width for a graph is known as its **treewidth**. Unfortunately finding the corresponding optimal elimination order is an NP-complete problem [Yan81; ACP87]. See Section 9.4.3 for a discussion of some approximate methods for finding good elimination orders.

### 9.4.3 Picking a good elimination order

Many algorithms take time (or space) which is exponential in the tree width of the corresponding graph. For example, this applies to Cholesky decompositions of sparse matrices, as well as to einsum contractions (see [https://github.com/dgasmith/opt\\_einsum](https://github.com/dgasmith/opt_einsum)). Hence we would like to find an elimination ordering that minimizes the width. We say that an ordering  $\pi$  is a **perfect elimination ordering** if it does not introduce any fill-in edges. Every graph that is already triangulated (e.g., a tree) has a perfect elimination ordering. We call such graphs **decomposable**.

In general, we will need to add fill-in edges to ensure the resulting graph is decomposable. Different orderings can introduce different numbers of fill-in edges, which affects the width of the resulting chordal graph; for example, compare Table 9.1 to Table 9.2.

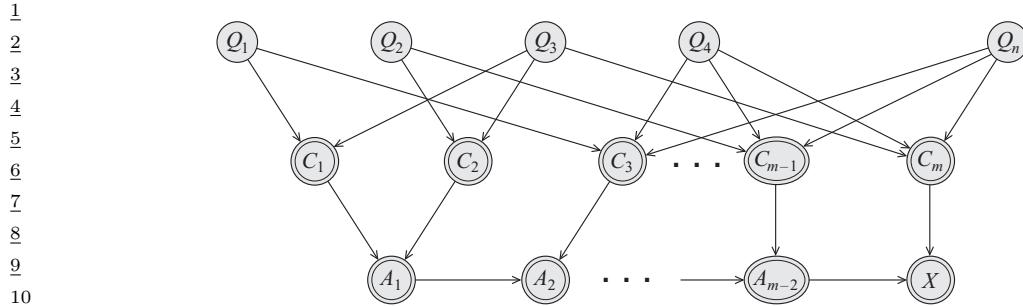
Choosing an elimination ordering with minimal width is NP-complete [Yan81; ACP87]. It is common to use greedy approximation known as the **min-fill heuristic**, which works as follows: eliminate any node which would not result in any fill-ins (i.e., all of whose uneliminated neighbors already form a clique); if there is no such node, eliminate the node which would result in the minimum number of fill-in edges. When nodes have different weights (e.g., representing different numbers of states), we can use the **min-weight heuristic**, where we try to minimize the weight of the created cliques at each step.

Of course, many other methods are possible. See [Heg06] for a general survey. [Kja90; Kja92] compared simulated annealing with the above greedy method, and found that it sometimes works better (although it is much slower). [MJ97] approximate the discrete optimization problem by a continuous optimization problem. [BG96] present a randomized approximation algorithm. [Gil88] present the nested dissection order, which is always within  $O(\log N)$  of optimal. [Ami01] discuss various constant-factor appoximation algorithms. [Dav+04] present the **AMD** or approximate minimum degree ordering algorithm, which is implemented in Matlab.<sup>3</sup> The **METIS** library can be used for finding elimination orderings for large graphs; this implements the **nested dissection** algorithm [GT86]. For a planar graph with  $N$  nodes, the resulting treewidth will have the optimal size of  $O(N^{3/2})$ .

### 9.4.4 Computational complexity of exact inference

We have seen that variable elimination takes  $O(NK^{w+1})$  time to compute the marginals for a graph with  $N$  nodes, and treewidth  $w$ , where each variable has  $K$  states. If the graph is densely connected, then  $w = O(N)$ , and so inference will take time exponential in  $N$ .

<sup>3</sup> See the description of the symamd command at <https://bit.ly/31N6E2b>. (“sym” stands for symbolic, “amd” stands approximate minimum degree.)



*Figure 9.13: Encoding a 3-SAT problem on  $n$  variables and  $m$  clauses as a DGM. The  $Q_s$  variables are binary random variables. The  $C_t$  variables are deterministic functions of the  $Q_s$ 's, and compute the truth value of each clause. The  $A_t$  nodes are a chain of AND gates, to ensure that the CPT for the final  $x$  node has bounded size. The double rings denote nodes with deterministic CPDs. From Figure 9.1 of [KF09a]. Used with kind permission of Daphne Koller.*

16

17

18

19 Of course, just because some particular algorithm is slow doesn't mean that there isn't some  
20 smarter algorithm out there. Unfortunately, this seems unlikely, since it is easy to show that exact  
21 inference for discrete graphical models is NP-hard [DL93]. The proof is a simple reduction from the  
22 satisfiability problem. In particular, note that we can encode any 3-SAT problem as a DPGM with  
23 deterministic links, as shown in Figure 9.13. We clamp the final node,  $x$ , to be on, and we arrange  
24 the CPTs so that  $p(x = 1) > 0$  iff there is a satisfying assignment. Computing any posterior marginal  
25 requires evaluating the normalization constant,  $p(x = 1)$ , so inference in this model implicitly solves  
26 the SAT problem.

27 In fact, exact inference is #P-hard [Rot96], which is even harder than NP-hard. The intuitive  
28 reason for this is that to compute the normalizing constant, we have to *count* how many satisfying  
29 assignments there are. (By contrast, MAP estimation is provably easier for some model classes  
30 [GPS89], since, intuitively speaking, it only requires finding one satisfying assignment, not counting  
31 all of them.) Furthermore, even approximate inference is computationally hard in general [DL93;  
32 Rot96].

33 The above discussion was just concerned with inferring the states of discrete hidden variables.  
34 When we have continuous hidden variables, the problem can be even harder, since even a simple  
35 two-node graph, of the form  $z \rightarrow y$ , can be intractable to invert if the variables are high dimensional  
36 and do not have a conjugate relationship (Section 3.2). Inference in mixed discrete-continuous models  
37 can also be hard [LP01].

38 As a consequence of these hardness results, we often have to resort to approximate inference  
39 methods, such as variational inference (Chapter 10) and Monte Carlo inference (Chapter 11).

40

#### 41 9.4.5 Drawbacks of VE

42

43 Consider using VE to compute all the marginals in a chain-structured graphical model, such as an  
44 HMM. We can easily compute the final marginal  $p(z_T | \mathbf{y})$  by eliminating all the nodes  $z_1$  to  $z_{T-1}$   
45 in order. This is equivalent to the forwards algorithm, and takes  $O(K^2 T)$  time, as we discussed in  
46 Section 8.2.4. But now suppose we want to compute  $p(z_{T-1} | \mathbf{y})$ . We have to run VE again, at a cost  
47

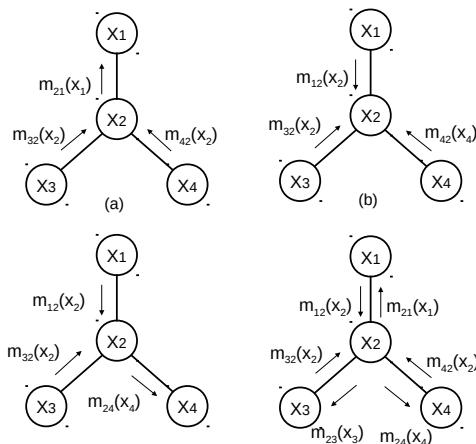


Figure 9.14: Sending multiple messages along a tree. (a)  $z_1$  is root. (b)  $z_2$  is root. (c)  $z_4$  is root. (d) All of the messages needed to compute all singleton marginals. Adapted from Figure 4.3 of [Jor07].

of  $O(K^2T)$  time. So the total cost to compute all the marginals is  $O(K^2T^2)$ . However, we know that we can solve this problem in  $O(K^2T)$  using the forwards-backwards, as we discussed in Section 8.2.4. The difference is that FB caches the messages computed on the forwards pass, so it can reuse them later. (Caching previously computed results is the core idea behind dynamic programming.)

The same problem arises when applying VE to trees. For example, consider the 4-node tree in Figure 9.14. We can compute  $p(z_1|\mathbf{y})$  by eliminating  $z_{2:4}$ ; this is equivalent to sending messages up to  $z_1$  (the messages correspond to the  $\tau$  factors created by VE). Similarly we can compute  $p(z_2|\mathbf{y})$ ,  $p(z_3|\mathbf{y})$  and then  $p(z_4|\mathbf{y})$ . We see that some of the messages used to compute the marginal on one node can be re-used to compute the marginals on the other nodes. By storing the messages for later re-use, we can compute all the marginals in  $O(K^2T)$  time, as we show in Section 9.2.

The question is: how do we get these benefits of message passing on a tree when the graph is not a tree? We give the answer in Section 9.5.

## 9.5 The junction tree algorithm (JTA)

The **junction tree algorithm** or **JTA** is a generalization of variable elimination that lets us efficiently compute all the posterior marginals without repeating redundant work, by using dynamic programming, thus avoiding the problems mentioned in Section 9.4.5. The basic idea is to convert the graph into a special kind of tree, known as a **junction tree** (also called a **join tree**, or **clique tree**), and then to run belief propagation (message passing) on this tree. We can create the join tree by running variable elimination “symbolically”, as discussed in Section 9.4.2, and adding the generated fill-in edges to the graph. The resulting chordal graph can then be converted to a tree, as explained in Supplementary Section 9.2.1. Once we have a tree, we can perform message passing on it, using a variant of the method Section 9.2.2. See Supplementary Section 9.2.2 for details.

---

1 **9.6 Inference as optimization**

3 In this section, we discuss how to perform posterior inference by solving an optimization problem,  
4 which is often computationally simpler. See also [Supplementary](#) Section 9.3.  
5

6 **9.6.1 Inference as backpropagation**

8 In this section, we discuss how to compute posterior marginals in a graphical model using automatic  
9 differentiation. For notational simplicity, we focus on undirected graphical models, where the joint  
10 can be represented as an exponential family (Section 2.3) follows:  
11

$$\underline{12} \quad p(\mathbf{x}) = \frac{1}{Z} \prod_c \psi_c(\mathbf{x}_c) = \exp\left(\sum_c \boldsymbol{\eta}_c^\top \mathcal{T}(\mathbf{x}_c) - \log A(\boldsymbol{\eta})\right) = \exp(\boldsymbol{\eta}^\top \mathcal{T}(\mathbf{x}) - \log A(\boldsymbol{\eta})) \quad (9.44)$$

13

15 where  $\psi_c$  is the potential function for clique  $c$ ,  $\boldsymbol{\eta}_c$  are the natural parameters for clique  $c$ ,  $\mathcal{T}(\mathbf{x}_c)$  are  
16 the corresponding sufficient statistics, and  $A = \log Z$  is the log partition function.

17 We will consider pairwise models (with node and edge potentials), and discrete variables. The  
18 natural parameters are the node and edge log potentials,  $\boldsymbol{\eta} = (\{\eta_{s;j}\}, \{\eta_{s,t;j,k}\})$ , and the sufficient  
19 statistics are node and edge indicator functions,  $\mathcal{T}(\mathbf{x}) = (\{\mathbb{I}(x_s = j)\}, \{\mathbb{I}(x_s = j, x_t = k)\})$ . (Note:  
20 we use  $s, t \in \mathcal{V}$  to index nodes and  $j, k \in \mathcal{X}$  to index states.)

21 The mean of the sufficient statistics are given by

$$\underline{23} \quad \boldsymbol{\mu} = \mathbb{E}[\mathcal{T}(\mathbf{x})] = (\{p(x_s = j)\}_s, \{p(x_s = j, x_t = k)\}_{s \neq t}) = (\{\mu_{s;j}\}_s, \{\mu_{s,t;j,k}\}_{s \neq t}) \quad (9.45)$$

24

25 The key result, from Equation (2.219), is that  $\boldsymbol{\mu} = \nabla_{\boldsymbol{\eta}} A(\boldsymbol{\eta})$ . Thus as long as we have a function that  
26 computes  $A(\boldsymbol{\eta}) = \log Z(\boldsymbol{\eta})$ , we can use automatic differentiation (Section 6.2) to compute gradients,  
27 and then we can extract the corresponding node marginals from the gradient vector. If we have  
28 evidence (known values) on some of the variables, we simply “clamp” the corresponding entries to 0  
29 or 1 in the node potentials.

30 The observation that probabilistic inference can be performed using automatic differentiation has  
31 been discovered independently by several groups (e.g., [[Dar03](#); [PD03](#); [Eis16](#); [ASM17](#)]). It also lends  
32 itself to the development of differentiable approximations to inference (see e.g., [[MB18](#)]).

33 **9.6.1.1 Example: inference in a small model**

35 As a concrete example, consider a small chain structured model  $x_1 - x_2 - x_3$ , where each node has  
36  $K$  states. We can represent the node potentials as  $K \times 1$  tensors (table of numbers), and the edge  
37 potentials by  $K \times K$  tensors. The partition function is given by  
38

$$\underline{39} \quad Z(\boldsymbol{\psi}) = \sum_{x_1, x_2, x_3} \psi_1(x_1) \psi_2(x_2) \psi_3(x_3) \psi_{12}(x_1, x_2) \psi_{23}(x_2, x_3) \quad (9.46)$$

40

42 Let  $\boldsymbol{\eta} = \log(\boldsymbol{\psi})$  be the log potentials, and  $A(\boldsymbol{\eta}) = \log Z(\boldsymbol{\eta})$  be the log partition function. We can  
43 compute the single node marginals  $\boldsymbol{\mu}_s = p(x_s = 1 : K)$  using  $\boldsymbol{\mu}_s = \nabla_{\boldsymbol{\eta}_s} A(\boldsymbol{\eta})$ , and the pairwise  
44 marginals  $\boldsymbol{\mu}_{s,t}(j, k) = p(x_s = j, x_t = k)$  using  $\boldsymbol{\mu}_{s,t} = \nabla_{\boldsymbol{\eta}_{s,t}} A(\boldsymbol{\eta})$ .

45 We can compute the partition function  $Z$  efficiently use numpy’s **einsum** function, which imple-  
46 ments tensor contraction using Einstein summation notation. We label each dimension of the tensors  
47

by A, B and C, so einsum knows how to match things up. We then compute gradients using an auto-diff library.<sup>4</sup> The result is that inference can be done in two lines of Python code, as shown in Listing 9.1:

*Listing 9.1: Computing marginals from derivative of log partition function*

```

import jax.numpy as jnp
from jax import grad

logZ_fun = lambda logpots: np.log(jnp.einsum("A,B,C,AB,BC",
                                              *[jnp.exp(lp) for lp in logpots]))
probs = grad(logZ_fun)(logpots)

```

To perform conditional inference, such as  $p(x_s = k|x_t = e)$ , we multiply in one-hot indicator vectors to clamp  $x_t$  to the value  $e$  so that the unnormalized joint only assigns non-zero probability to state combinations that are valid. We then sum over all values of the unclamped variables to get the constrained partition function  $Z_e$ . The gradients will now give us the marginals conditioned on the evidence [Dar03].

## 9.6.2 Perturb and MAP

In this section, we discuss how to draw posterior samples from a graphical model by leveraging optimization as a subroutine. The basic idea is to make  $S$  copies of the model, each of which has slightly perturbed versions of the parameters,  $\theta_s = \theta_s + \epsilon_s$ , and then to compute the MAP estimate,  $\mathbf{x}_s = \text{argmax } p(\mathbf{x}|\mathbf{y}; \theta_s)$ . For a suitably chosen noise distribution for  $\epsilon_s$ , this technique — known as **perturb-and-MAP** — can be shown that this gives exact posterior samples [PY10; PY11; PY14].

### 9.6.2.1 Gaussian case

We first consider the case of a Gaussian MRF. Let  $\mathbf{x} \in \mathbb{R}^N$  be the vector of hidden states with prior

$$p(\mathbf{x}) \propto \mathcal{N}(\mathbf{G}\mathbf{x}|\boldsymbol{\mu}_p, \boldsymbol{\Sigma}_p) \propto \exp\left(-\frac{1}{2}\mathbf{x}^\top \mathbf{K}_x \mathbf{x} + \mathbf{h}_x^\top \mathbf{x}\right) \quad (9.47)$$

where  $\mathbf{G} \in \mathbb{R}^{K \times N}$  is a matrix that represents prior dependencies (e.g., pairwise correlations),  $\mathbf{K}_x = \mathbf{G}^\top \boldsymbol{\Sigma}_p^{-1} \mathbf{G}$ , and  $\mathbf{h}_x = \mathbf{G}^\top \boldsymbol{\Sigma}_p^{-1} \boldsymbol{\mu}_p$ . Let  $\mathbf{y} \in \mathbb{R}^M$  be the measurements with likelihood

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{H}\mathbf{x} + \mathbf{c}, \boldsymbol{\Sigma}_n) \propto \exp\left(-\frac{1}{2}\mathbf{x}^\top \mathbf{K}_{y|x} \mathbf{x} + \mathbf{h}_{y|x}^\top \mathbf{x} - \frac{1}{2}\mathbf{y}^\top \boldsymbol{\Sigma}_n^{-1} \mathbf{y}\right) \quad (9.48)$$

where  $\mathbf{H} \in \mathbb{R}^{M \times N}$  represents dependencies between the hidden and visible variables,  $\mathbf{K}_{y|x} = \mathbf{H}^\top \boldsymbol{\Sigma}_n^{-1} \mathbf{H}$  and  $\mathbf{h}_{y|x} = \mathbf{H}^\top \boldsymbol{\Sigma}_n^{-1} (\mathbf{y} - \mathbf{c})$ . The posterior is given by the following (c.f., one step of the information filter in Section 8.3.4)

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (9.49)$$

$$\boldsymbol{\Sigma}^{-1} = \mathbf{K} = \mathbf{G}^\top \boldsymbol{\Sigma}_p^{-1} \mathbf{G} + \mathbf{H}^\top \boldsymbol{\Sigma}_n^{-1} \mathbf{H} \quad (9.50)$$

$$\boldsymbol{\mu} = \mathbf{K}(\mathbf{G}^\top \boldsymbol{\Sigma}_p^{-1} \boldsymbol{\mu}_p + \mathbf{H}^\top \boldsymbol{\Sigma}_n^{-1} (\mathbf{y} - \mathbf{c})) \quad (9.51)$$

4. See [ugm\\_inf\\_autodiff.py](#) for the full (JAX) code, and see <https://github.com/srush/ProbTalk> for a (PyTorch) version by Sasha Rush.

where we have assumed  $\mathbf{K} = \mathbf{K}_x + \mathbf{K}_{y|x}$  is invertible (although the prior or likelihood on their own may be singular).

The  $K$  rows of  $\mathbf{G} = [\mathbf{g}_1^\top; \dots; \mathbf{g}_K^\top]$  and the  $M$  rows of  $\mathbf{H} = [\mathbf{h}_1^\top; \dots; \mathbf{h}_M^\top]$  can be combined into the  $L$  rows of  $\mathbf{F} = [\mathbf{f}_1^\top; \dots; \mathbf{f}_L^\top]$ , which define the linear constraints of the system. If we assume that  $\Sigma_p$  and  $\Sigma_n$  are diagonal, then the structure of the graphical model is uniquely determined by the sparsity of  $\mathbf{F}$ . The resulting posterior factorizes as a product of  $L$  Gaussian “experts”:

$$p(\mathbf{x}|\mathbf{y}) \propto \prod_{l=1}^L \exp\left(-\frac{1}{2} \mathbf{x}^\top \mathbf{K}_l \mathbf{x} + \mathbf{h}_l^\top \mathbf{x}\right) \propto \prod_{l=1}^L \mathcal{N}(\mathbf{f}_l^\top \mathbf{x}; \mu_l, \Sigma_l) \quad (9.52)$$

where  $\Sigma_l$  equals  $\Sigma_{p,l,l}$  for  $l = 1 : K$  and equals  $\Sigma_{n,l',l'}$  for  $l = K+1 : L$  where  $l' = l - K$ . Similarly  $\mu_l = \mu_{p,l}$  for  $l = 1 : K$  and  $\mu_l = (y_{l'} - c_{l'})$  for  $l = K+1 : L$ .

To apply perturb-and-MAP, we proceed as follows. First perturb the prior mean by sampling  $\tilde{\boldsymbol{\mu}}_p \sim \mathcal{N}(\boldsymbol{\mu}_p, \Sigma_p)$ , and perturb the measurements by sampling  $\tilde{\mathbf{y}} \sim \mathcal{N}(\mathbf{y}, \Sigma_n)$ . (Note that this is equivalent to first perturbing the linear term in each information form potential, using  $\tilde{\mathbf{h}}_l = \mathbf{h}_l + \mathbf{f}_l \Sigma_l^{-\frac{1}{2}} \epsilon_l$ , where  $\epsilon_l \sim \mathcal{N}(0, 1)$ .) Then compute the MAP estimate for  $\mathbf{x}$  using the perturbed parameters:

$$\tilde{\mathbf{x}} = \mathbf{K}^{-1} \mathbf{G}^\top \Sigma_p^{-1} \tilde{\boldsymbol{\mu}}_p + \mathbf{K}^{-1} \mathbf{H}^\top \Sigma_n^{-1} (\tilde{\mathbf{y}} - \mathbf{c}) \quad (9.53)$$

$$= \underbrace{\mathbf{K}^{-1} \mathbf{G}^\top \Sigma_p^{-1} (\boldsymbol{\mu}_p + \boldsymbol{\epsilon}_\mu)}_{\mathbf{A}} + \underbrace{\mathbf{K}^{-1} \mathbf{H}^\top \Sigma_n^{-1} (\mathbf{y} + \boldsymbol{\epsilon}_y - \mathbf{c})}_{\mathbf{B}} \quad (9.54)$$

$$= \boldsymbol{\mu} + \mathbf{A} \boldsymbol{\epsilon}_\mu + \mathbf{B} \boldsymbol{\epsilon}_y \quad (9.55)$$

We see that  $\mathbb{E}[\tilde{\mathbf{x}}] = \boldsymbol{\mu}$  and  $\mathbb{E}[(\tilde{\mathbf{x}} - \boldsymbol{\mu})(\tilde{\mathbf{x}} - \boldsymbol{\mu})^\top] = \mathbf{K}^{-1} = \Sigma$ , so the method produces exact samples.

This approach is very scalable, since computing the MAP estimate of sparse GMRFs (i.e., posterior mean) can be done efficiently using conjugate gradient solvers. Alternatively we can use loopy belief propagation (Section 9.3), which can often compute the exact posterior mean (see e.g., [WF01a; JMW06; Bic09; Du+18]).

### 9.6.2.2 Discrete case

In [PY11; PY14] they extend perturb-and-MAP to the case of discrete graphical models. This setup is more complicated, and requires the use of Gumbel noise, which can be sampled using  $\epsilon = -\log(-\log(u))$ , where  $u \sim \text{Unif}(0, 1)$ . This noise should be added to all the potentials in the model, but as a simple approximation, it can just be added to the unary terms, i.e., the local evidence potentials. Let the score, or unnormalized log probability, of configuration  $\mathbf{x}$  given inputs  $\mathbf{c}$  be

$$S(\mathbf{x}; \mathbf{c}) = \log p(\mathbf{x}|\mathbf{c}) + \text{const} = \sum_i \log \phi_i(x_i) + \sum_{ij} \log \psi_{ij}(x_{i,j}) \quad (9.56)$$

where we have assumed a pairwise CRF for notational simplicity. If we perturb the local evidence potentials  $\phi_i(k)$  by adding  $\epsilon_{ik}$  to each entry, where  $k$  indexes the discrete latent states, we get  $\tilde{S}(\mathbf{x}; \mathbf{c})$ . We then compute a sample  $\tilde{\mathbf{x}}$  by solving  $\tilde{\mathbf{x}} = \text{argmax } \tilde{S}(\mathbf{x}; \mathbf{c})$ . The advantage of this approach is that it can leverage efficient MAP solvers for discrete models, such as those discussed in Supplementary Section 9.3. This can in turn be used for parameter learning, and estimating the partition function [HJ12; Erm+13].

# 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; Zha+19b; Bro18].

### 10.1.1 Variational free energy

Consider a model with unknown variables  $\mathbf{z}$ , known variables  $\mathbf{x}$ , and fixed parameters  $\boldsymbol{\theta}$ . (If the parameters are unknown, they can be added to  $\mathbf{z}$ .) Since computing the true posterior  $p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$  is assumed intractable, we will use an approximation  $q(\mathbf{z})$ , which we choose to 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  $\boldsymbol{\psi}$  to represent the **variational parameters**. We compute the best variational parameters (for given  $\mathbf{x}$ ) as follows:

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

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

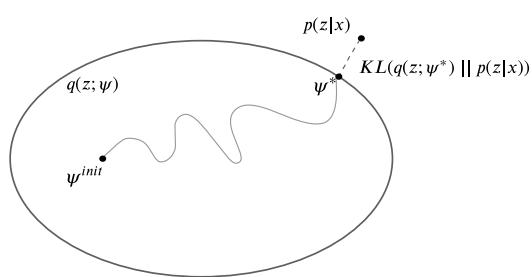
$$= \underset{\boldsymbol{\psi}}{\operatorname{argmin}} \underbrace{\mathbb{E}_{q(\mathbf{z}|\boldsymbol{\psi})} [\log q(\mathbf{z}|\boldsymbol{\psi}) - \log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z}) - \log p_{\boldsymbol{\theta}}(\mathbf{z})]}_{\mathcal{L}(\boldsymbol{\psi}|\boldsymbol{\theta},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

The final term  $\log p_{\boldsymbol{\theta}}(\mathbf{x}) = \int p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})d\mathbf{z}$  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 using the following shorthand:<sup>1</sup>

$$\mathcal{L}(\boldsymbol{\psi}|\boldsymbol{\theta}, \mathbf{x}) = D_{\text{KL}}(q(\mathbf{z}|\boldsymbol{\psi}) \parallel p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})) \triangleq \mathbb{E}_{q(\mathbf{z}|\boldsymbol{\psi})} [\log q(\mathbf{z}|\boldsymbol{\psi}) - \log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})] \quad (10.5)$$

---

1. Technically speaking  $D_{\text{KL}}(q(\mathbf{z}|\boldsymbol{\psi}) \parallel p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}))$  is not a KL divergence, since these are distributions over different spaces: the first is a conditional distribution over  $\mathbf{z}$ , the second is a joint distribution over  $\mathbf{z}$  and  $\mathbf{x}$ . However, hopefully this shorthand is clear.



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

1617

18 If we define  $\mathcal{E}(z) = -\log p_{\theta}(z, x)$  as the energy, then we can write

19

$$\mathcal{L}(\psi|\theta, x) = \mathbb{E}_{q(z|\psi)} [\mathcal{E}(z)] - \mathbb{H}(q) \quad (10.6)$$

21

22 where  $\mathbb{H}(q)$  is the entropy. In physics, this is known as the **variational free energy**. We can  
23 interpret this as the expected energy minus the entropy:

24

$$\text{VFE} = \text{expected energy} - \text{entropy} \quad (10.7)$$

25

26 This is an upper bound on the **free energy**,  $-\log p_{\theta}(x)$ , which follows from the fact that

27

$$D_{\text{KL}}(q \parallel p) = \text{VFE}(\psi|\theta, x) + \log p_{\theta}(x) \geq 0 \quad (10.8)$$

29

30 Our goal is to minimize the VFE. See Figure 10.1 for an illustration.

31

### 32 10.1.2 Evidence lower bound (ELBO)

33 The negative of the VFE is known as the **evidence lower bound** or **ELBO** function [BKM16]:  
34

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

36

37 where  $\log p_{\theta}(x, z)$  is the unnormalized log joint. The name “ELBO” arises because

38

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

39

40 where  $\log p_{\theta}(x)$  is called the “evidence”. The inequality follows from Equation (10.8). Therefore  
41 maximizing the ELBO wrt  $\psi$  will decrease the original KL, since  $\log p_{\theta}(x)$  is a constant wrt  $\psi$ .  
42 (Note: we use the symbol  $\mathcal{L}$  for the ELBO, rather than  $\mathcal{L}$ , since the latter denotes a loss we want to  
43 minimize.)

44

45 We can rewrite the ELBO as follows:

46

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

47

1 We can interpret this  
2

3 ELBO = expected log joint + entropy of posterior (10.12)  
4

5 The second term encourages the posterior to be maximum entropy, while the first term encourages it  
6 to be a joint MAP configuration,  
7

We can also rewrite the ELBO in the following equivalent way:  
8

9  $\mathcal{L}(\psi|\theta, \mathbf{x}) = \mathbb{E}_{q(z|\psi)} [\log p_\theta(\mathbf{x}|z) + \log p_\theta(z) - \log q(z|\psi)]$  (10.13)  
10

11  $= \mathbb{E}_{q(z|\psi)} [\log p_\theta(\mathbf{x}|z)] - D_{\text{KL}}(q(z|\psi) \parallel p_\theta(z))$  (10.14)  
12

We can interpret this as follows:  
13

14 ELBO = expected log likelihood – KL from posterior to prior (10.15)  
15

15 The KL term acts like a regularizer, preventing the posterior from diverging too much from the prior.  
16 (See also Section 21.2.2, where we discuss the ELBO in more detail, in the context of variational  
17 autoencoders.)  
18

## 19 10.2 Mean field VI 20

21 A common approximation in variational inference is to assume that all the latent variables are  
22 independent, i.e.,  
23

24  $q(z|\psi) = \prod_{j=1}^J q_j(z_j)$  (10.16)  
25

27 where  $J$  is the number of hidden variables, and  $q_j(z_j)$  is shorthand for  $q_{\psi_j}(z_j)$ , where  $\psi_j$  are the  
28 variational parameters for the  $j$ 'th distribution. This is called the **mean field** approximation.  
29

From Equation (10.11), the ELBO becomes  
30

31  $\mathcal{L}(\psi) = \int q(z|\psi) \log p_\theta(\mathbf{x}, z) dz + \sum_{j=1}^J \mathbb{H}(q_j)$  (10.17)  
32

34 since the entropy of a product distribution is the sum of entropies of each component in the product.  
35 We can either directly optimize this (see e.g., [Baq+16]), or use a coordinate-wise optimization  
36 scheme, as we discuss in Section 10.2.1.  
37

### 38 10.2.1 Coordinate ascent variational inference (CAVI) 39

40 We now discuss a coordinate ascent method for optimizing the mean field objective, which we call  
41 **coordinate ascent variational inference** or **CAVI**.  
42

To derive the update equations, we initially assume there are just 3 discrete latent variables, to  
43 simplify notation. In this case the ELBO is given by  
44

45  $\mathcal{L}(q_1, q_2, q_3) = \sum_{z_1} \sum_{z_2} \sum_{z_3} q_1(z_1) q_2(z_2) q_3(z_3) \log \tilde{p}(z_1, z_2, z_3) + \sum_{j=1}^3 \mathbb{H}(q_j)$  (10.18)  
46

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  $q_3$ :

$$\mathcal{L}_3(q_3) = \sum_{z_3} q_3(z_3) \left[ \sum_{z_1} \sum_{z_2} q_1(z_1) q_2(z_2) \log \tilde{p}(z_1, z_2, z_3) \right] + \mathbb{H}(q_3) + \text{const} \quad (10.19)$$

$$= \sum_{z_3} q_3(z_3) [g_3(z_3) - \log q_3(z_3)] + \text{const} \quad (10.20)$$

where

$$g_3(z_3) \triangleq \sum_{z_1} \sum_{z_2} q_1(z_1) q_2(z_2) \log \tilde{p}(z_1, z_2, z_3) = \mathbb{E}_{\mathbf{z}_{-3}} [\log \tilde{p}(z_1, z_2, z_3)] \quad (10.21)$$

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

$$\tilde{f}_3(z_3) = \exp(g_3(z_3)) \quad (10.22)$$

which we can normalize to get

$$f_3(z_3) = \frac{\tilde{f}_3(z_3)}{\sum_{z'_3} \tilde{f}_3(z'_3)} \propto \exp(g_3(z_3)) \quad (10.23)$$

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

$$\mathcal{L}_3(q_3) = \sum_{z_3} q_3(z_3) [\log f_3(z_3) - \log q_3(z_3)] + \text{const} = -D_{\text{KL}}(q_3 \parallel f_3) + \text{const} \quad (10.24)$$

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

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

$$\mathcal{L}_3(q_3) = \sum_{z_3} q_3(z_3) \left[ \sum_{z_2} q_2(z_2) \log \tilde{p}(z_2, z_3) \right] + \mathbb{H}(q_3) + \text{const} \quad (10.25)$$

$$= \sum_{z_3} q_3(z_3) [\log f_3(z_3) - \log q_3(z_3)] + \text{const} \quad (10.26)$$

where

$$f_3(z_3) \propto \exp \left[ \sum_{z_2} q_2(z_2) \log \tilde{p}(z_2, z_3) \right] = \exp \left[ \mathbb{E}_{\mathbf{z}_{\text{mb}_3}} [\log \tilde{p}(z_2, z_3)] \right] \quad (10.27)$$

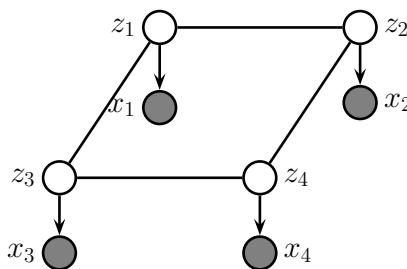


Figure 10.2: A grid-structured MRF with hidden nodes  $z_i$  and local evidence nodes  $x_i$ . The prior  $p(\mathbf{z})$  is an undirected Ising model, and the likelihood  $p(\mathbf{x}|\mathbf{z}) = \prod_i p(x_i|z_i)$  is a directed fully factored model.

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

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

$$q_j(\mathbf{z}_j) \propto \exp \left[ \mathbb{E}_{\mathbf{z}_{\text{mb}_j}} [\log \tilde{p}(\mathbf{z}_j, \mathbf{z}_{\text{mb}_j})] \right] \quad (10.28)$$

(Compare to the equation for Gibbs sampling in Equation (12.19).) The CAVI method simply computes  $q_j$  for each dimension  $j$  in turn, in an iterative fashion (see Algorithm 19). Convergence is guaranteed since the bound is convex wrt each of the factors  $q_i$  [Bis06, p. 466].

---

**Algorithm 19:** Coordinate Ascent Variational Inference (CAVI).

---

```

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

```

---

Note that the functional form of the  $q_i$  distributions does not need to be specified in advance, but will be determined by the form of the log joint. This is therefore called **free-form VI**, as opposed to fixed-form, where we explicitly choose a convenient distributional type for  $q$  (we discuss fixed-form VI in Section 10.3). We give some examples below that will make this clearer.

### 10.2.2 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  $i$ ,

1 and  $x_i \in \mathbb{R}$  be the observed noisy value. See Figure 10.2 for the graphical model.  
2

3 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  $i$   
4 in the graphical model). The overall likelihood has the form

$$\underline{5} \quad p(\mathbf{x}|\mathbf{z}) = \prod_i p(x_i|z_i) = \exp\left(\sum_i L_i(z_i)\right) \quad (10.29)$$

6 Our goal is to approximate the posterior  $p(\mathbf{z}|\mathbf{x})$ . We will use an Ising model for the prior:  
7

$$\underline{10} \quad p(\mathbf{z}) = \frac{1}{Z_0} \exp(-\mathcal{E}_0(\mathbf{z})) \quad (10.30)$$

$$\underline{12} \quad \mathcal{E}_0(\mathbf{z}) = -\sum_{i \sim j} W_{ij} z_i z_j \quad (10.31)$$

15 where we sum over each  $i - j$  edge. Therefore the posterior has the form  
16

$$\underline{17} \quad p(\mathbf{z}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp(-\mathcal{E}(\mathbf{z})) \quad (10.32)$$

$$\underline{19} \quad \mathcal{E}(\mathbf{z}) = \mathcal{E}_0(\mathbf{z}) - \sum_i L_i(z_i) \quad (10.33)$$

21 We will now make the following fully factored approximation:  
22

$$\underline{23} \quad q(\mathbf{z}) = \prod_i q_i(z_i) = \prod_i \text{Ber}(z_i|\mu_i) \quad (10.34)$$

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

$$\underline{30} \quad \log \tilde{p}(\mathbf{z}) = z_i \sum_{j \in \text{nbr}_i} W_{ij} z_j + L_i(z_i) + \text{const} \quad (10.35)$$

32 This only depends on the states of the neighboring nodes. Hence  
33

$$\underline{35} \quad q_i(z_i) \propto \exp(\mathbb{E}_{q_{-i}(\mathbf{z})} [\log \tilde{p}(\mathbf{z})]) = \exp\left(z_i \sum_{j \in \text{nbr}_i} W_{ij} \mu_j + L_i(z_i)\right) \quad (10.36)$$

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

40 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  $i$ .

41 Also, let  $L_i^+ \triangleq L_i(+1)$  and  $L_i^- \triangleq L_i(-1)$ . The approximate marginal posterior is given by  
42

$$\underline{43} \quad q_i(z_i = 1) = \frac{e^{m_i + L_i^+}}{e^{m_i + L_i^+} + e^{-m_i + L_i^-}} = \frac{1}{1 + e^{-2m_i + L_i^- - L_i^+}} = \sigma(2a_i) \quad (10.37)$$

$$\underline{45} \quad a_i \triangleq m_i + 0.5(L_i^+ - L_i^-) \quad (10.38)$$

47

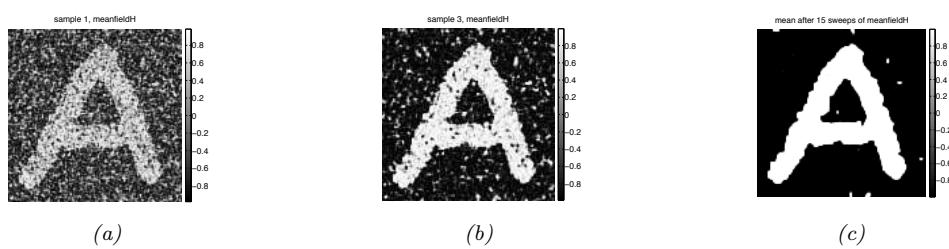


Figure 10.3: Example of image denoising using mean field (with parallel updates and a damping factor of 0.5). We use an Ising prior with  $W_{ij} = 1$  and a Gaussian noise model with  $\sigma = 2$ . We show the results after 1, 3 and 15 iterations across the image. Compare to Figure 12.3, which shows the results of using Gibbs sampling. Generated by [ising\\_image\\_denoise\\_demo.ipynb](#).

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

$$\mu_i = \mathbb{E}_{q_i}[z_i] = q_i(z_i = +1) \cdot (+1) + q_i(z_i = -1) \cdot (-1) \quad (10.39)$$

$$= \frac{1}{1 + e^{-2a_i}} - \frac{1}{1 + e^{2a_i}} = \frac{e^{a_i}}{e^{a_i} + e^{-a_i}} - \frac{e^{-a_i}}{e^{-a_i} + e^{a_i}} = \tanh(a_i) \quad (10.40)$$

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

$$\mu_i^t = \tanh \left( \sum_{j \in \text{nbr}_i} W_{ij} \mu_j^{t-1} + 0.5(L_i^+ - L_i^-) \right) \quad (10.41)$$

Following [MJW99], we can use **damped updates** of the following form to improve convergence:

$$\mu_i^t = (1 - \lambda) \mu_i^{t-1} + \lambda \tanh \left( \sum_{j \in \text{nbr}_i} W_{ij} \mu_j^{t-1} + 0.5(L_i^+ - L_i^-) \right) \quad (10.42)$$

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

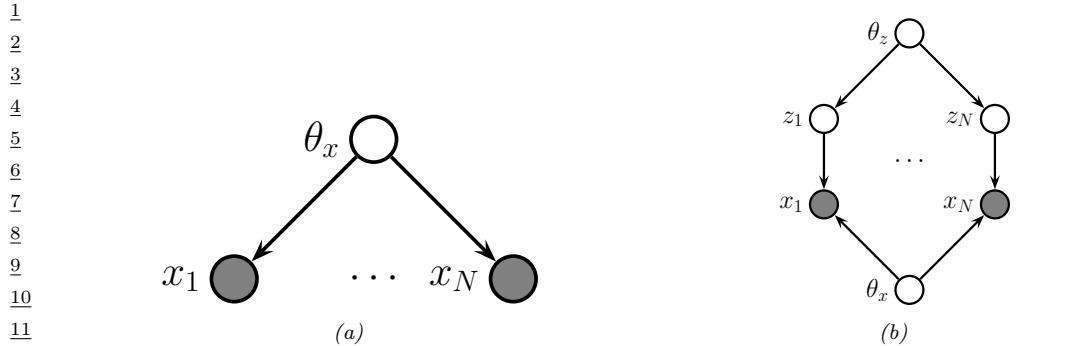
Figure 10.3 shows the method in action, applied to a 2d Ising model with homogeneous attractive potentials,  $W_{ij} = 1$ . We use parallel updates with a damping factor of  $\lambda = 0.5$ . (If we don't use damping, we tend to get "checkerboard" artefacts.)

### 10.2.3 Variational Bayes

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

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

$$p(\theta, \mathcal{D}) = p(\theta) \prod_{n=1}^N p(\mathcal{D}_n | \theta) \quad (10.43)$$



*Figure 10.4: Graphical models with (a) Global hidden variable  $\theta_x$  and observed variables  $\mathbf{x}_{1:N}$ . (b) Local hidden variables  $\mathbf{z}_{1:N}$ , global hidden variables  $\theta_x, \theta_z$ , and observed variables  $\mathbf{x}_{1:N}$ .*

These conditional independencies are illustrated in Figure 10.4a.

We will fit the variational posterior by maximizing the ELBO

$$\mathcal{L}(\psi_{\theta} | \mathcal{D}) = \mathbb{E}_{q(\theta | \psi_{\theta})} [\log p(\theta, \mathcal{D})] + \mathbb{H}(q(\theta | \psi_{\theta})) \quad (10.44)$$

We will assume the variational posterior factorizes over the parameters:

$$q(\theta | \psi_{\theta}) = \prod_j q(\theta_j | \psi_{\theta_j}) \quad (10.45)$$

We can then update each  $\psi_{\theta_j}$  using CAVI (Section 10.2.1).

### 10.2.4 Example: VB for a univariate Gaussian

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

$$p(\mu, \lambda) = \mathcal{N}(\mu | \mu_0, (\kappa_0 \lambda)^{-1}) \text{Ga}(\lambda | a_0, b_0) \quad (10.46)$$

It is possible to derive the posterior  $p(\mu, \lambda | \mathcal{D})$  for this model exactly, as shown in Section 3.2.3.3. However, here we use the VB method with the following factored approximate posterior:

$$q(\mu, \lambda) = q(\mu | \psi_{\mu})q(\lambda | \psi_{\lambda}) \quad (10.47)$$

We do not need to specify the forms for the distributions  $q(\mu | \psi_{\mu})$  and  $q(\lambda | \psi_{\lambda})$ ; 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  
2 **10.2.4.1 Target distribution**

3 The unnormalized log posterior has the form  
4

5  $\log \tilde{p}(\mu, \lambda) = \log p(\mu, \lambda, \mathcal{D}) = \log p(\mathcal{D}|\mu, \lambda) + \log p(\mu|\lambda) + \log p(\lambda)$  (10.48)  
6

7  $= \frac{N_{\mathcal{D}}}{2} \log \lambda - \frac{\lambda}{2} \sum_{n=1}^{N_{\mathcal{D}}} (x_n - \mu)^2 - \frac{\kappa_0 \lambda}{2} (\mu - \mu_0)^2$   
8  
9  $+ \frac{1}{2} \log(\kappa_0 \lambda) + (a_0 - 1) \log \lambda - b_0 \lambda + \text{const}$  (10.49)  
10  
11

12  
13 **10.2.4.2 Updating  $q(\mu|\psi_\mu)$**

14 The optimal form for  $q(\mu|\psi_\mu)$  is obtained by averaging over  $\lambda$ :  
15

16  $\log q(\mu|\psi_\mu) = \mathbb{E}_{q(\lambda|\psi_\lambda)} [\log p(\mathcal{D}|\mu, \lambda) + \log p(\mu|\lambda)] + \text{const}$  (10.50)  
17

18  $= -\frac{\mathbb{E}_{q(\lambda|\psi_\lambda)} [\lambda]}{2} \left\{ \kappa_0(\mu - \mu_0)^2 + \sum_{n=1}^{N_{\mathcal{D}}} (x_n - \mu)^2 \right\} + \text{const}$  (10.51)  
19  
20

21 By completing the square one can show that  $q(\mu|\psi_\mu) = \mathcal{N}(\mu|\mu_N, \kappa_N^{-1})$ , where  
22

23  $\mu_N = \frac{\kappa_0 \mu_0 + N_{\mathcal{D}} \bar{x}}{\kappa_0 + N_{\mathcal{D}}}, \quad \kappa_N = (\kappa_0 + N_{\mathcal{D}}) \mathbb{E}_{q(\lambda|\psi_\lambda)} [\lambda]$  (10.52)  
24  
25

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.  
28

29  
30 **10.2.4.3 Updating  $q(\lambda|\psi_\lambda)$**

31 The optimal form for  $q(\lambda|\psi_\lambda)$  is given by  
32

33  $\log q(\lambda|\psi_\lambda) = \mathbb{E}_{q(\mu|\psi_\mu)} [\log p(\mathcal{D}|\mu, \lambda) + \log p(\mu|\lambda) + \log p(\lambda)] + \text{const}$  (10.53)  
34

35  $= (a_0 - 1) \log \lambda - b_0 \lambda + \frac{1}{2} \log \lambda + \frac{N_{\mathcal{D}}}{2} \log \lambda$   
36  
37  $- \frac{\lambda}{2} \mathbb{E}_{q(\mu|\psi_\mu)} \left[ \kappa_0(\mu - \mu_0)^2 + \sum_{n=1}^{N_{\mathcal{D}}} (x_n - \mu)^2 \right] + \text{const}$  (10.54)  
38  
39

40 We recognize this as the log of a Gamma distribution, hence  $q(\lambda|\psi_\lambda) = \text{Ga}(\lambda|a_N, b_N)$ , where  
41

42  $a_N = a_0 + \frac{N+1}{2}$  (10.55)  
43

44  $b_N = b_0 + \frac{1}{2} \mathbb{E}_{q(\mu|\psi_\mu)} \left[ \kappa_0(\mu - \mu_0)^2 + \sum_{n=1}^{N_{\mathcal{D}}} (x_n - \mu)^2 \right]$  (10.56)  
45  
46  
47

1 **10.2.4.4 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

$$\underline{6} \quad \mathbb{E}_{q(\mu)} [\mu] = \mu_N \quad (10.57)$$

$$\underline{7} \quad \mathbb{E}_{q(\mu)} [\mu^2] = \frac{1}{\kappa_N} + \mu_N^2 \quad (10.58)$$

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

$$\underline{11} \quad \mathbb{E}_{q(\lambda)} [\lambda] = \frac{a_N}{b_N} \quad (10.59)$$

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

$$\underline{16} \quad \mu_N = \frac{\kappa_0 \mu_0 + N_{\mathcal{D}} \bar{x}}{\kappa_0 + N_{\mathcal{D}}} \quad (10.60)$$

$$\underline{18} \quad \kappa_N = (\kappa_0 + N_{\mathcal{D}}) \frac{a_N}{b_N} \quad (10.61)$$

20 and for  $q(\lambda)$  we have  
21

$$\underline{22} \quad a_N = a_0 + \frac{N+1}{2} \quad (10.62)$$

$$\underline{24} \quad b_N = b_0 + \frac{1}{2} \kappa_0 (\mathbb{E} [\mu^2] + \mu_0^2 - 2\mathbb{E} [\mu] \mu_0) + \frac{1}{2} \sum_{n=1}^{N_{\mathcal{D}}} (x_n^2 + \mathbb{E} [\mu^2] - 2\mathbb{E} [\mu] x_n) \quad (10.63)$$

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

31 **10.2.4.5 Illustration**

33 Figure 10.5 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 5.1.3.3.  
38

39 **10.2.4.6 Lower bound**

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

$$\underline{42} \quad \underline{43} \quad \mathcal{L}(\psi_{\theta} | \mathcal{D}) \leq \log p(\mathcal{D}) = \log \iint p(\mathcal{D} | \mu, \lambda) p(\mu, \lambda) d\mu d\lambda \quad (10.64)$$

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

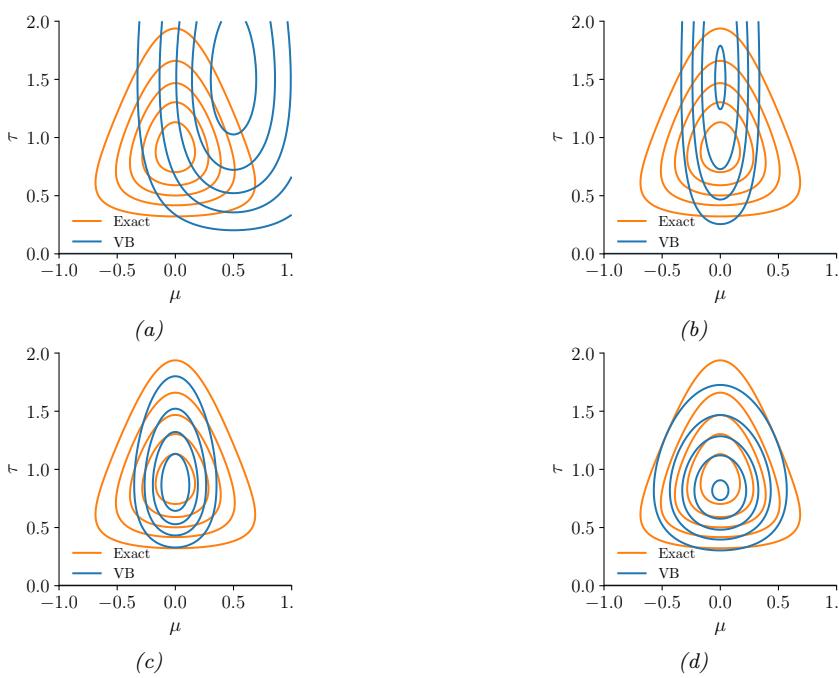


Figure 10.5: Factored variational approximation (red) to the Gaussian-Gamma distribution (green). (a) Initial guess. (b) After updating  $q(\mu|\psi_\mu)$ . (c) After updating  $q(\lambda|\psi_\lambda)$ . (d) At convergence (after 5 iterations). Adapted from Fig. 10.4 of [Bis06]. Generated by [unigauss\\_vb\\_demo.ipynb](#).

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. One can show that the lower bound has the following form:

$$\bar{L} = \text{const} + \frac{1}{2} \ln \frac{1}{\kappa_N} + \ln \Gamma(a_N) - a_N \ln b_N \quad (10.65)$$

### 10.2.5 Variational Bayes EM

In Bayesian latent variable models, we have two forms of hidden variables: local (or per example) hidden variables  $\mathbf{z}_n$ , and global (shared) hidden variables  $\boldsymbol{\theta}$ , which represent the parameters of the model. See Figure 10.4b 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  $\mathbf{h} = (\boldsymbol{\theta}, \mathbf{z}_{1:N})$  represents all the hidden variables, then the joint distribution is given by

$$p(\mathbf{h}, \mathcal{D}) = p(\boldsymbol{\theta}, \mathbf{z}_{1:N}, \mathcal{D}) = p(\boldsymbol{\theta}) \prod_{n=1}^N p(\mathbf{z}_n|\boldsymbol{\theta})p(\mathbf{x}_n|\mathbf{z}_n, \boldsymbol{\theta}) \quad (10.66)$$

1 We will make the following mean field assumption:  
2

3

$$\underline{q}(\boldsymbol{\theta}, \mathbf{z}_{1:N} | \psi_{1:N}, \psi_{\boldsymbol{\theta}}) = q(\boldsymbol{\theta} | \psi_{\boldsymbol{\theta}}) \prod_{n=1}^N q(\mathbf{z}_n | \psi_n) \quad (10.67)$$

4

5 Henceforth we will usually omit the variational parameters  $(\psi_{1:N}, \psi_{\boldsymbol{\theta}})$  for brevity.  
6

We will use VI to maximize the ELBO:

7

$$\underline{L}(\psi_{1:N}, \psi_{\boldsymbol{\theta}} | \mathcal{D}) = \mathbb{E}_{q(\boldsymbol{\theta}, \mathbf{z}_{1:N} | \psi_{1:N}, \psi_{\boldsymbol{\theta}})} [\log p(\mathbf{z}_{1:N}, \boldsymbol{\theta}, \mathcal{D}) - \log q(\boldsymbol{\theta}, \mathbf{z}_{1:N})] \quad (10.68)$$

8

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

17 Now suppose we approximate  $q(\boldsymbol{\theta})$  by a delta function,  $q(\boldsymbol{\theta}) = \delta(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})$ . The Bayesian LVM ELBO  
18 objective from Equation (10.68) simplifies to the the “LVM ELBO”:  
19

20

$$\underline{L}(\psi_{1:N}, \boldsymbol{\theta} | \mathcal{D}) = \mathbb{E}_{q(\mathbf{z}_{1:N} | \psi_{1:N})} [\log p(\boldsymbol{\theta}, \mathcal{D}, \mathbf{z}_{1:N}) - \log q(\mathbf{z}_{1:N} | \psi_{1:N})] \quad (10.69)$$

21

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

24 VEM is simpler than VBEM since in the the variational E step, we compute  $q(\mathbf{z}_n | \mathbf{x}_n, \hat{\boldsymbol{\theta}})$ , instead  
25 of  $\mathbb{E}_{\boldsymbol{\theta}}[q(\mathbf{z}_n | \mathbf{x}_n, \boldsymbol{\theta})]$ ; that is, we plugin a point estimate of the model parameters, rather than averaging  
26 over the parameters. For more details on VEM, see Section 6.6.6.1.  
27

### 28 10.2.6 Example: VBEM for a GMM

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

31

$$\underline{p}(\mathbf{z}, \mathbf{x} | \boldsymbol{\theta}) = \prod_n \prod_k \pi_k^{z_{nk}} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1})^{z_{nk}} \quad (10.70)$$

32

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

36

$$\underline{p}(\boldsymbol{\theta}) = \text{Dir}(\boldsymbol{\pi} | \check{\boldsymbol{\alpha}}) \prod_k \mathcal{N}(\boldsymbol{\mu}_k | \check{\mathbf{m}}, (\check{\boldsymbol{\Lambda}}_k)^{-1}) \text{Wi}(\boldsymbol{\Lambda}_k | \check{\mathbf{L}}, \check{\nu}) \quad (10.71)$$

37

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

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 10.3.5.3, where we discuss a different variational approximation for this model.)  
45

1  
2 **10.2.6.1 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  $q_n(z_n) = \text{Cat}(z_n | \mathbf{r}_n)$  (10.72)  
8

9  $q(\boldsymbol{\theta}) = \text{Dir}(\boldsymbol{\pi} | \hat{\boldsymbol{\alpha}}) \prod_k \mathcal{N}(\boldsymbol{\mu}_k | \hat{\mathbf{m}}_k, (\hat{\kappa}_k \boldsymbol{\Lambda}_k)^{-1}) \text{Wi}(\boldsymbol{\Lambda}_k | \hat{\mathbf{L}}_k, \hat{\nu}_k)$  (10.73)  
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 **10.2.6.2 Derivation of  $q(\boldsymbol{\theta})$  (variational M step)**  
16

17 Using the mean field recipe in Algorithm 19, we write down the log joint, and take expectations over  
18 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 
$$\log q(\boldsymbol{\theta}) = \log p(\boldsymbol{\pi}) + \underbrace{\sum_n \mathbb{E}_{q(z_n)} [\log p(\mathbf{z}_n | \boldsymbol{\pi})]}_{L_{\boldsymbol{\pi}}}$$
  
20  
21  
22  
23  
24  
25  $+ \sum_k \left[ \underbrace{\log p(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k) \sum_n \mathbb{E}_{q(z_n)} [z_{nk}] \log \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1})}_{L_{\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k}} \right] + \text{const}$  (10.74)  
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  $q(\boldsymbol{\theta}) = q(\boldsymbol{\pi}) \prod_k q(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)$  (10.75)  
33  
34

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

37  $\log q(\boldsymbol{\pi}) = (\alpha_0 - 1) \sum_k \log \pi_k + \sum_k \sum_n r_{nk} \log \pi_k + \text{const}$  (10.76)  
38  
39

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

42  $q(\boldsymbol{\pi}) = \text{Dir}(\boldsymbol{\pi} | \hat{\boldsymbol{\alpha}})$  (10.77)  
43

44  $\hat{\alpha}_k = \alpha_0 + N_k$  (10.78)

45  $N_k = \sum_n r_{nk}$  (10.79)  
46

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

$$q(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k) = \mathcal{N}(\boldsymbol{\mu}_k | \widehat{\boldsymbol{m}}_k, (\widehat{\kappa}_k \boldsymbol{\Lambda}_k)^{-1}) \text{Wi}(\boldsymbol{\Lambda}_k | \widehat{\mathbf{L}}_k, \widehat{\nu}_k) \quad (10.80)$$

$$\widehat{\kappa}_k = \check{\kappa} + N_k \quad (10.81)$$

$$\widehat{\boldsymbol{m}}_k = (\check{\kappa} \check{\boldsymbol{m}} + N_k \bar{\boldsymbol{x}}_k) / \widehat{\kappa}_k \quad (10.82)$$

$$\widehat{\mathbf{L}}_k^{-1} = \check{\mathbf{L}}^{-1} + N_k \mathbf{S}_k + \frac{\check{\kappa} N_k}{\check{\kappa} + N_k} (\bar{\boldsymbol{x}}_k - \check{\boldsymbol{m}})(\bar{\boldsymbol{x}}_k - \check{\boldsymbol{m}})^\top \quad (10.83)$$

$$\widehat{\nu}_k = \check{\nu} + N_k \quad (10.84)$$

$$\bar{\boldsymbol{x}}_k = \frac{1}{N_k} \sum_n r_{nk} \boldsymbol{x}_n \quad (10.85)$$

$$\mathbf{S}_k = \frac{1}{N_k} \sum_n r_{nk} (\boldsymbol{x}_n - \bar{\boldsymbol{x}}_k)(\boldsymbol{x}_n - \bar{\boldsymbol{x}}_k)^\top \quad (10.86)$$

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  $\hat{\boldsymbol{\theta}}$ .

### 10.2.6.3 Derivation of $q(z)$ (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

$$\begin{aligned} \log q(\boldsymbol{z}) &= \sum_n \sum_k z_{nk} \left( \mathbb{E}_{q(\boldsymbol{\pi})} [\log \pi_k] + \frac{1}{2} \mathbb{E}_{q(\boldsymbol{\Lambda}_k)} [\log |\boldsymbol{\Lambda}_k|] - \frac{D}{2} \log(2\pi) \right. \\ &\quad \left. - \frac{1}{2} \mathbb{E}_{q(\boldsymbol{\theta})} [(\boldsymbol{x}_n - \boldsymbol{\mu}_k)^\top \boldsymbol{\Lambda}_k (\boldsymbol{x}_n - \boldsymbol{\mu}_k)] \right) + \text{const} \end{aligned} \quad (10.87)$$

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

$$\exp(\mathbb{E}_{q(\boldsymbol{\pi})} [\log \pi_k]) = \frac{\exp(\psi(\widehat{\alpha}_k))}{\exp(\psi(\sum_{k'} \widehat{\alpha}_{k'}))} \triangleq \tilde{\pi}_k \quad (10.88)$$

where  $\psi$  is the **digamma function**

$$\psi(x) = \frac{d}{dx} \log \Gamma(x) \quad (10.89)$$

This takes care of the first term.

For the second term, one can show

$$\mathbb{E}_{q(\boldsymbol{\Lambda}_k)} [\log |\boldsymbol{\Lambda}_k|] = \sum_{j=1}^D \psi\left(\frac{\widehat{\nu}_k + 1 - j}{2}\right) + D \log 2 + \log |\widehat{\mathbf{L}}_k| \quad (10.90)$$

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

$$\mathbb{E}_{q(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)} [(\boldsymbol{x}_n - \boldsymbol{\mu}_k)^\top \boldsymbol{\Lambda}_k (\boldsymbol{x}_n - \boldsymbol{\mu}_k)] = D \widehat{\kappa}_k^{-1} + \widehat{\nu}_k (\boldsymbol{x}_n - \widehat{\boldsymbol{m}}_k)^\top \widehat{\mathbf{L}}_k (\boldsymbol{x}_n - \widehat{\boldsymbol{m}}_k) \triangleq \tilde{\Lambda}_k \quad (10.91)$$

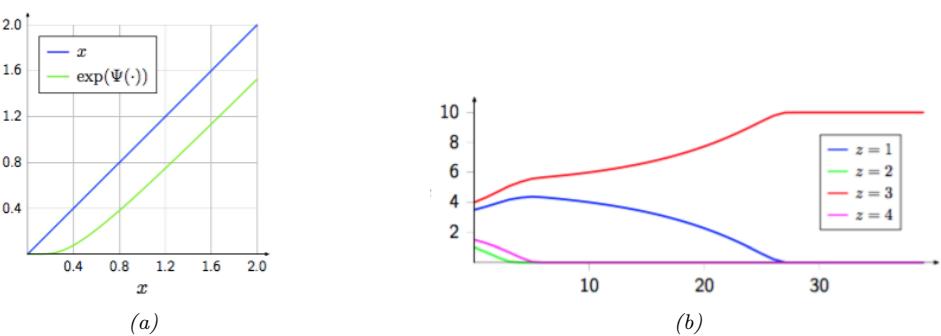


Figure 10.6: (a) We plot  $\exp(\psi(x))$  vs  $x$ . We see that this function performs a form of shrinkage, so that small values get set to zero. (b) We plot  $N_k$  vs time for 4 different states ( $z$  values), starting from random initial values. We perform a series of VBEM updates, ignoring the 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

$$r_{nk} \propto \tilde{\pi}_k \tilde{\Lambda}_k^{\frac{1}{2}} \exp\left(-\frac{D}{2 \tilde{\kappa}_k} - \frac{\hat{\nu}_k}{2} (\mathbf{x}_n - \hat{\boldsymbol{\mu}}_k)^T \hat{\boldsymbol{\Lambda}}_k (\mathbf{x}_n - \hat{\boldsymbol{\mu}}_k)\right) \quad (10.92)$$

Compare this to the expression used in regular EM:

$$r_{nk}^{EM} \propto \hat{\pi}_k |\hat{\boldsymbol{\Lambda}}_k|^{\frac{1}{2}} \exp\left(-\frac{1}{2} (\mathbf{x}_n - \hat{\boldsymbol{\mu}}_k)^T \hat{\boldsymbol{\Lambda}}_k (\mathbf{x}_n - \hat{\boldsymbol{\mu}}_k)\right) \quad (10.93)$$

where  $\hat{\pi}_k$  is the MAP estimate for  $\pi_k$ . The significance of this difference is discussed in Section 10.2.6.4.

#### 10.2.6.4 Automatic sparsity inducing effects of VBEM

In regular EM, the E step has the form given in Equation (10.93), whereas in VBEM, the E step has the form given in Equation (10.92). 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 (10.88) we have

$$r_{nk}^{VB} = \tilde{\pi}_k = \frac{\exp(\psi(\hat{\alpha}_k))}{\exp(\psi(\sum_{k'} \hat{\alpha}_{k'}))} \quad (10.94)$$

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

$$r_{nk}^{EM} = \hat{\pi}_k = \frac{\hat{\alpha}_k - 1}{\sum_{k'} (\hat{\alpha}_{k'} - 1)} \quad (10.95)$$

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 2.8 that using  $\alpha_0 \ll 1$  causes  $\pi$  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 data points 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:

$$\tilde{\pi}_k = \frac{\exp(\psi(N_k))}{\exp(\psi(\sum_{k'} N_{k'}))} \quad (10.96)$$

$$\hat{\pi}_k = \frac{N_k - 1}{\sum_{k'} (N_{k'} - 1)} \quad (10.97)$$

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 10.6(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 10.6(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 Section 31.3, when we discuss Dirichlet process mixture models.

The reason for this effect is shown in Figure 10.6(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  $\ell_1$ -regularization (see Section 15.2.5). 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 10.7. In Figure 10.8, we plot the initial and final values of  $\alpha_k$ ; we see that  $\hat{\alpha}_k = 0$  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 10.2.6.5 Lower bound on the marginal likelihood 39

40 The VBEM algorithm is maximizing the following lower bound  
41

$$42 \quad \mathcal{L} = \sum_{\mathbf{z}} \int d\boldsymbol{\theta} q(\mathbf{z}, \boldsymbol{\theta}) \log \frac{p(\mathbf{x}, \mathbf{z}, \boldsymbol{\theta})}{q(\mathbf{z}, \boldsymbol{\theta})} \leq \log p(\mathbf{x}) \quad (10.98)$$

45 This quantity increases monotonically with each iteration, as shown in Figure 10.9.  
46

47

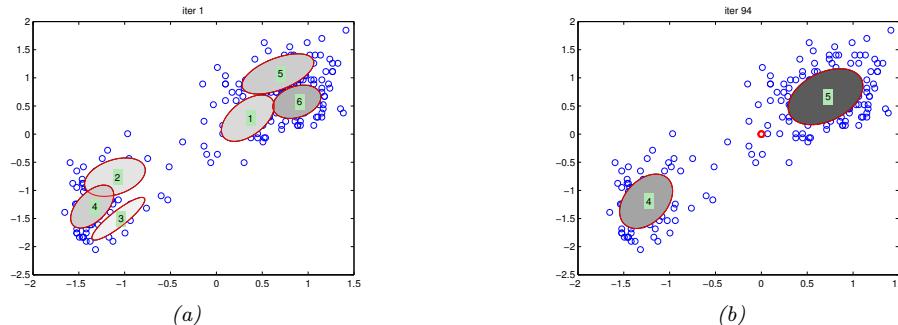


Figure 10.7: We visualize the posterior mean parameters at various stages of the VBEM algorithm applied to a mixture of Gaussians model on the Old Faithful data. Shading intensity is proportional to the mixing weight. We initialize with K-means and use  $\alpha_0 = 0.001$  as the Dirichlet hyper-parameter. (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](#).

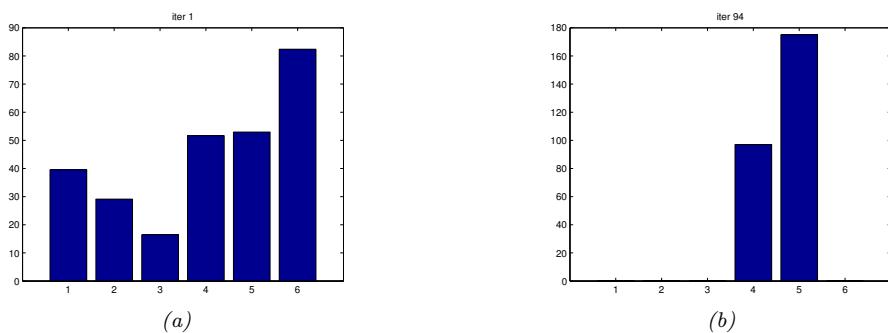


Figure 10.8: We visualize the posterior values of  $\alpha_k$  for the model in Figure 10.7 after the first and last iteration of the algorithm. We see that unnecessary components get “killed off”. (Interestingly, the initially large cluster 6 gets “replaced” by cluster 5.) Generated by [gmm\\_vb\\_em.ipynb](#).

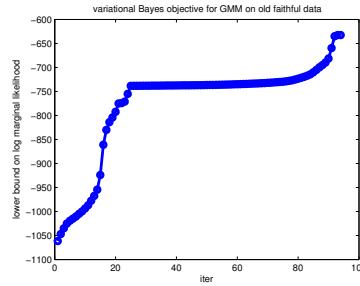


Figure 10.9: Lower bound vs iterations for the VB algorithm in Figure 10.7. The steep parts of the curve correspond to places where the algorithm figures out that it can increase the bound by “killing off” unnecessary mixture components, as described in Section 10.2.6.6. The plateaus correspond to slowly moving the clusters around. Generated by [gmm\\_vb\\_em.ipynb](#).

1 **10.2.6.6 Model selection using VBEM**

3 Section 10.2.6.4 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

$$\underline{8} \quad p(K|\mathcal{D}) = \frac{p(\mathcal{D}|K)}{\sum_{K'} p(\mathcal{D}|K')} \approx \frac{e^{\mathcal{L}(K)}}{\sum_{K'} e^{\mathcal{L}(K')}} \quad (10.99)$$

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 **10.2.7 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.,

$$\underline{22} \quad p(\mathbf{z}_j|\mathbf{z}_{-j}, \mathbf{x}) = h(\mathbf{z}_j) \exp[\boldsymbol{\eta}_j(\mathbf{z}_{-j}, \mathbf{x})^\top \mathcal{T}(\mathbf{z}_j) - A_j(\boldsymbol{\eta}_j(\mathbf{z}_{-j}, \mathbf{x}))] \quad (10.100)$$

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 (10.100) holds, the mean field update node  $j$  becomes  
28

$$\underline{29} \quad q_j(\mathbf{z}_j) \propto \exp [\mathbb{E} [\log p(\mathbf{z}_j|\mathbf{z}_{-j}, \mathbf{x})]] \quad (10.101)$$

$$\underline{30} \quad = \exp \left[ \log h(\mathbf{z}_j) + \mathbb{E} [\boldsymbol{\eta}_j(\mathbf{z}_{-j}, \mathbf{x})]^\top \mathcal{T}(\mathbf{z}_j) - \mathbb{E} [A_j(\boldsymbol{\eta}_j(\mathbf{z}_{-j}, \mathbf{x}))] \right] \quad (10.102)$$

$$\underline{32} \quad \propto h(\mathbf{z}_j) \exp \left[ \mathbb{E} [\boldsymbol{\eta}_j(\mathbf{z}_{-j}, \mathbf{x})]^\top \mathcal{T}(\mathbf{z}_j) \right] \quad (10.103)$$

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

$$\underline{37} \quad \boldsymbol{\psi}_j = \mathbb{E} [\boldsymbol{\eta}_j(\mathbf{z}_{-j}, \mathbf{x})] \quad (10.104)$$

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 9.2): 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.

### 10.2.8 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).

## 10.3 Fixed-form VI

Recall that the goal of variational inference is to maximize the following lower bound wrt the variational parameters  $\psi$ :

$$\mathcal{L}(\psi|\mathcal{D}) = \mathbb{E}_{q_\psi(\mathbf{z})} \left[ \log \frac{p(\mathbf{z})p(\mathcal{D}|\mathbf{z})}{q_\psi(\mathbf{z})} \right] = \mathbb{E}_{q_\psi} [\ell_\psi(\mathbf{z})] \quad (10.105)$$

where

$$\ell_\psi(\mathbf{z}) = \log p(\mathbf{z}, \mathcal{D}) - \log q_\psi(\mathbf{z}) \quad (10.106)$$

Here  $\mathbf{z}$  are the unknown latent variables (or parameters), and  $\psi$  are the variational parameters. In the mean field method of Section 10.2, we assumed that  $q(\mathbf{z})$  factorized across (groups of) variables,  $q(\mathbf{z}|\psi) = \prod_{j=1}^J q_j(z_j)$ . We do not specify the form of  $q_j$ , so this approach is called “free-form” VI; however, the form of the optimal  $q_j$  can be derived analytically. We then optimized each  $q_j$  using coordinate ascent.

In this section, we take a different approach: we pick any convenient form we like for  $q(\mathbf{z})$ , such as a multivariate Gaussian, and then we directly maximize the ELBO using gradient ascent. This is called **fixed-form VI**.

### 10.3.1 Stochastic variational inference

In many models, the likelihood factorizes into a product of terms, in which case we have

$$\ell_\psi(\mathbf{z}) = \left[ \sum_{n=1}^N \log p(\mathbf{x}_n|\mathbf{z}) \right] + \log p(\mathbf{z}) - \log q_\psi(\mathbf{z}) \quad (10.107)$$

If  $N$  is large, we can compute an unbiased minibatch approximation to this expression as follows:

$$\hat{\ell}_\psi(\mathbf{z}) = \left[ \frac{N}{B} \sum_{b=1}^B \log p(\mathbf{x}_b|\mathbf{z}) \right] + \log p(\mathbf{z}) - \log q_\psi(\mathbf{z}) \quad (10.108)$$

1 We can use this to create an unbiased Monte Carlo approximation to the ELBO:  
2

3  $\hat{L}(\psi|\mathcal{D}) = \mathbb{E}_{q_\psi} [\hat{\ell}_\psi(z)]$  (10.109)  
4

5 This is called **stochastic variational inference** or **SVI** [Hof+13], and allows VI to scale to large  
6 datasets.  
7

### 8 10.3.2 Black-box variational inference

9 In this section, we assume that we can evaluate  $\ell_\psi(z)$  pointwise, but we do not assume we can take  
10 gradients of this function. (For example,  $z$  may contain discrete variables.) We are thus treating the  
11 model as a “blackbox”. Hence this approach is called **black box variational inference** or **BBVI**  
12 [RGB14; ASD20].

13 To estimate the gradient of the ELBO, we will use the **score function estimator**, also called  
14 the **REINFORCE** estimator (Section 6.5.3). To derive this, we use the fact that  $\nabla \log q = \frac{\nabla q}{q}$  to  
15 conclude that  $\nabla q = q \nabla \log q$ . (This is called the **log derivative trick**.) We also exploit the fact  
16 that  $\int q(z) dz = 1$ . With this, the gradient of the ELBO can be derived as follows:

17  $\nabla_\psi \hat{L}(\psi) = \nabla_\psi \int q_\psi(z) \log \frac{p(z, \mathcal{D})}{q_\psi(z)} dz$  (10.110)  
18

19  $= \int [\nabla_\psi q_\psi(z)] \left[ \log \frac{p(z, \mathcal{D})}{q_\psi(z)} \right] dz + \int [q_\psi(z)] \left[ \nabla_\psi \log \frac{p(z)p(\mathcal{D}|z)}{q_\psi(z)} \right] dz$  (10.111)  
20

21  $= \int [\nabla_\psi q_\psi(z)] \left[ \log \frac{p(z)p(\mathcal{D}|z)}{q_\psi(z)} \right] dz - \int q_\psi(z) \nabla_\psi \log q_\psi(z) dz$  (10.112)  
22

23  $= \int [\nabla_\psi q_\psi(z)] \left[ \log \frac{p(z)p(\mathcal{D}|z)}{q_\psi(z)} \right] dz - \int \nabla_\psi q_\psi(z) dz$  (10.113)  
24

25  $= \int [\nabla_\psi q_\psi(z)] \left[ \log \frac{p(z)p(\mathcal{D}|z)}{q_\psi(z)} \right] dz - \nabla_\psi \int q_\psi(z) dz$  (10.114)  
26

27  $= \int [\nabla_\psi q_\psi(z)] \left[ \log \frac{p(z)p(\mathcal{D}|z)}{q_\psi(z)} \right] dz$  (10.115)  
28

29  $= \int q_\psi(z) (\nabla_\psi \log q_\psi(z) \ell_\psi(z)) dz$  (10.116)  
30

31  $= \mathbb{E}_{q_\psi(z)} [\nabla_\psi \log q_\psi(z) \ell_\psi(z)]$  (10.117)  
32

33 We can compute a stochastic approximation to this gradient by sampling  $z_s \sim q_\psi(z)$  and then  
34 computing  
35

36  $\widehat{\nabla_\psi \hat{L}}(\psi_t) = \frac{1}{S} \sum_{s=1}^S \nabla_\psi \log q_\psi(z_s) \ell_\psi(z_s) |_{\psi=\psi_t}$  (10.118)  
37

38 We can pass this to any kind of gradient optimizer, such as SGD or Adam. If we further approximate  
39  $\ell_\psi(z_s)$  using the SVI minibatch approximation, we get a “**doubly stochastic**” approximation  
40 [TLG14].  
41

42

In practice, the variance of this estimator is quite large, so it is important to use methods such as **control variates** or **CV** (Section 6.5.3.1). To see how this works, consider the naive gradient estimator in Equation (10.118), which for the  $i$ 'th component we can write as

$$\nabla_{\psi_i} \widehat{L}(\psi_t)^{\text{naive}} = \frac{1}{S} \sum_{s=1}^S \tilde{g}_i(\mathbf{z}_s) \quad (10.119)$$

$$\tilde{g}_i(\mathbf{z}_s) = g_i(\mathbf{z}_s) \times \ell_{\psi}(\mathbf{z}_s) \quad (10.120)$$

$$g_i(\mathbf{z}_s) = \nabla_{\psi_i} \log q_{\psi}(\mathbf{z}_s) \quad (10.121)$$

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

$$\tilde{g}_i^{cv}(\mathbf{z}) = \tilde{g}_i(\mathbf{z}) + c_i(\mathbb{E}[b_i(\mathbf{z})] - b_i(\mathbf{z})) \quad (10.122)$$

where  $b_i(\mathbf{z})$  is a baseline function and  $c_i$  is some constant, to be specified below. A convenient baseline is the score function,  $b_i(\mathbf{z}) = \nabla_{\psi_i} \log q_{\psi_i}(\mathbf{z})$ , since this is correlated with  $\tilde{g}_i(\mathbf{z})$ , and has the property that  $\mathbb{E}[b_i(\mathbf{z})] = \mathbf{0}$ , since the expected value of the score function is zero, as we showed in Equation (2.245). Hence

$$\tilde{g}_i^{cv}(\mathbf{z}) = \tilde{g}_i(\mathbf{z}) - c_i g_i(\mathbf{z}) = g_i(\mathbf{z})(\ell_{\psi}(\mathbf{z}) - c_i) \quad (10.123)$$

so the CV estimator is given by

$$\nabla_{\psi_i} \widehat{L}(\psi_t)^{\text{cv}} = \frac{1}{S} \sum_{s=1}^S g_i(\mathbf{z}_s) \times (\ell_{\psi}(\mathbf{z}_s) - c_i) \quad (10.124)$$

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

$$c_i = \frac{\text{Cov}[g_i(\mathbf{z})\ell_{\psi}(\mathbf{z}), g_i(\mathbf{z})]}{\mathbb{V}[g_i(\mathbf{z})]} \quad (10.125)$$

which can be estimated by sampling  $\mathbf{z} \sim q_{\psi}(\mathbf{z})$ . Thus the overall algorithm is as shown in Algorithm 20.

We can stop the algorithm when the lower bound stops increasing. We can compute a stochastic approximation to the lower bound using

$$\hat{L}(\psi) = \frac{1}{S} \sum_{s=1}^S \ell_{\psi}(\mathbf{z}_s) \quad (10.126)$$

where  $\mathbf{z}_s \sim q_{\psi}(\mathbf{z})$ . To smooth out the noise, we can use a running average over the last  $w$  observations to get

$$\bar{L}(\psi_t) = \frac{1}{w} \sum_{k=1}^w \hat{L}(\psi_{t-k+1}) \quad (10.127)$$

If the moving average does not improve after  $P$  consecutive iterations, we declare convergence, where  $P$  is the **patience** parameter. Typical values are  $P = 20$  and  $w = 20$  [TND21].

---

1  
2   **Algorithm 20:** Blackbox VI with control variates  
3   1 Initialize  $\psi_0$   
4   2  $\mathbf{z}_s \sim q_{\psi_0}(\mathbf{z})$ ,  $s = 1 : S$   
5   3  $h_{\psi_0}(\mathbf{z}_s) = \log p(\mathbf{z}_s, \mathcal{D}) - \log q_{\psi_0}(\mathbf{z}_s)$ ,  $s = 1 : S$   
6   4  $\mathbf{g}_0 = \frac{1}{S} \sum_{s=1}^S [\nabla_{\psi} \log q_{\psi_0}(\mathbf{z}_s)] h_{\psi_0}(\mathbf{z}_s)$   
7   5 Compute  $\mathbf{c}_1$  using Equation (10.125) applied to  $\mathbf{z}_s$   
8   6 **for**  $t = 1 : T$  **do**  
9   7    $\mathbf{z}_s \sim q_{\psi_t}(\mathbf{z})$ ,  $s = 1 : S$   
10   8    $h_{\psi_t}(\mathbf{z}_s) = \log p(\mathbf{z}_s, \mathcal{D}) - \log q_{\psi_t}(\mathbf{z}_s)$ ,  $s = 1 : S$   
11   9    $\mathbf{g}_t = \frac{1}{S} \sum_{s=1}^S [\nabla_{\psi} \log q_{\psi_t}(\mathbf{z}_s)] \odot (h_{\psi_t}(\mathbf{z}_s) - \mathbf{c}_t)$   
12   10   Compute  $\mathbf{c}_{t+1}$  using Equation (10.125) applied to  $\mathbf{z}_s$   
13   11    $\psi_t = \text{gradient-update}(\psi_{t-1}, \mathbf{g}_t)$

---

15  
16  
17  
18  
19  
20

### 21 10.3.3 Reparameterization VI

22 In this section, we exploit the **reparameterization trick** from Section 6.5.4 to get a lower variance  
23 estimator for the gradient. This assumes that  $\ell_{\psi}(\mathbf{z})$  is a differentiable function of  $\mathbf{z}$ . It also assumes  
24 that we can sample  $\mathbf{z} \sim q_{\psi}(\mathbf{z})$  by first sampling a noise term  $\epsilon \sim q_0(\epsilon)$ , and then transforming it to  
25 compute the latent random variables  $\mathbf{z} = r(\psi, \epsilon)$ . In this case, the ELBO becomes  
26

27  
28  
29

$$\mathbb{L}(\psi) = \mathbb{E}_{q_0(\epsilon)} [\ell_{\psi}(r(\psi, \epsilon))] = \mathbb{E}_{q_0(\epsilon)} [\log p(r(\psi, \epsilon), \mathcal{D}) - \log q_{\psi}(r(\psi, \epsilon))] \quad (10.128)$$

30  
31 Since the sampling distribution  $q_0(\epsilon)$  is independent of the variational parameters  $\psi$ , we can push  
32 the gradient operator inside the expectation, and thus we can estimate the gradient using standard  
33 automatic differentiation methods, as shown in Algorithm 21. This is called **reparameterized VI**  
34 or **RVI**, and has provably lower variance than BBVI in certain cases [Xu+19].  
35  
36

---

37   **Algorithm 21:** Estimate of ELBO gradient

---

38   1 def elbo( $\psi$ ):  
39   2    $\epsilon \sim q_0(\epsilon)$   
40   3    $\mathbf{z} = r(\psi, \epsilon)$   
41   4   return  $\log p(\mathbf{z}, \mathcal{D}) - q_{\psi}(\mathbf{z}|\mathcal{D})$   
42   5 def elbo-grad( $\psi$ ):  
43   6   return grad(elbo( $\psi$ ))

---

44  
45  
46  
47

---

1    **10.3.3.1 “Sticking the landing” estimator**

2    Applying the results from Section 6.5.4.2, we can derive the gradient estimate of the reparameterized  
3    ELBO, for a single Monte Carlo sample, as follows:

4     $\nabla_{\psi} \hat{L}(\psi, \epsilon) = \nabla_{\psi} [\log p(z, \mathcal{D}) - \log q_{\psi}(z|\mathcal{D})] \quad (10.129)$

5     $= \underbrace{\nabla_z [\log p(z|\mathcal{D}) - \log q_{\psi}(z|\mathcal{D})]}_{\text{path derivative}} \mathbf{J} - \underbrace{\nabla_{\psi} \log q_{\psi}(z|\mathcal{D})}_{\text{score function}} \quad (10.130)$

6    where  $z = r(\psi, \epsilon)$  and  $\mathbf{J} = \nabla_{\psi} r(\psi, \epsilon)$  is the Jacobian matrix of the noise transformation.

7    The first term is the indirect effect of  $\psi$  on the objective via the generated samples  $z$ . The second  
8    term is the direct effect of  $\psi$  on the objective. The second term is zero in expectation since it is  
9    the score function (see Equation (2.245)), but it may be non-zero for a finite number of samples,  
10   even if  $q_{\psi}(z|\mathcal{D}) = p(z|\mathcal{D})$  is the true posterior. In a paper called “**sticking the landing**”, [RWD17]  
11   propose to drop the second term to create a lower variance estimator.<sup>2</sup> In practice, this means we  
12   compute the gradient using Algorithm 22 instead of Algorithm 21.<sup>3</sup>

---

13   **Algorithm 22:** “Sticking the landing” estimator of ELBO gradient

---

14   1 def elbo-pd( $\psi$ ):  
15   2      $\epsilon \sim q_0(\epsilon)$   
16   3      $z = r(\psi, \epsilon)$   
17   4      $\psi' = \text{stop-gradient}(\psi)$   
18   5     return  $\log p(z, \mathcal{D}) - q_{\psi'}(z|\mathcal{D})$   
19   6 def elbo-grad-pd( $\psi$ ):  
20   7     return grad(elbo-pd( $\psi$ ))

---

21   Note that the STL estimator is not always better than the “standard” estimator. In [GD20], they  
22   propose to use a weighted combination of estimators, where the weights are optimized so as to reduce  
23   variance for a fixed amount of compute.

24   **10.3.3.2 Example: reparameterized SVI for GMMs**

25   In this section, we use reparameterized SVI to fit a Gaussian mixture model. We will marginalize out  
26   the discrete latent variables, so just need to approximate  $p(\theta|\mathcal{D})$ . We choose a factored variational  
27   posterior that is conjugate to the likelihood, but is also reparameterizable, so we can fit the posterior  
28   with SGD instead of having to use coordinate ascent (Section 10.2.1).

29   For simplicity, we assume diagonal covariance matrices. Thus the likelihood for one data point,  
30    $x \in \mathbb{R}^D$ , is

31    $p(x|\theta) = \sum_{k=1}^K \pi_k \mathcal{N}(x|\mu_k, \text{diag}(\lambda_k)^{-1}) \quad (10.131)$

---

32   2. The expression “to stick a landing” means to land firmly on one’s feet after performing a gymnastics move. In the  
33   current context, the analogy is this: if the variational posterior is optimal, so  $q_{\psi}(z|\mathcal{D}) = p(z|\mathcal{D})$ , then we want our  
34   objective to be 0, and not to “wobble” with Monte Carlo noise.

35   3. The difference is that the path derivative version ignores the score function. This can be achieved by using  
36    $\log q_{\psi'}(z|\mathcal{D})$ , where  $\psi'$  is a “disconnected” copy of  $\psi$  that does not affect the gradient.

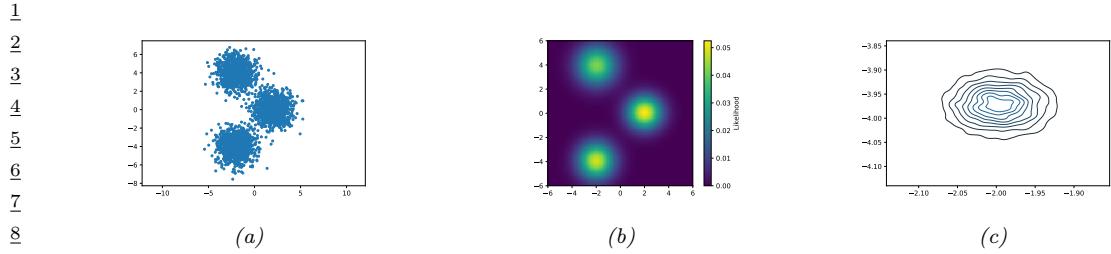


Figure 10.10: SVI for fitting a mixture of 3 Gaussians in 2d. (a) 3000 training points. (b) Fitted density, plugging in the posterior mean parameters. (c) Kernel density estimate fit to 10,000 samples from  $q(\mu_1|\psi_\theta)$ . Generated by [svi\\_gmm\\_demo\\_2d.ipynb](#).

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:

$$p(\boldsymbol{\theta}) = \left[ \prod_{k=1}^K \prod_{d=1}^D \mathcal{N}(\mu_{kd}|0, 1) \text{Ga}(\lambda_{kd}|5, 5) \right] \text{Dir}(\boldsymbol{\pi}|\mathbf{1}) \quad (10.132)$$

We assume the following mean field posterior:

$$q(\boldsymbol{\theta}|\psi_\theta) = \left[ \prod_{k=1}^K \prod_{d=1}^D \mathcal{N}(\mu_{kd}|m_{kd}, s_{kd}) \text{Ga}(\lambda_{kd}|\alpha_{kd}, \beta_{kd}) \right] \text{Dir}(\boldsymbol{\pi}|\mathbf{c}) \quad (10.133)$$

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

We can compute the ELBO using

$$\mathcal{L}(\psi_\theta|\mathcal{D}) = \mathbb{E}_{q(\boldsymbol{\theta}|\psi_\theta)} \left[ \sum_{n=1}^N \log p(\mathbf{x}_n|\boldsymbol{\theta}) \right] - D_{\text{KL}}(q(\boldsymbol{\theta}|\psi_\theta) \parallel p(\boldsymbol{\theta})) \quad (10.134)$$

We can approximate the first term using minibatching. Since  $q(\boldsymbol{\theta}|\psi_\theta)$  is reparameterizable (see Section 10.3.3), we can sample from it and push gradients inside. If we use a single posterior sample per minibatch,  $\boldsymbol{\theta}^s \sim q(\boldsymbol{\theta}|\psi_\theta)$ , we get

$$\nabla_{\psi_\theta} \mathcal{L}(\psi_\theta|\mathcal{D}) \approx \frac{N}{B} \sum_{b=1}^B \nabla_{\psi_\theta} \log p(\mathbf{x}_b|\boldsymbol{\theta}^s) - \nabla_{\psi_\theta} D_{\text{KL}}(q(\boldsymbol{\theta}|\psi_\theta) \parallel p(\boldsymbol{\theta})) \quad (10.135)$$

We can now optimize this with SGD.

Figure 10.10 gives an example of this in practice. We generate a dataset from a mixture of 3 Gaussians in 2d, using  $\boldsymbol{\mu}_1^* = [2, 0]$ ,  $\boldsymbol{\mu}_2^* = [-2, -4]$ ,  $\boldsymbol{\mu}_3^* = [-2, 4]$ , precisions  $\lambda_{dk}^* = 1$ , and uniform mixing weights,  $\boldsymbol{\pi}^* = [1/3, 1/3, 1/3]$ . Figure 10.10a 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 10.10b shows the predictions of the fitted model. More precisely, it shows  $p(\mathbf{x}|\bar{\boldsymbol{\theta}})$ , where  $\bar{\boldsymbol{\theta}} = \mathbb{E}_{q(\boldsymbol{\theta}|\psi_\theta)}[\boldsymbol{\theta}]$ . Figure 10.10c shows a kernel density estimate fit to 10,000 samples from  $q(\mu_1|\psi_\theta)$ . We see that the posterior mean is  $\mathbb{E}[\mu_1] \approx [-2, -4]$ . Due to label switching unidentifiability, we see this matches  $\boldsymbol{\mu}_2^*$  rather than  $\boldsymbol{\mu}_1^*$ .

---

### 10.3.4 Gaussian VI

The most widely used RVI approximation is when  $q_\psi(\mathbf{z})$  is a Gaussian, where  $\psi = (\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . Following [TND21], we call this **Gaussian VI**. We give a brief summary below. For a theoretical analysis of the computational / statistical tradeoff with using this form of posterior approximation, see [Bha+22].

#### 10.3.4.1 Full-rank Gaussian VI

In this section, we represent the covariance using its Cholesky decomposition,  $\boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}^\top$ , where  $\mathbf{L} = \text{tril}(\mathbf{l})$  is a lower triangular matrix, as in [TLG14; TN18]. The variational parameters are  $\psi = (\boldsymbol{\mu}, \mathbf{l})$ . The noise transformation has the form

$$13 \quad \mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \iff \mathbf{z} = \boldsymbol{\mu} + \mathbf{L}\boldsymbol{\epsilon} \quad (10.136)$$

15 where  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , so  $r(\psi, \boldsymbol{\epsilon}) = \boldsymbol{\mu} + \mathbf{L}\boldsymbol{\epsilon}$ .

16 For any given sample  $\mathbf{z}$ , we can use automatic differentiation to compute the derivative of the  
17 sampled ELBO:

$$19 \quad \nabla_\psi \ell_\psi(\mathbf{z}(\psi)) = \nabla_\psi \log p(\mathbf{z}(\psi), \mathcal{D}) - \nabla_\psi \log q_\psi(\mathbf{z}(\psi)) \quad (10.137)$$

#### 10.3.4.2 Low-rank Gaussian VI

23 In high dimensions, an efficient alternative to using a Cholesky decomposition is the factor decompo-  
24 sition

$$26 \quad \boldsymbol{\Sigma} = \mathbf{B}\mathbf{B}^\top + \mathbf{C}^2 \quad (10.138)$$

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

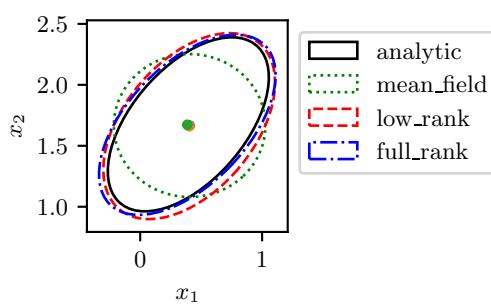
32 In the special case where  $f = 1$ , the covariance matrix becomes

$$34 \quad \boldsymbol{\Sigma} = \mathbf{b}\mathbf{b}^\top + \text{diag}(\mathbf{c}^2) \quad (10.139)$$

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

40 In particular, let  $\mathbf{g} = (\mathbf{g}_1, \mathbf{g}_2, \mathbf{g}_3)$  be the regular gradient of the ELBO wrt  $\boldsymbol{\mu}$ ,  $\mathbf{b}$  and  $\mathbf{c}$ . Let  
41  $\mathbf{v}_1 = \mathbf{c}^2 - 2\mathbf{b}^2 \odot \mathbf{c}^{-4}$ ,  $\mathbf{v}_2 = \mathbf{b}^2 \odot \mathbf{c}^{-3}$ ,  $\kappa_1 = \sum_{i=1}^d b_i^2/c_i^2$ , and  $\kappa_2 = \frac{1}{2}(1 + \sum_{i=1}^d v_{2i}^2/v_{1i})^{-1}$ . Then the  
42 natural gradient is given by

$$44 \quad \mathbf{g}^{\text{nat}} = \begin{pmatrix} (\mathbf{g}_1^\top \mathbf{b})\mathbf{b} + \mathbf{c}^2 \odot \mathbf{g}_1 \\ \frac{1+\kappa_1}{2\kappa_1}((\mathbf{g}_2^\top \mathbf{b})\mathbf{b} + \mathbf{c}^2 \odot \mathbf{g}_2) \\ \frac{1}{2}\mathbf{v}_1^{-1} \odot \mathbf{g}_3 + \kappa_2[(\mathbf{v}_1^{-1} \odot \mathbf{v}_2)^\top \mathbf{g}_3](\mathbf{v}_1^{-1} \odot \mathbf{v}_2) \end{pmatrix} \quad (10.140)$$



11 Figure 10.11: Gaussian variational approximation to a Gaussian posterior for the mean of a 2d Gaussian.  
12 Generated by [gaussian\\_2d\\_vi.ipynb](#).

13

#### 14 10.3.4.3 Example: GVI for a linear Gaussian system

15 As a sanity check, we try to approximate the posterior mean of a multivariate Gaussian with fixed  
16 covariance. We use a Gaussian prior, so the joint distribution has the form

$$\frac{19}{20} p(\mathbf{z}|\mathcal{D}) \propto p(\mathbf{z}) \prod_{n=1}^N \mathcal{N}(\mathbf{y}_n|\mathbf{z}, \boldsymbol{\Sigma}) \quad (10.141)$$

22 where the measurement noise  $\boldsymbol{\Sigma}$  is a fixed, non-diagonal matrix. If we use a conjugate prior,  
23  $\mathcal{N}(\mathbf{z}|\tilde{\mathbf{m}}, \tilde{\mathbf{V}})$ , then the exact posterior,  $p(\mathbf{z}|\mathcal{D}) = \mathcal{N}(\mathbf{z}|\hat{\mathbf{m}}, \hat{\mathbf{V}})$ , can be computed analytically, as  
24 discussed in Section 3.3.1.

25 We also compute a Gaussian variational approximation,  $q(\mathbf{z}|\psi_\mu, \psi_\Sigma)$ , where  $\psi_\Sigma$  is either full-rank,  
26 diagonal, or rank-1 plus diagonal. The results are shown in Figure 10.11. We see that the full-rank  
27 approximation matches the true posterior, as expected. However, the diagonal approximation is  
28 overconfident, which is a well-known flaw of variational inference.  
29

#### 30 10.3.5 Automatic differentiation VI

32 To apply Gaussian VI, we need to transform constrained parameters (such as variance terms) to  
33 unconstrained form, so they live in  $\mathbb{R}^D$ . This technique can be used for any distribution for which  
34 we can define a bijection to  $\mathbb{R}^D$ . This approach is called **automatic differentiation variational**  
35 **inference** or **ADVI** [Kuc+16]. We give the details below.  
36

##### 37 10.3.5.1 Basic idea

38 Our goal is to approximate the posterior  $p(\boldsymbol{\theta}|\mathcal{D}) \propto p(\boldsymbol{\theta})p(\mathcal{D}|\boldsymbol{\theta})$ , where  $\boldsymbol{\theta} \in \Theta$  lives in some  $D$ -  
39 dimensional constrained parameter space. Let  $T : \Theta \rightarrow \mathbb{R}^D$  be a bijective mapping that maps from  
40 the constrained space to some unconstrained space, with inverse  $T^{-1} : \mathbb{R}^D \rightarrow \Theta$  that maps to the  
41 constrained domain of the prior. Let  $\mathbf{z} = T(\boldsymbol{\theta})$  be the unconstrained latent variables. We will use a  
42 block-factored Gaussian variational approximation to the posterior for  $\mathbf{z}$ , i.e.,:

$$\frac{44}{45} q(\mathbf{z}|\psi) = \prod_{b=1}^B \mathcal{N}(z_b|\mu_b, \Sigma_b) \quad (10.142)$$

1 where  $\psi = (\mu_{1:B}, \Sigma_{1:B})$ .  
 2

3 By the change of variable formula Equation (2.286), we have  
 4

$$5 p(\mathbf{z}) = p(T^{-1}(\mathbf{z})) |\det(\mathbf{J}_{T^{-1}}(\mathbf{z}))| \quad (10.143)$$

6 where  $\mathbf{J}_{T^{-1}}$  is the Jacobian of the inverse mapping  $\mathbf{z} \rightarrow \theta$ . Hence the ELBO becomes  
 7

$$8 \mathcal{L}(\psi) = E_{\mathbf{z} \sim q(\mathbf{z}|\psi)} [p(\mathcal{D}|T^{-1}(\mathbf{z})) + \log p(T^{-1}(\mathbf{z})) + \log |\det(\mathbf{J}_{T^{-1}}(\mathbf{z}))|] + \mathbb{H}(\psi) \quad (10.144)$$

9 We can use a Monte Carlo approximation of the expectation over  $\mathbf{z}$ , together with the reparameterization trick, which (in the fully diagonal case) replaces  $\mathbf{z} \sim q(\mathbf{z}|\psi)$  with  $\mathbf{z} = \mu + \sigma \odot \epsilon$  where  
 10  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and  $\sigma = \sigma_{1:D}$ . We can compute the final entropy term in closed form. In the case of a  
 11 fully diagonal approximation, it follows from Equation (5.96) that  
 12

$$13 \mathbb{H}(\psi) = \sum_{i=1}^D 2 \log(\sigma_i) + \text{const} \quad (10.145)$$

14 Since the objective is stochastic, we can use SGD to optimize it. However, [Ing20] propose  
 15 **deterministic ADVI**, in which the samples  $\epsilon_s \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  are held fixed during the optimization  
 16 process. This is called the common random numbers trick (Section 11.6.1), and makes the objective  
 17 a deterministic function; this allows for the use of more powerful second-order optimization methods,  
 18 such as BFGS. (Of course, if the dataset is large, we might need to use minibatch subsampling, which  
 19 reintroduces stochasticity.)  
 20

### 21 10.3.5.2 Example: ADVI for Beta-Binomial model

22 To illustrate ADVI, we consider the 1d beta-binomial model from Section 7.4.4. We want to  
 23 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  
 24 the sufficient statistics are  $N_1 = 10$ ,  $N_0 = 1$ , and the prior is uninformative,  $a = b = 1$ . We use the  
 25 transformation  $\theta = T^{-1}(z) = \sigma(z)$ , and optimize the ELBO with SGD. The results of this method  
 26 are shown in Figure 7.3 and show that the Gaussian fit is a good approximation, despite the skewed  
 27 nature of the posterior.  
 28

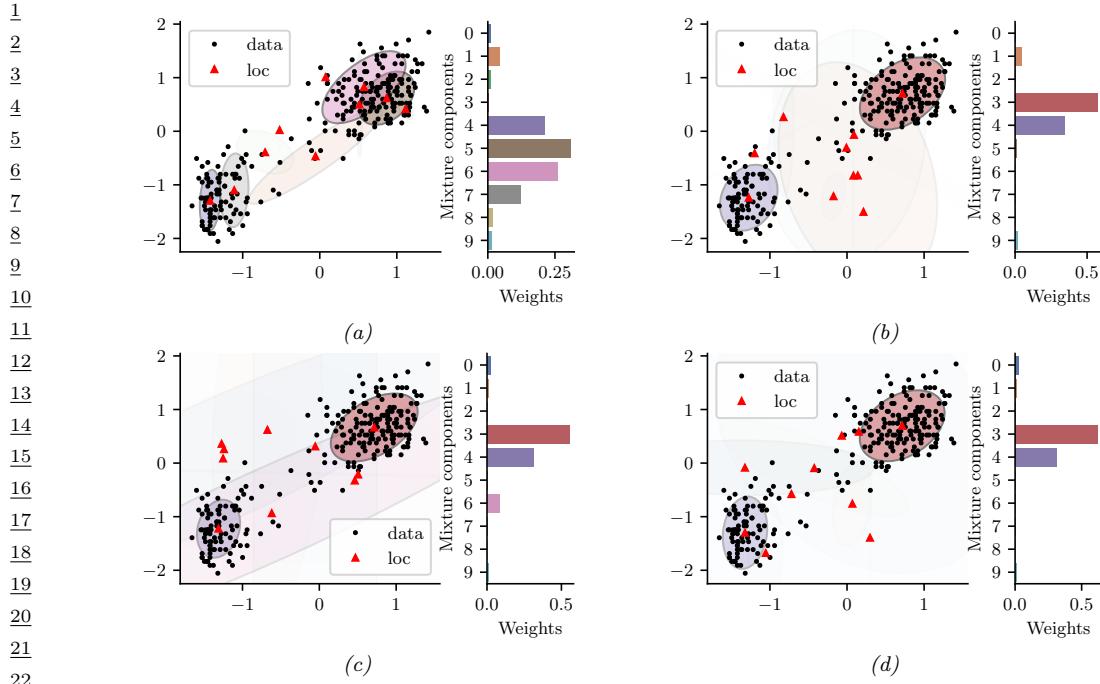
### 29 10.3.5.3 Example: ADVI for GMMs

30 In this section, we use ADVI to approximate the posterior of the parameters of a mixture of Gaussians.  
 31 The difference from the VBEM algorithm of Section 10.2.6 is that we use ADVI combined with a  
 32 Gaussian variational posterior, rather than using a mean field approximation defined by a product of  
 33 conjugate distributions.  
 34

35 To apply ADVI, we marginalize out the discrete local discrete latents  $m_n \in \{1, \dots, K\}$  analytically,  
 36 so the likelihood has the form  
 37

$$38 p(\mathcal{D}|\theta) = \prod_{n=1}^N \left[ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{y}_n | \boldsymbol{\mu}_k, \text{diag}(\boldsymbol{\Sigma}_k)) \right] \quad (10.146)$$

39 We use an uninformative Gaussian prior for the  $\boldsymbol{\mu}_k$ , a uniform LKJ prior for the  $\mathbf{L}_k$ , a log-Normal  
 40 prior for the  $\boldsymbol{\sigma}_k$ , and a uniform Dirichlet prior for the mixing weights  $\boldsymbol{\pi}$ . (See [Kuc+16, Fig 21]  
 41



23 *Figure 10.12: Posterior over the mixing weights (histogram) and the means and covariances of each Gaussian*  
 24 *mixture component, using  $K = 10$ , when fitting the model to the Old Faithful dataset from Figure 10.7.* (a)  
 25 *MAP approximation. (b-d) 3 samples from the Gaussian approximation. The intensity of the shading is*  
 26 *proportional to the mixture weight. Generated by [gmm\\_advi\\_bijax.ipynb](#).*

31 for a definition of the model in STAN syntax.) The posterior approximation for the unconstrained  
 32 parameters is a block-diagonal gaussian,  $q(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\psi}_\mu, \boldsymbol{\psi}_\Sigma)$ , where the unconstrained parameters  
 33 are computed using suitable bijections (see code for details).

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

41

42

#### 43 10.3.5.4 More complex posteriors

44

45 We can combine ADVI with any of the improved posterior approximations that we discuss in  
 46 Section 10.4, such as Gaussian mixtures [Mor+21b] or normalizing flows [ASD20].

47

---

### 10.3.6 Amortized inference

Suppose we want to perform parameter estimation in a model with local latent variables:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \sum_{n=1}^N \log \sum_{z_n} p_{\theta}(x_n, z_n) \quad (10.147)$$

To compute the marginal likelihood, we first need to compute the posteriors  $p_{\theta}(z_n|x_n)$  for each example  $n$ ; the normalization constant then gives us  $\log p_{\theta}(x_n)$ . When using BBVI or RVI for this, we therefore have to solve an optimization problem for each  $q(z_n|\psi_n)$ , which can be slow.

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

$$q(z_n|\psi_n) = q(z_n|f_{\phi}^{\text{inf}}(x_n)) = q_{\phi}(z_n|x_n) \quad (10.148)$$

The “**amortized ELBO**”, for a model with local latents and fixed global parameters, becomes

$$\mathcal{L}(\phi, \theta | \mathcal{D}) = \frac{1}{N} \sum_{n=1}^N [\mathbb{E}_{q_{\phi}(z_n|x_n)} [\log p_{\theta}(x_n, z_n) - \log q_{\phi}(z_n|x_n)]] \quad (10.149)$$

We can approximate this by sampling a single data point  $x_n \sim p_{\mathcal{D}}$ , and then sampling a single latent  $z_n \sim q_{\phi}(z_n|x_n)$ , as in DSVI, to get

$$\tilde{\mathcal{L}}(\phi, \theta | x_n, z_n) = \log p_{\theta}(x_n, z_n) - \log q_{\phi}(z_n) \quad (10.150)$$

(We call this the “**per-sample ELBO**”, although [Ble17] call it the **instantaneous ELBO**.) If the posteriors are reparameterizable, we can push gradients inside and then apply SGD. See Algorithm 23 for the resulting pseudocode.

---

#### Algorithm 23: Amortized SVI

---

```

1 Initialize  $\theta, \phi$ 
2 repeat
3   Sample  $x_n \sim p_{\mathcal{D}}$ 
4   Sample  $z_n \sim q_{\phi}(z|x_n)$ 
5    $\theta := \theta + \eta \nabla_{\theta} \tilde{\mathcal{L}}(\phi, \theta | x_n, z_n)$ 
6    $\phi := \phi + \eta \nabla_{\phi} \tilde{\mathcal{L}}(\phi, \theta | x_n, z_n)$ 
7   Update learning rate  $\eta$ 
8 until converged

```

---

This method is very 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  $\psi_n$ .

<sup>1</sup> This is called the **amortization gap** [CLD18]. We can close this gap by using the inference network  
<sup>2</sup> to warm-start an optimizer for  $\psi_n$ ; this is known as **semi-amortized VI** [Kim+18c]. The key  
<sup>3</sup> insight is that the local SVI procedure is itself differentiable, so the inference network and generative  
<sup>4</sup> model can be trained end-to-end. (See also [MYM18], who propose a closely related method called  
<sup>5</sup> **iterative amortized inference**.)  
<sup>6</sup>

<sup>7</sup> An alternative approach is to use the inference network as a proposal distribution. If we combine  
<sup>8</sup> this with importance sampling, we get the IWAE bound of Section 10.5.1. If we use this with  
<sup>9</sup> Metropolis Hastings, we get a VI-MCMC hybrid (see Section 10.4.5).  
<sup>10</sup>

## <sup>12</sup> 10.4 More accurate variational posteriors

<sup>13</sup>  
<sup>14</sup> In general, we can improve the tightness of the ELBO lower bound, and hence reduce the KL  
<sup>15</sup> divergence of our posterior approximation, if we use more flexible posterior families (although  
<sup>16</sup> optimizing within more flexible families may be slower, and can incur statistical error if the sample  
<sup>17</sup> size is low [Bha+21]). In this section, we give several examples of more accurate variational posteriors,  
<sup>18</sup> going beyond fully factored mean field approximations.  
<sup>19</sup>

### <sup>21</sup> 10.4.1 Structured mean field

<sup>22</sup>  
<sup>23</sup> The mean field assumption is quite strong, and can sometimes give poor results. Fortunately,  
<sup>24</sup> sometimes we can exploit **tractable substructure** in our problem, so that we can efficiently handle  
<sup>25</sup> some kinds of dependencies between the variables in the posterior in an analytic way, rather than  
<sup>26</sup> assuming they are all independent. This is called the **structured mean field** approach [SJ95].  
<sup>27</sup>

<sup>28</sup> A common example arises when applying VI to time series models, such as HMMs, where the  
<sup>29</sup> latent variables within each sequence are usually highly correlated across time. Rather than as-  
<sup>30</sup> suming a fully factorized posterior, we can treat each sequence  $\mathbf{z}_{n,1:T}$  as a block, and just assume  
<sup>31</sup> 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  
<sup>32</sup>  $q(\mathbf{z}_{n,1:T}) = \prod_t q(\mathbf{z}_{n,t} | \mathbf{z}_{n,t-1})$ . We can compute the joint distribution  $q(\mathbf{z}_{n,1:T})$ , taking into account  
<sup>33</sup> the dependence between time steps, using the forwards-backwards algorithm. For details, see  
<sup>34</sup> [JW14; Fot+14]. A similar approach was applied to the factorial HMM model, as we discuss in  
<sup>35</sup> Supplementary Section 10.3.2.

<sup>36</sup> An automatic way to derive a structured variational approximation to a probabilistic model,  
<sup>37</sup> specified by a probabilistic programming language, is discussed in [AHG20].  
<sup>38</sup>

### <sup>39</sup> 10.4.2 Hierarchical (auxiliary variable) posteriors

<sup>40</sup>  
<sup>41</sup> Suppose  $q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}) = \prod_k q_{\boldsymbol{\phi}}(z_k|\mathbf{x})$  is a factorized distribution, such as a diagonal Gaussian. This does  
<sup>42</sup> not capture dependencies between the latent variables (components of  $\mathbf{z}$ ). We could of course use a  
<sup>43</sup> full covariance matrix, but this might be too expensive.  
<sup>44</sup>

<sup>45</sup> An alternative approach is to use a hierarchical model, in which we add **auxiliary latent variables**  
<sup>46</sup>  $\mathbf{a}$ , which are used to increase the flexibility of the variational posterior. In particular, we can still  
<sup>47</sup> assume  $q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}, \mathbf{a})$  is conditionally factorized, but when we marginalize out  $\mathbf{a}$ , we induce dependencies

1 between the elements of  $\mathbf{z}$ , i.e.,  
2

3

$$\underline{4} \quad q_{\phi}(\mathbf{z}|\mathbf{x}) = \int q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{a}) q_{\phi}(\mathbf{a}|\mathbf{x}) d\mathbf{a} \neq \prod_k q_{\phi}(z_k|\mathbf{x}) \quad (10.151)$$

5

6 This is called a **hierarchical variational model** [Ran16], or an **auxiliary variable deep gener-  
7 ative model** [Maa+16].

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

### 12 10.4.3 Normalizing flow posteriors

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

### 19 10.4.4 Implicit posteriors

21 In Chapter 26, we discuss implicit probability distributions, which are models which we can sample  
22 from, but which we cannot evaluate pointwise. For example, consider passing a Gaussian noise term,  
23  $\mathbf{z}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , through a nonlinear, *non-invertible* mapping  $f$  to create  $\mathbf{z} = f(\mathbf{z}_0)$ ; it is easy to sample  
24 from  $q(\mathbf{z})$ , but it is intractable to evaluate the density  $q(\mathbf{z})$  (unlike with flows). This makes it hard to  
25 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  
26 can use the same method as is used in GANs (generative adversarial networks, Chapter 26), in which  
27 we train a classifier that discriminates prior samples from samples from the variational posterior by  
28 evaluating  $T(\mathbf{x}, \mathbf{z}) = \log q_{\psi}(\mathbf{z}|\mathbf{x}) - \log p_{\theta}(\mathbf{z})$ . See e.g., [TR19] for details.

### 30 10.4.5 Combining VI with MCMC inference

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

37

$$\nabla_{\mathbf{z}} \log p_{\theta}(\mathbf{z}|\mathbf{x}) = \nabla_{\mathbf{z}} \log [p_{\theta}(\mathbf{x}, \mathbf{z}) - \log p_{\theta}(\mathbf{x})] = \nabla_{\mathbf{z}} \log p_{\theta}(\mathbf{x}, \mathbf{z}) \quad (10.152)$$

38

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

42 A simpler approach is proposed in [Hof17]. Here they train an inference network to initialize an  
43 HMC sampler, using the standard ELBO for  $\phi$ , but they optimize the generative parameters  $\theta$  using  
44 a stochastic approximation to the log marginal likelihood, given by  $\log p_{\theta}(\mathbf{z}, \mathbf{x})$  where  $\mathbf{z}$  is a sample  
45 from the HMC chain. This does not require learning a reverse inference network, and avoids problems  
46 with variational pruning, since it does not use the ELBO for training the generative model.

## 10.5 Tighter bounds

3 Another way to improve the quality of the posterior approximation is to optimize  $q$  wrt a bound that  
4 is a tighter approximation to the log marginal likelihood compared to the standard ELBO. We give  
5 some examples below.  
6

### 10.5.1 Multi-sample ELBO (IWAE bound)

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

13 Let the inference network  $q_\phi(\mathbf{z}|\mathbf{x})$  be viewed as a proposal distribution for the target posterior  
14  $p_\theta(\mathbf{z}|\mathbf{x})$ . 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  $w_s =$   
15  $w_s^*/(\sum_{s'=1}^S w_s^*)$  as the normalized importance weights. From Equation (11.43) we can compute an  
16 estimate of the marginal likelihood  $p(\mathbf{x})$  using  
17

$$\hat{p}_S(\mathbf{x}|\mathbf{z}_{1:S}) \triangleq \frac{1}{S} \sum_{k=1}^S \frac{p_\theta(\mathbf{x}, \mathbf{z}_s)}{q_\phi(\mathbf{z}_s|\mathbf{x})} = \frac{1}{S} \sum_{k=1}^S w_s \quad (10.153)$$

22 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  $q_\phi(\mathbf{z}_{1:S}|\mathbf{x}) = \prod_{k=1}^S q_\phi(\mathbf{z}_s|\mathbf{x})$ . In addition,  
23 since the estimator is always positive, we can take logarithms, and thus obtain a stochastic lower  
24 bound on the log likelihood:  
25

$$\mathbb{L}_S(\phi, \theta|\mathbf{x}) \triangleq \mathbb{E}_{q_\phi(\mathbf{z}_{1:S}|\mathbf{x})} \left[ \log \left( \frac{1}{S} \sum_{s=1}^S w_s \right) \right] = \mathbb{E}_{q_\phi(\mathbf{z}_{1:S}|\mathbf{x})} [\log \hat{p}_S(\mathbf{z}_{1:S})] \quad (10.154)$$

$$\leq \log \mathbb{E}_{q_\phi(\mathbf{z}_{1:S}|\mathbf{x})} [\hat{p}_S(\mathbf{z}_{1:S})] = \log p(\mathbf{x}) \quad (10.155)$$

30 where we used Jensen's inequality in the penultimate line, and the unbiased property in the last line.  
31 This is called the **multi-sample ELBO** or **IWAE bound** [BGS16]. If  $S = 1$ ,  $\mathbb{L}_S$  reduces to the  
32 standard ELBO:  
33

$$\mathbb{L}_1(\phi, \theta|\mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log w] = \int q_\phi(\mathbf{z}|\mathbf{x}) \log \frac{p_\theta(\mathbf{z}, \mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} d\mathbf{z} \quad (10.156)$$

36 One can show [BGS16] that increasing the number of samples  $S$  is guaranteed to make the bound  
37 tighter, thus making it a better proxy for the log likelihood. Intuitively, averaging the  $S$  samples  
38 inside the log removes the need for every sample  $\mathbf{z}_s$  to explain the data  $\mathbf{x}$ . This encourages the  
39 proposal distribution  $q$  to be less concentrated than the single-sample variational posterior.  
40

#### 10.5.1.1 Pathologies of optimizing the IWAE bound

43 Unfortunately, increasing the number of samples in the IWAE bound can decrease the signal to noise  
44 ratio, resulting in learning a worse model [Rai+18a]. Intuitively, the reason this happens is that  
45 increasing  $S$  reduces the dependence of the bound on the quality of the inference network, which  
46 makes the gradient of the ELBO wrt  $\phi$  less informative (higher variance).  
47

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

### 10.5.2 The thermodynamic variational objective (TVO)

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

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

### 10.5.3 Minimizing the evidence upper bound

Recall that the evidence lower bound or ELBO is given by

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi} | \mathbf{x}) = \log p_{\boldsymbol{\theta}}(\mathbf{x}) - D_{\text{KL}}(q_{\boldsymbol{\phi}}(\mathbf{z} | \mathbf{x}) \| p_{\boldsymbol{\theta}}(\mathbf{z} | \mathbf{x})) \leq \log p_{\boldsymbol{\theta}}(\mathbf{x}) \quad (10.157)$$

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

$$\text{EUBO}(\boldsymbol{\theta}, \boldsymbol{\phi} | \mathbf{x}) = \log p_{\boldsymbol{\theta}}(\mathbf{x}) + D_{\text{KL}}(p_{\boldsymbol{\theta}}(\mathbf{z} | \mathbf{x}) \| q_{\boldsymbol{\phi}}(\mathbf{z} | \mathbf{x})) \geq \log p_{\boldsymbol{\theta}}(\mathbf{x}) \quad (10.158)$$

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

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

Now suppose we sample  $\mathbf{x}$  from the empirical distribution, and minimize  $\mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [\text{EUBO}(\boldsymbol{\theta}, \boldsymbol{\phi} | \mathbf{x})]$  wrt  $\boldsymbol{\phi}$ . To approximate the expectation, we can use self-normalized importance sampling, as in Equation (10.173), to get

$$\nabla_{\boldsymbol{\phi}} \text{EUBO}(\boldsymbol{\theta}, \boldsymbol{\phi} | \mathbf{x}) = \sum_{s=1}^S \bar{w}_s \nabla_{\boldsymbol{\phi}} \log q_{\boldsymbol{\phi}}(\mathbf{z}^s | \mathbf{x}) \quad (10.159)$$

where  $\bar{w}_s = w^{(s)} / (\sum_s w^{(s)})$ , and  $w^{(s)} = \frac{p(\mathbf{x}, \mathbf{z}^s)}{q(\mathbf{z}^s | \boldsymbol{\phi}_t)}$ . This is equivalent to the “daydream” update (aka “wake-phase  $\boldsymbol{\phi}$  update”) of the wake-sleep algorithm (see Section 10.6.3).

## 1 2 10.6 Wake-sleep algorithm

3 So far in this chapter we have focused on fitting latent variable models by maximizing the ELBO.  
4 This has two main drawbacks. First, it does not work well when we have discrete latent variables,  
5 because in such cases we cannot use the reparameterization trick; thus we have to use higher variance  
6 estimators, such as REINFORCE (see Section 10.3.2). Second, even in the case where we can use  
7 the reparameterization trick, the lower bound may not be very tight. We can improve the tightness  
8 by using the IWAE multi-sample bound (Section 10.5.1), but paradoxically this may not result in  
9 learning a better model, for reasons discussed in Section 10.5.1.1.

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

### 21 22 10.6.1 Wake phase

23 In the **wake phase**, we minimize the KL divergence from the empirical distribution to the model's  
24 distribution:

$$\mathcal{L}(\theta) = D_{\text{KL}}(p_{\mathcal{D}}(\mathbf{x}) \parallel p_{\theta}(\mathbf{x})) = \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})}[-\log p_{\theta}(\mathbf{x})] + \text{const} \quad (10.160)$$

25 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  
26 data:

$$\ell(\theta) = \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})}[\log p_{\theta}(\mathbf{x})] \quad (10.161)$$

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

$$\ell(\theta|\phi, \mathbf{x}) = \log \left( \frac{1}{S} \sum_{s=1}^S w_s \right) \quad (10.162)$$

32 where  $w_s = \frac{p_{\theta}(\mathbf{x}, \mathbf{z}_s)}{q_{\phi}(\mathbf{z}_s|\mathbf{x})}$ .

33 We now discuss how to compute the gradient of this objective. Using the log-derivative trick, we  
34 have that

$$\nabla_{\theta} \log w_s = \frac{1}{w_s} \nabla_{\theta} w_s = \nabla_{\theta} \log p_{\theta}(\mathbf{x}, \mathbf{z}_s) \quad (10.163)$$

1  
2 Hence

3  
4  $\nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta} | \boldsymbol{\phi}, \mathbf{x}) = \frac{1}{\sum_{s=1}^S w_s} \left( \frac{1}{S} \sum_{s=1}^S \nabla_{\boldsymbol{\theta}} w_s \right)$  (10.164)  
5

6  
7  $= \frac{1}{\sum_{s=1}^S w_s} \left( \sum_{s=1}^S \frac{p_{\boldsymbol{\theta}}(\mathbf{z}_s, \mathbf{x})}{q_{\boldsymbol{\phi}}(\mathbf{z}_s | \mathbf{x})} \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{z}_s, \mathbf{x}) \right)$  (10.165)  
8

9  
10  $= \sum_{s=1}^S \bar{w}_s \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{z}_s, \mathbf{x})$  (10.166)  
11

12  
13 where  $\bar{w}_s = w_s / (\sum_{s'=1}^S w_{s'})$ .

### 14 10.6.2 Sleep phase

15 In the **sleep phase**, we try to minimize the KL divergence between the true posterior (under the  
16 current model) and the inference network’s approximation to that posterior:

17  
18  $\mathcal{L}(\boldsymbol{\phi}) = \mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{x})} [D_{\text{KL}}(p_{\boldsymbol{\theta}}(\mathbf{z} | \mathbf{x}) \| q_{\boldsymbol{\phi}}(\mathbf{z} | \mathbf{x}))] = \mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{z}, \mathbf{x})} [-\log q_{\boldsymbol{\phi}}(\mathbf{z} | \mathbf{x})] + \text{const}$  (10.167)  
19

20 Equivalently, we can maximize the following loglikelihood objective:

21  
22  $\ell(\boldsymbol{\phi} | \boldsymbol{\theta}) = \mathbb{E}_{(\mathbf{z}, \mathbf{x}) \sim p_{\boldsymbol{\theta}}(\mathbf{z}, \mathbf{x})} [\log q_{\boldsymbol{\phi}}(\mathbf{z} | \mathbf{x})]$  (10.168)  
23

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

28  
29  $\ell(\boldsymbol{\phi} | \boldsymbol{\theta}) = \frac{1}{S} \sum_{s=1}^S \log q_{\boldsymbol{\phi}}(\mathbf{z}'_s | \mathbf{x}'_s)$  (10.169)  
30

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

32  
33  $\nabla_{\boldsymbol{\phi}} \ell(\boldsymbol{\phi} | \boldsymbol{\theta}) = \frac{1}{S} \sum_{s=1}^S \nabla_{\boldsymbol{\phi}} \log q_{\boldsymbol{\phi}}(\mathbf{z}'_s | \mathbf{x}'_s)$  (10.170)  
34

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

### 41 10.6.3 Daydream phase

42 The disadvantage of the sleep phase is that the inference network,  $q_{\boldsymbol{\phi}}(\mathbf{z} | \mathbf{x})$ , is trying to follow a  
43 moving target,  $p_{\boldsymbol{\theta}}(\mathbf{z} | \mathbf{x})$ . Furthermore, it is only being trained on synthetic data from the model,  
44 not on real data. The reweighted wake-sleep algorithm of [BB15b] proposed to learn the inference  
45

<sup>1</sup> network by using real data from the empirical distribution, in addition to fantasy data. They call  
<sup>2</sup> the case where you use real data the “**wake-phase  $q$  update**”, but we will call it the “**daydream**  
<sup>3</sup> **phase**”, since, unlike sleeping, the system uses real data  $\mathbf{x}$  to update the inference model, instead of  
<sup>4</sup> fantasies.<sup>4</sup> [Le+19] went further, and proposed to only use the wake and daydream phases, and to  
<sup>5</sup> skip the sleep phase entirely.

<sup>6</sup> In more detail, the new objective which we want to minimize becomes

$$\mathcal{L}(\phi|\theta) = \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [D_{\text{KL}}(p_{\theta}(\mathbf{z}|\mathbf{x}) \| q_{\phi}(\mathbf{z}|\mathbf{x}))] \quad (10.171)$$

<sup>10</sup> We can compute a single sample approximation to the negative of the above expression as follows:

$$\ell(\phi|\theta, \mathbf{x}) = \mathbb{E}_{p_{\theta}(\mathbf{z}|\mathbf{x})} [\log q_{\phi}(\mathbf{z}|\mathbf{x})] \quad (10.172)$$

<sup>14</sup> where  $\mathbf{x} \sim p_{\mathcal{D}}$ . We can approximate this expectation using importance sampling, with  $q_{\phi}$  as the  
<sup>15</sup> proposal. This results in the following estimator of the gradient for each datapoint:

$$\nabla_{\phi} \ell(\phi|\theta, \mathbf{x}) = \int p_{\theta}(\mathbf{z}|\mathbf{x}) \nabla_{\phi} \log q_{\phi}(\mathbf{z}|\mathbf{x}) d\mathbf{z} \approx \sum_{s=1}^S \bar{w}_s \nabla_{\phi} \log q_{\phi}(z_s|\mathbf{x}) \quad (10.173)$$

<sup>20</sup> where  $\mathbf{z}_s \sim q_{\phi}(\mathbf{z}_s|\mathbf{x})$  and  $\bar{w}_s$  are the normalized weights.

<sup>21</sup> We see that Equation (10.173) is very similar to Equation (10.170). The key difference is that in  
<sup>22</sup> 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 data point, whereas  
<sup>23</sup> 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 data point.

#### <sup>25</sup> 10.6.4 Summary of algorithm

---

<sup>28</sup> **Algorithm 24:** One SGD update using wake-sleep algorithm.

---

- <sup>29</sup> 1 Sample  $\mathbf{x}_n$  from dataset
  - <sup>30</sup> 2 Draw  $S$  samples from inference network:  $\mathbf{z}_s \sim q(\mathbf{z}|\mathbf{x}_n)$
  - <sup>31</sup> 3 Compute unnormalized weights:  $w_s = \frac{p(\mathbf{x}_n, \mathbf{z}_s)}{q(\mathbf{z}_s|\mathbf{x}_n)}$
  - <sup>32</sup> 4 Compute normalized weights:  $\bar{w}_s = \frac{w_s}{\sum_{s'=1}^S w_{s'}}$
  - <sup>34</sup> 5 Optional: Compute estimate of log likelihood:  $\log p(\mathbf{x}_n) = \log(\frac{1}{S} \sum_{s=1}^S w_s)$
  - <sup>35</sup> 6 Wake phase: Update  $\theta$  using  $\sum_{s=1}^S \bar{w}_s \nabla_{\theta} \log p_{\theta}(\mathbf{z}_s, \mathbf{x}_n)$
  - <sup>36</sup> 7 Daydream phase: Update  $\phi$  using  $\sum_{s=1}^S \bar{w}_s \nabla_{\phi} \log q_{\phi}(\mathbf{z}_s|\mathbf{x}_n)$
  - <sup>37</sup> 8 Optional sleep phase: Draw  $S$  samples from model,  $(\mathbf{x}'_s, \mathbf{z}'_s) \sim p_{\theta}(\mathbf{x}, \mathbf{z})$  and update  $\phi$  using  
<sup>38</sup>  $\frac{1}{S} \sum_{s=1}^S \nabla_{\phi} \log q_{\phi}(\mathbf{z}'_s|\mathbf{x}'_s)$
- 

<sup>41</sup> We summarize the RWS algorithm in Algorithm 24. The disadvantage of the RWS algorithm is  
<sup>42</sup> that it does not optimize a single well-defined objective, so it is not clear if the method will converge,  
<sup>43</sup> in contrast to ELBO maximization. On the other hand, the method is fairly simple, since it consists  
<sup>44</sup> of two alternating weighted maximum likelihood problems. It can also be shown to “sandwich” a  
<sup>45</sup>

<sup>46</sup> 4. We thank Rif Saurous for suggesting this term.

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

$$\text{wake phase } \min_{\theta} D_{\text{KL}}(q_{\mathcal{D}, \phi}(\mathbf{x}, \mathbf{z}) \| p_{\theta}(\mathbf{x}, \mathbf{z})) \quad (10.174)$$

$$\text{daydream phase } \min_{\phi} D_{\text{KL}}(p_{\theta}(\mathbf{x}, \mathbf{z}) \| q_{\mathcal{D}, \phi}(\mathbf{x}, \mathbf{z})) \quad (10.175)$$

## 10.7 Expectation propagation (EP)

One problem with lower bound maximization (i.e., standard VI) is that we are minimizing  $D_{\text{KL}}(q \| p)$ , which induces **zero-forcing** behavior, as we discussed in Section 5.1.3.3. 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  $D_{\text{KL}}(p \| q)$ , which is zero-avoiding, as we discussed in Section 5.1.3.3. 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  $D_{\text{KL}}(p \| q)$ .

### 10.7.1 Algorithm

We assume the exact posterior can be written as follows:

$$p(\boldsymbol{\theta}|\mathcal{D}) = \frac{1}{Z_p} \hat{p}(\boldsymbol{\theta}), \quad \hat{p}(\boldsymbol{\theta}) = p_0(\boldsymbol{\theta}) \prod_{k=1}^K f_k(\boldsymbol{\theta}) \quad (10.176)$$

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

$$q(\boldsymbol{\theta}) = \frac{1}{Z_q} \hat{q}(\boldsymbol{\theta}), \quad \hat{q}(\boldsymbol{\theta}) = p_0(\boldsymbol{\theta}) \prod_{k=1}^K \tilde{f}_k(\boldsymbol{\theta}) \quad (10.177)$$

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

$$q_{-i}^{\text{cavity}}(\boldsymbol{\theta}) = \frac{q(\boldsymbol{\theta})}{\tilde{f}_i(\boldsymbol{\theta})} \propto \prod_{k \neq i} \tilde{f}_k(\boldsymbol{\theta}) \quad (10.178)$$

This division operation can be implemented by subtracting the natural parameters, as explained in Section 2.2.7.2. 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:

$$q_i^{\text{tilted}}(\boldsymbol{\theta}) = \frac{1}{Z_i} f_i(\boldsymbol{\theta}) q_{-i}^{\text{cavity}}(\boldsymbol{\theta}) \quad (10.179)$$

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:

$$q_i^{\text{proj}}(\boldsymbol{\theta}) = \text{proj}(q_i^{\text{tilted}}) \triangleq \underset{\tilde{q} \in \mathcal{Q}}{\operatorname{argmin}} D(q_i^{\text{tilted}} \parallel \tilde{q}) \quad (10.180)$$

This can be thought of as projecting the tilted distribution into the approximation family. If  $D(q_i^{\text{tilted}} \parallel q) = D_{\text{KL}}(q_i^{\text{tilted}} \parallel q)$ , this can be done by moment matching, as shown in Section 5.1.3.4.

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

$$\log q_i^{\text{tilted}}(\boldsymbol{\theta}) = \alpha \log f_i(\boldsymbol{\theta}) - \frac{1}{2} \boldsymbol{\theta}^\top \mathbf{Q}_{-i} \boldsymbol{\theta} + \mathbf{r}_{-i}^\top \boldsymbol{\theta} + \text{const} \quad (10.181)$$

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:

$$\mathbf{Q}_{\setminus i} = -\nabla_{\boldsymbol{\theta}}^2 \log q_i^{\text{tilted}}(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}}, \quad \mathbf{r}_{\setminus i} = \mathbf{Q}_{\setminus i} \hat{\boldsymbol{\theta}} \quad (10.182)$$

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

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

$$\tilde{f}_i^{\text{new}}(\boldsymbol{\theta}) = \frac{q_i^{\text{proj}}(\boldsymbol{\theta})}{q_{-i}^{\text{cavity}}(\boldsymbol{\theta})} \quad (10.183)$$

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

41

## 10.7.2 Example

Figure 10.13 illustrates the process of combining a very non-Gaussian likelihood  $f_i$  with a Gaussian cavity prior  $q_{-i}^{\text{cavity}}$  to yield a nearly Gaussian tilted distribution  $q_i^{\text{tilted}}$ , which can then be approximated by a Gaussian using projection.

47

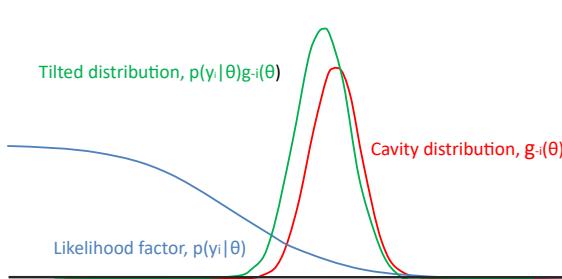


Figure 10.13: Combining a logistic likelihood factor  $f_i = p(y_i | \theta)$  with the cavity prior,  $q_{-i}^{cavity} = g_{-i}(\theta)$ , to get the tilted distribution,  $q_i^{tilted} = p(y_i | \theta)g_{-i}(\theta)$ . Adapted from Figure 2 of [Gel+14b].

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

### 10.7.3 EP as generalized ADF

We can view EP as a generalization of the ADF algorithm discussed in Section 8.10. ADF is a form of sequential Bayesian inference. At each step, it maintains a tractable approximation to the posterior,  $q_t(\mathbf{z}) \in \mathcal{Q}$ , updates it with the likelihood from the next observation,  $\hat{p}_{t+1}(\mathbf{z}) \propto q_t(\mathbf{z})p(\mathbf{x}_t | \mathbf{z})$ , and then projects the resulting updated posterior back to the tractable family using  $q_{t+1} = \operatorname{argmin}_{q \in \mathcal{Q}} D_{\text{KL}}(\hat{p}_{t+1} \| q)$ . ADF minimizes KL in the desired direction. However, it is a sequential algorithm, designed for the online setting. In the batch setting, the method can give 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  $\delta$ . More precisely, we change the final step to be the following:

$$\tilde{f}_i^{\text{new}}(\theta) = \left( \tilde{f}_i(\theta) \right)^{1-\delta} \left( \frac{q_i^{\text{proj}}(\theta)}{q_{-i}^{\text{cavity}}} \right)^\delta \quad (10.184)$$

This can be implemented by scaling the natural parameters by  $\delta$ . [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.

1 In addition to numerical stability, there is no guarantee that EP will converge in its vanilla form,  
2 although empirically it can work well, especially with log-concave factors  $f_i$  (e.g., as in GP classifiers).  
3

### 4 5 10.7.5 Power EP and $\alpha$ -divergence

6 We also have a choice about what divergence measure  $D(q_i^{\text{tilted}} \parallel q)$  to use when we approximate  
7 the tilted distribution. If we use  $D_{\text{KL}}(q_i^{\text{tilted}} \parallel q)$ , we recover classic EP, as described above. If  
8 we use  $D_{\text{KL}}(q \parallel q_i^{\text{tilted}})$ , we recover the reverse KL used in standard variational inference. We can  
9 generalize the above results by using  $\alpha$ -divergences (Section 2.7.1.2), which allow us to interpolate  
10 between mode seeking and mode covering behavior, as shown in Figure 2.19. We can optimize the  
11  $\alpha$ -divergence by using the **power EP** method of [Min04].  
12

13 Algorithmically, this is a fairly small modification to regular EP. In particular, we first compute the  
14 cavity distribution,  $q_{-i}^{\text{cavity}} \propto \frac{q}{f_i^\alpha}$ ; we then approximate the tilted distribution,  $q_i^{\text{proj}} = \text{proj}(q_{-i}^{\text{cavity}} f_i^\alpha)$ ;  
15 and finally we compute the new factor  $\tilde{f}_i^{\text{new}} \propto \left(\frac{q_i^{\text{proj}}}{q_{-i}^{\text{cavity}}}\right)^{1/\alpha}$ .  
16

### 17 18 10.7.6 Stochastic EP

19 The main disadvantage of EP in the big data setting is that we need to store the  $\tilde{f}_n(\boldsymbol{\theta})$  terms for  
20 each datapoint  $n$ , so we can compute the cavity distribution. If  $\boldsymbol{\theta}$  has  $D$  dimensions, and we use full  
21 covariance Gaussians, this requires  $O(ND^2)$  memory.  
22

23 The idea behind **stochastic EP** [HLHT15] is to approximate the local factors with a shared factor  
24 that acts like an aggregated likelihood, i.e.,  
25

$$\prod_{n=1}^N f_n(\boldsymbol{\theta}) \approx \tilde{f}(\boldsymbol{\theta})^N \quad (10.185)$$

26 where typically  $f_n(\boldsymbol{\theta}) = p(\mathbf{x}_n | \boldsymbol{\theta})$ . This exploits the fact that the posterior only cares about approxi-  
27 mating the product of the likelihoods, rather than each likelihood separately. Hence it suffices for  
28  $\tilde{f}(\boldsymbol{\theta})$  to approximate the average likelihood.  
29

30 We can modify EP to this setting as follows. First, when computing the cavity distribution, we  
31 use  
32

$$q_{-1}(\boldsymbol{\theta}) \propto q(\boldsymbol{\theta}) / \tilde{f}(\boldsymbol{\theta}) \quad (10.186)$$

33 We then compute the tilted distribution  
34

$$q_{\setminus n}(\boldsymbol{\theta}) \propto f_n(\boldsymbol{\theta}) q_{-1}(\boldsymbol{\theta}) \quad (10.187)$$

35 Next we derive the new local factor for this datapoint using moment matching:  
36

$$\tilde{f}_n(\boldsymbol{\theta}) = \text{proj}(q_{\setminus n}(\boldsymbol{\theta})) / q_{-1}(\boldsymbol{\theta}) \quad (10.188)$$

37 Finally we perform a damped update of the average likelihood  $\tilde{f}(\boldsymbol{\theta})$  using this new local factor:  
38

$$\tilde{f}_{\text{new}}(\boldsymbol{\theta}) = \tilde{f}_{\text{old}}(\boldsymbol{\theta})^{1-1/N} \tilde{f}_n(\boldsymbol{\theta})^{1/N} \quad (10.189)$$

1      The ADF algorithm is similar to SEP, in that we compute the tilted distribution  $q_{\setminus t} \propto f_t q_{t-1}$  and  
2      then project it, without needing to keep the  $f_t$  factors. The difference is that instead of using the  
3      cavity distribution  $q_{-1}(\theta)$  as prior, it uses the posterior from the previous time step,  $q_{t-1}$ . This  
4      avoids the need to compute and store  $\tilde{f}$ , but results in overconfidence in the batch setting.  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
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 inference

## 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:

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

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

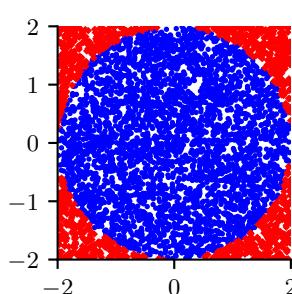
An alternative approach is to draw multiple random samples,  $\mathbf{x}_n \sim p(\mathbf{x})$ , and then to compute

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

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].



11 *Figure 11.1: Estimating  $\pi$  by Monte Carlo integration using 5000 samples. Blue points are inside the circle,*  
12 *red points are outside. Generated by `mc_estimate_pi.ipynb`.*

13

14  
15 need to uniformly cover the entire space. In particular, it can be shown that the accuracy is in  
16 principle independent of the dimensionality of  $\mathbf{x}$ , and only depends on the number of samples  $N_s$   
17 (see Section 11.2.2 for details). The catch is that we need a way to generate the samples  $\mathbf{x}_n \sim p(\mathbf{x})$   
18 in the first place. In addition, the estimator may have high variance. We will discuss this topic at  
19 length in the sections below.

20

### 21 11.2.1 Example: estimating $\pi$ by Monte Carlo integration

22

23 MC integration can be used for many applications, not just in ML and statistics. For example,  
24 suppose we want to estimate  $\pi$ . We know that the area of a circle with radius  $r$  is  $\pi r^2$ , but it is also  
25 equal to the following definite integral:

$$\underline{26} \quad I = \int_{-r}^r \int_{-r}^r \mathbb{I}(x^2 + y^2 \leq r^2) dx dy \quad (11.3)$$

27 Hence  $\pi = I/(r^2)$ . Let us approximate this by Monte Carlo integration. Let  $f(x, y) = \mathbb{I}(x^2 + y^2 \leq r^2)$   
28 be an indicator function that is 1 for points inside the circle, and 0 outside, and let  $p(x)$  and  $p(y)$  be  
29 uniform distributions on  $[-r, r]$ , so  $p(x) = p(y) = 1/(2r)$ . Then

$$\underline{30} \quad I = (2r)(2r) \int \int f(x, y) p(x) p(y) dx dy \quad (11.4)$$

$$\underline{31} \quad = 4r^2 \int \int f(x, y) p(x) p(y) dx dy \quad (11.5)$$

$$\underline{32} \quad \approx 4r^2 \frac{1}{N_s} \sum_{n=1}^{N_s} f(x_n, y_n) \quad (11.6)$$

33 Using 5000 samples, we find  $\hat{\pi} = 3.10$  with standard error 0.09 compared to the true value of  $\pi = 3.14$ .

34 We can plot the points that are accepted or rejected as in Figure 11.1.

35

### 36 11.2.2 Accuracy of Monte Carlo integration

37

38 The accuracy of an MC approximation increases with sample size. This is illustrated in Figure 11.2.

39 On the top line, we plot a histogram of samples from a Gaussian distribution. On the bottom line,

40

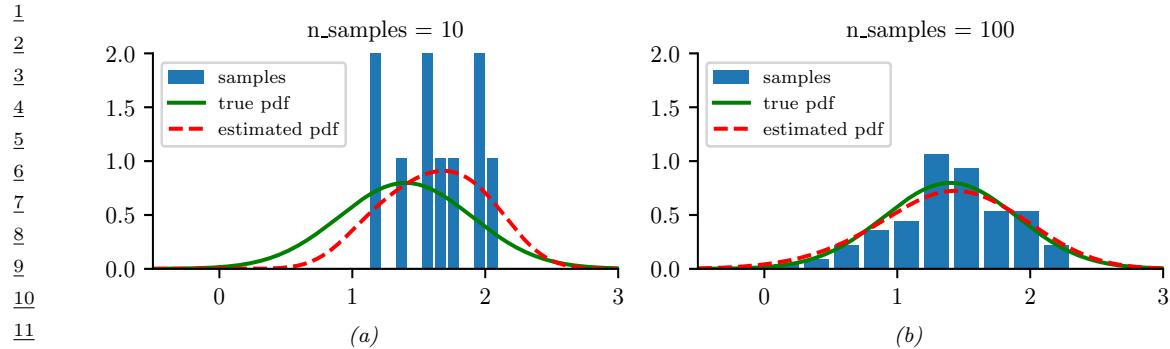


Figure 11.2: 10 and 100 samples from a Gaussian distribution,  $\mathcal{N}(\mu = 1.5, \sigma^2 = 0.25)$ . A dotted red line denotes kernel density estimate derived from the samples. Generated by [mc\\_accuracy\\_demo.ipynb](#).

we plot a smoothed version of these samples, created using a kernel density estimate. This smoothed distribution is then evaluated on a dense grid of points and plotted. Note that this smoothing is just for the purposes of plotting, it is not used for the Monte Carlo estimate itself.

If we denote the exact mean by  $\mu = \mathbb{E}[f(X)]$ , and the MC approximation by  $\hat{\mu}$ , one can show that, with independent samples,

$$(\hat{\mu} - \mu) \rightarrow \mathcal{N}\left(0, \frac{\sigma^2}{N_s}\right) \quad (11.7)$$

where

$$\sigma^2 = \mathbb{V}[f(X)] = \mathbb{E}[f(X)^2] - \mathbb{E}[f(X)]^2 \quad (11.8)$$

This is a consequence of the central limit theorem. Of course,  $\sigma^2$  is unknown in the above expression, but it can be estimated by MC:

$$\hat{\sigma}^2 = \frac{1}{N_s} \sum_{n=1}^{N_s} (f(x_n) - \hat{\mu})^2 \quad (11.9)$$

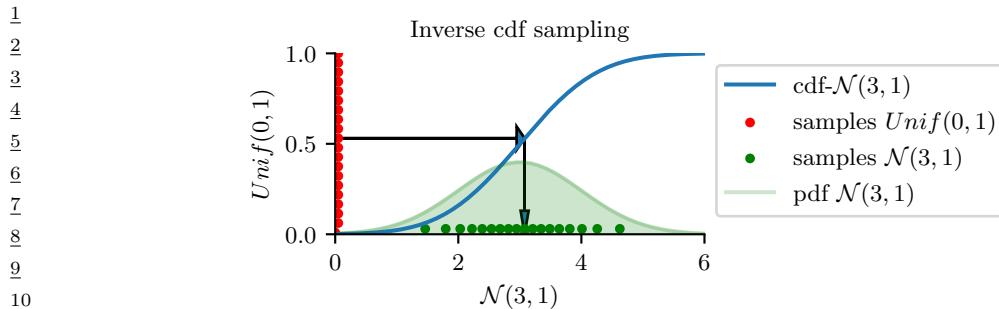
Thus for large enough  $N_s$  we have

$$P\left\{\hat{\mu} - 1.96 \frac{\hat{\sigma}}{\sqrt{N_s}} \leq \mu \leq \hat{\mu} + 1.96 \frac{\hat{\sigma}}{\sqrt{N_s}}\right\} \approx 0.95 \quad (11.10)$$

The term  $\sqrt{\hat{\sigma}^2/N_s}$  is called the (numerical or empirical) **standard error**, and is an estimate of our uncertainty about our estimate of  $\mu$ .

If we want to report an answer which is accurate to within  $\pm\epsilon$  with probability at least 95%, we need to use a number of samples  $N_s$  which satisfies  $1.96\sqrt{\hat{\sigma}^2/N_s} \leq \epsilon$ . We can approximate the 1.96 factor by 2, yielding  $N_s \geq \frac{4\hat{\sigma}^2}{\epsilon^2}$ .

The remarkable thing to note about the above results is that the error in the estimate,  $\sigma^2/N_s$ , is theoretically independent of the dimensionality of the integral. The catch is that sampling from high dimensional distributions can be hard. We turn to that topic next.

Figure 11.3: Sampling from  $\mathcal{N}(3, 1)$  using an inverse CDF.

### 11.3 Generating random samples from simple distributions

We saw in Section 11.2 how we can evaluate  $\mathbb{E}[f(X)]$  for different functions  $f$  of a random variable  $X$  using Monte Carlo integration. The main computational challenge is to efficiently generate samples from the probability distribution  $p^*(\mathbf{x})$  (which may be a posterior,  $p^*(\mathbf{x}) \propto p(\mathbf{x}|\mathcal{D})$ ). In this section, we discuss sampling methods that are suitable for parametric univariate distributions. These can be used as building blocks for sampling from more complex multivariate distributions.

#### 11.3.1 Sampling using the inverse cdf

The simplest method for sampling from a univariate distribution is based on the **inverse probability transform**. Let  $F$  be a cdf of some distribution we want to sample from, and let  $F^{-1}$  be its inverse. Then we have the following result.

**Theorem 11.3.1.** If  $U \sim U(0, 1)$  is a uniform rv, then  $F^{-1}(U) \sim F$ .

*Proof.*

$$\Pr(F^{-1}(U) \leq x) = \Pr(U \leq F(x)) \quad (\text{applying } F \text{ to both sides}) \quad (11.11)$$

$$= F(x) \quad (\text{because } \Pr(U \leq y) = y) \quad (11.12)$$

where the first line follows since  $F$  is a monotonic function, and the second line follows since  $U$  is uniform on the unit interval.  $\square$

Hence we can sample from any univariate distribution, for which we can evaluate its inverse cdf, as follows: generate a random number  $u \sim U(0, 1)$  using a **pseudo random number generator** (see e.g., [Pre+88] for details). Let  $u$  represent the height up the  $y$  axis. Then “slide along” the  $x$  axis until you intersect the  $F$  curve, and then “drop down” and return the corresponding  $x$  value. This corresponds to computing  $x = F^{-1}(u)$ . See Figure 11.3 for an illustration.

For example, consider the exponential distribution

$$\text{Expon}(x|\lambda) \triangleq \lambda e^{-\lambda x} \mathbb{I}(x \geq 0) \quad (11.13)$$

The cdf is

$$F(x) = 1 - e^{-\lambda x} \mathbb{I}(x \geq 0) \quad (11.14)$$

whose inverse is the quantile function

$$F^{-1}(p) = -\frac{\ln(1-p)}{\lambda} \quad (11.15)$$

By the above theorem, if  $U \sim \text{Unif}(0, 1)$ , we know that  $F^{-1}(U) \sim \text{Expon}(\lambda)$ . So we can sample from the exponential distribution by first sampling from the uniform and then transforming the results using  $-\ln(1-u)/\lambda$ . (In fact, since  $1-U \sim \text{Unif}(0, 1)$ , we can just use  $-\ln(u)/\lambda$ .)

### 11.3.2 Sampling from a Gaussian (Box-Muller method)

In this section, we describe a method to sample from a Gaussian. The idea is we sample uniformly from a unit radius circle, and then use the change of variables formula to derive samples from a spherical 2d Gaussian. This can be thought of as two samples from a 1d Gaussian.

In more detail, sample  $z_1, z_2 \in (-1, 1)$  uniformly, and then discard pairs that do not satisfy  $z_1^2 + z_2^2 \leq 1$ . The result will be points uniformly distributed inside the unit circle, so  $p(\mathbf{z}) = \frac{1}{\pi} \mathbb{I}(z \text{ inside circle})$ .

Now define

$$x_i = z_i \left( \frac{-2 \ln r^2}{r^2} \right)^{\frac{1}{2}} \quad (11.16)$$

for  $i = 1 : 2$ , where  $r^2 = z_1^2 + z_2^2$ . Using the multivariate change of variables formula, we have

$$p(x_1, x_2) = p(z_1, z_2) \left| \frac{\partial(z_1, z_2)}{\partial(x_1, x_2)} \right| = \left[ \frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}x_1^2) \right] \left[ \frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}x_2^2) \right] \quad (11.17)$$

Hence  $x_1$  and  $x_2$  are two independent samples from a univariate Gaussian. This is known as the **Box-Muller** method.

To sample from a multivariate Gaussian, we first compute the Cholesky decomposition of its covariance matrix,  $\Sigma = \mathbf{L}\mathbf{L}^\top$ , where  $\mathbf{L}$  is lower triangular. Next we sample  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  using the Box-Muller method. Finally we set  $\mathbf{y} = \mathbf{L}\mathbf{x} + \boldsymbol{\mu}$ . This is valid since

$$\text{Cov}[\mathbf{y}] = \mathbf{LCov}[\mathbf{x}]\mathbf{L}^\top = \mathbf{L} \mathbf{I} \mathbf{L}^\top = \Sigma \quad (11.18)$$

## 11.4 Rejection sampling

Suppose we want to sample from the **target distribution**

$$p(\mathbf{x}) = \tilde{p}(\mathbf{x})/Z_p \quad (11.19)$$

where  $\tilde{p}(\mathbf{x})$  is the unnormalized version, and

$$Z_p = \int \tilde{p}(\mathbf{x}) d\mathbf{x} \quad (11.20)$$

is the (possibly unknown) normalization constant. One of the simplest approaches to this problem is **rejection sampling**, which we now explain.

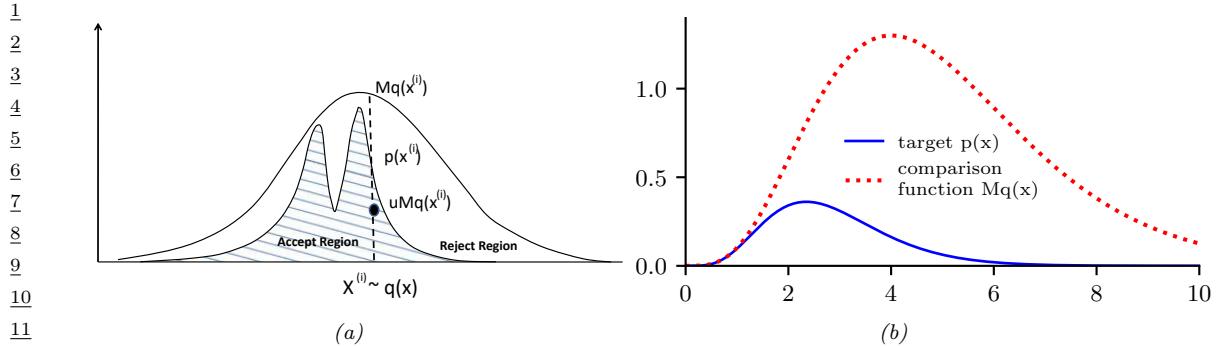


Figure 11.4: (a) Schematic illustration of rejection sampling. From Figure 2 of [And+03]. Used with kind permission of Nando de Freitas. (b) Rejection sampling from a  $\text{Ga}(\alpha = 5.7, \lambda = 2)$  distribution (solid blue) using a proposal of the form  $\text{MGa}(k, \lambda - 1)$  (dotted red), where  $k = \lfloor 5.7 \rfloor = 5$ . The curves touch at  $\alpha - k = 0.7$ . Generated by [rejection\\_sampling\\_demo.ipynb](#).

### 11.4.1 Basic idea

In rejection sampling, we require access to a **proposal distribution**  $q(\mathbf{x})$  which satisfies  $Cq(\mathbf{x}) \geq \tilde{p}(\mathbf{x})$ , for some constant  $C$ . The function  $Cq(\mathbf{x})$  provides an upper envelope for  $\tilde{p}$ .

We can use the proposal distribution to generate samples from the target distribution as follows. We first sample  $\mathbf{x}_0 \sim q(\mathbf{x})$ , which corresponds to picking a random  $\mathbf{x}$  location, and then we sample  $u_0 \sim \text{Unif}(0, Cq(\mathbf{x}_0))$ , which corresponds to picking a random height ( $y$  location) under the envelope. If  $u_0 > \tilde{p}(\mathbf{x}_0)$ , we reject the sample, otherwise we accept it. This process is illustrated in 1d in Figure 11.4(a): the acceptance region is shown shaded, and the rejection region is the white region between the shaded zone and the upper envelope.

We now prove this procedure is correct. First note that the probability of any given sample  $\mathbf{x}_0$  being accepted equals the probability of a sample  $u_0 \sim \text{Unif}(0, Cq(\mathbf{x}_0))$  being less than or equal to  $\tilde{p}(\mathbf{x}_0)$ , i.e.,

$$q(\text{accept}|\mathbf{x}_0) = \int_0^{\tilde{p}(\mathbf{x}_0)} \frac{1}{Cq(\mathbf{x}_0)} du = \frac{\tilde{p}(\mathbf{x}_0)}{Cq(\mathbf{x}_0)} \quad (11.21)$$

Therefore

$$q(\text{propose and accept } \mathbf{x}_0) = q(\mathbf{x}_0)q(\text{accept}|\mathbf{x}_0) = q(\mathbf{x}_0) \frac{\tilde{p}(\mathbf{x}_0)}{Cq(\mathbf{x}_0)} = \frac{\tilde{p}(\mathbf{x}_0)}{C} \quad (11.22)$$

Integrating both sides give

$$\int q(\mathbf{x}_0)q(\text{accept}|\mathbf{x}_0) d\mathbf{x}_0 = q(\text{accept}) = \frac{\int \tilde{p}(\mathbf{x}_0) d\mathbf{x}_0}{C} = \frac{Z_p}{C} \quad (11.23)$$

Hence we see that the distribution of accepted points is given by the target distribution:

$$q(\mathbf{x}_0|\text{accept}) = \frac{q(\mathbf{x}_0, \text{accept})}{q(\text{accept})} = \frac{\tilde{p}(\mathbf{x}_0)}{C} \frac{C}{Z_p} = \frac{\tilde{p}(\mathbf{x}_0)}{Z_p} = p(\mathbf{x}_0) \quad (11.24)$$

How efficient is this method? If  $\tilde{p}$  is a normalized target distribution, the acceptance probability is  $1/C$ . Hence we want to choose  $C$  as small as possible while still satisfying  $Cq(x) \geq \tilde{p}(x)$ .

### 11.4.2 Example

For example, suppose we want to sample from a Gamma distribution:<sup>3</sup>

$$\text{Ga}(x|\alpha, \lambda) = \frac{1}{\Gamma(\alpha)} x^{\alpha-1} \lambda^\alpha \exp(-\lambda x) \quad (11.25)$$

where  $\Gamma(\alpha)$  is the gamma function. One can show that if  $X_i \stackrel{iid}{\sim} \text{Expon}(\lambda)$ , and  $Y = X_1 + \dots + X_k$ , then  $Y \sim \text{Ga}(k, \lambda)$ . For non-integer shape parameters  $\alpha$ , we cannot use this trick. However, we can use rejection sampling using a  $\text{Ga}(k, \lambda - 1)$  distribution as a proposal, where  $k = \lfloor \alpha \rfloor$ . The ratio has the form

$$\frac{p(x)}{q(x)} = \frac{\text{Ga}(x|\alpha, \lambda)}{\text{Ga}(x|k, \lambda - 1)} = \frac{x^{\alpha-1} \lambda^\alpha \exp(-\lambda x)/\Gamma(\alpha)}{x^{k-1} (\lambda - 1)^k \exp(-(\lambda - 1)x)/\Gamma(k)} \quad (11.26)$$

$$= \frac{\Gamma(k) \lambda^\alpha}{\Gamma(\alpha) (\lambda - 1)^k} x^{\alpha-k} \exp(-x) \quad (11.27)$$

This ratio attains its maximum when  $x = \alpha - k$ . Hence

$$C = \frac{\text{Ga}(\alpha - k|\alpha, \lambda)}{\text{Ga}(\alpha - k|k, \lambda - 1)} \quad (11.28)$$

See Figure 11.4(b) for a plot.

### 11.4.3 Adaptive rejection sampling

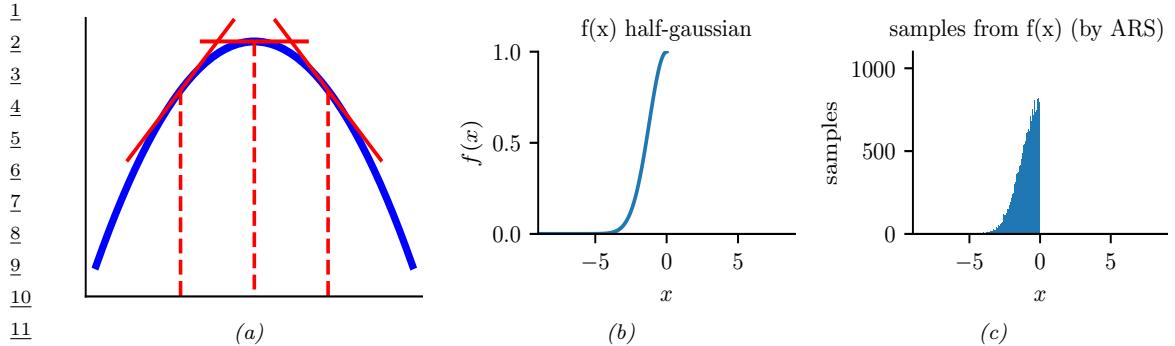
We now describe a method that can automatically come up with a tight upper envelope  $q(x)$  to any log concave 1d density  $p(x)$ . The idea is to upper bound the log density with a piecewise linear function, as illustrated in Figure 11.5(a). We choose the initial locations for the pieces based on a fixed grid over the support of the distribution. We then evaluate the gradient of the log density at these locations, and make the lines be tangent at these points.

Since the log of the envelope is piecewise linear, the envelope itself is piecewise exponential:

$$q(x) = C_i \lambda_i \exp(-\lambda_i(x - x_{i-1})), \quad x_{i-1} < x \leq x_i \quad (11.29)$$

where  $x_i$  are the grid points. It is relatively straightforward to sample from this distribution. If the sample  $x$  is rejected, we create a new grid point at  $x$ , and thereby refine the envelope. As the number of grid points is increased, the tightness of the envelope improves, and the rejection rate goes down. This is known as **adaptive rejection sampling** (ARS) [GW92]. Figure 11.5(b-c) gives an example of the method in action. As with standard rejection sampling, it can be applied to unnormalized distributions.

<sup>3</sup> This section is based on notes by Ioana A. Cosma, available at <http://users.aims.ac.za/~ioana/cp2.pdf>.



*Figure 11.5: (a) Idea behind adaptive rejection sampling. We place piecewise linear upper (and lower) bounds on the log-concave density. Adapted from Figure 1 of [GW92]. Generated by `ars_envelope.ipynb`. (b-c) Using ARS to sample from a half-Gaussian. Generated by `ars_demo.ipynb`.*

16

17

#### 18 11.4.4 Rejection sampling in high dimensions

19 It is clear that we want to make our proposal  $q(\mathbf{x})$  as close as possible to the target distribution  $p(\mathbf{x})$ ,  
 20 while still being an upper bound. But this is quite hard to achieve, especially in high dimensions. To  
 21 see this, consider sampling from  $p(\mathbf{x}) = \mathcal{N}(\mathbf{0}, \sigma_p^2 \mathbf{I})$  using as a proposal  $q(\mathbf{x}) = \mathcal{N}(\mathbf{0}, \sigma_q^2 \mathbf{I})$ . Obviously  
 22 we must have  $\sigma_q^2 \geq \sigma_p^2$  in order to be an upper bound. In  $D$  dimensions, the optimum value is given  
 23 by  $C = (\sigma_q/\sigma_p)^D$ . The acceptance rate is  $1/C$  (since both  $p$  and  $q$  are normalized), which decreases  
 24 exponentially fast with dimension. For example, if  $\sigma_q$  exceeds  $\sigma_p$  by just 1%, then in 1000 dimensions  
 25 the acceptance ratio will be about  $1/20,000$ . This is a fundamental weakness of rejection sampling.  
 26

27

#### 28 11.5 Importance sampling

29

30 In this section, we describe a Monte Carlo method known as **importance sampling** for approximating integrals of the form  
 31

$$33 \quad \mathbb{E}[\varphi(\mathbf{x})] = \int \varphi(\mathbf{x})\pi(\mathbf{x})d\mathbf{x} \quad (11.30)$$

35

36 where  $\varphi$  is called a **target function**, and  $\pi(\mathbf{x})$  is the **target distribution**, often a conditional  
 37 distribution of the form  $\pi(\mathbf{x}) = p(\mathbf{x}|\mathbf{y})$ . Since in general it is difficult to draw from the target  
 38 distribution, we will instead draw from some **proposal distribution**  $q(\mathbf{x})$  (which will usually  
 39 depend on  $\mathbf{y}$ ). We then adjust for the inaccuracies of this by associating weights with each sample,  
 40 so we end up with a weighted MC approximation:

$$41 \quad \mathbb{E}[\varphi(\mathbf{x})] \approx \sum_{n=1}^N W_n \varphi(\mathbf{x}_n) \quad (11.31)$$

44

45 We discuss two cases, first when the target is normalized, and then when it is unnormalized. This  
 46 will affect the ways the weights are computed, as well as statistical properties of the estimator.  
 47

1

### 11.5.1 Direct importance sampling

2

In this section, we assume that we can *evaluate* the normalized target distribution  $\pi(\mathbf{x})$ , but we  
3 cannot sample from it. So instead we will sample from the proposal  $q(\mathbf{x})$ . We can then write  
4

5

$$\int \varphi(\mathbf{x})\pi(\mathbf{x})d\mathbf{x} = \int \varphi(\mathbf{x})\frac{\pi(\mathbf{x})}{q(\mathbf{x})}q(\mathbf{x})d\mathbf{x} \quad (11.32)$$

6

We require that the proposal be non-zero whenever the target is non-zero, i.e., the support of  $q(\mathbf{x})$   
7 needs to be greater or equal to the support of  $\pi(\mathbf{x})$ . If we draw  $N_s$  samples  $\mathbf{x}_n \sim q(\mathbf{x})$ , we can write  
8

9

$$\mathbb{E}[\varphi(\mathbf{x})] \approx \frac{1}{N_s} \sum_{n=1}^{N_s} \frac{\pi(\mathbf{x}_n)}{q(\mathbf{x}_n)} \varphi(\mathbf{x}_n) = \frac{1}{N_s} \sum_{n=1}^{N_s} \tilde{w}_n \varphi(\mathbf{x}_n) \quad (11.33)$$

10

where we have defined the **importance weights** as follows:

11

$$\tilde{w}_n = \frac{\pi(\mathbf{x}_n)}{q(\mathbf{x}_n)} \quad (11.34)$$

12

The result is an unbiased estimate of the true mean  $\mathbb{E}[\varphi(\mathbf{x})]$ .

13

### 11.5.2 Self-normalized importance sampling

14

The disadvantage of direct importance sampling is that we need a way to evaluate the normalized  
15 target distribution  $\pi$  in order to compute the weights. It is often much easier to evaluate the  
16 **unnormalized target distribution**

17

$$\tilde{\gamma}(\mathbf{x}) = Z\pi(\mathbf{x}) \quad (11.35)$$

18

where

19

$$Z = \int \tilde{\gamma}(\mathbf{x})d\mathbf{x} \quad (11.36)$$

20

is the normalization constant. (For example, if  $\pi(\mathbf{x}) = p(\mathbf{x}|\mathbf{y})$ , then  $\tilde{\gamma}(\mathbf{x}) = p(\mathbf{x}, \mathbf{y})$  and  $Z = p(\mathbf{y})$ .)  
21 The key idea is to also approximate the normalization constant  $Z$  with importance sampling. This  
22 method is called **self-normalized importance sampling**. The resulting estimate is a ratio of  
23 two estimates, and hence is biased. However as  $N_s \rightarrow \infty$ , the bias goes to zero, under some weak  
24 assumptions (see e.g., [RC04] for details).

25

In more detail, SNIS is based on this approximation:

26

$$\mathbb{E}[\varphi(\mathbf{x})] = \int \varphi(\mathbf{x})\pi(\mathbf{x})d\mathbf{x} = \frac{\int \varphi(\mathbf{x})\tilde{\gamma}(\mathbf{x})d\mathbf{x}}{\int \tilde{\gamma}(\mathbf{x})d\mathbf{x}} = \frac{\int \left[ \frac{\tilde{\gamma}(\mathbf{x})}{q(\mathbf{x})} \right] \varphi(\mathbf{x})d\mathbf{x}}{\int \left[ \frac{\tilde{\gamma}(\mathbf{x})}{q(\mathbf{x})} \right] q(\mathbf{x})d\mathbf{x}} \quad (11.37)$$

27

$$\approx \frac{\frac{1}{N_s} \sum_{n=1}^{N_s} \tilde{w}_n \varphi(\mathbf{x}_n)}{\frac{1}{N_s} \sum_{n=1}^{N_s} \tilde{w}_n} \quad (11.38)$$

1 where we have defined the **unnormalized weights**  
2

3  
4  $\tilde{w}_n = \frac{\tilde{\gamma}(\mathbf{x}_n)}{q(\mathbf{x}_n)}$  (11.39)  
5

6 We can write Equation (11.38) more compactly as  
7

8  
9  $\mathbb{E}[\varphi(\mathbf{x})] \approx \sum_{n=1}^{N_s} W_n \varphi(\mathbf{x}_n)$  (11.40)  
10

11 where we have defined the **normalized weights** by  
12

13  
14  $W_n = \frac{\tilde{w}_n}{\sum_{n'=1}^{N_s} \tilde{w}_{n'}}$  (11.41)  
15

16 This is equivalent to approximating the target distribution using a weighted sum of delta functions:  
17

18  
19  $\pi(\mathbf{x}) \approx \sum_{n=1}^{N_s} W_n \delta(\mathbf{x} - \mathbf{x}_n) \triangleq \hat{\pi}(\mathbf{x})$  (11.42)  
20

21 As a byproduct of this algorithm we get the following appoximation to the normalization constant:  
22

23  
24  $Z \approx \frac{1}{N_s} \sum_{n=1}^{N_s} \tilde{w}_n \triangleq \hat{Z}$  (11.43)  
25

### 26 11.5.3 Choosing the proposal 27

28 The performance of importance sampling depends crucially on the quality of the proposal distribution.  
29 As we mentioned, we require that the support of  $q$  cover the support of the target (i.e.,  $\tilde{\gamma}(\mathbf{x}) > 0 \implies q(\mathbf{x}) > 0$ ). However, we also want the proposal to not be too “loose” or a “covering”. Ideally it should also take into account properties of the target function  $\varphi$  as well, as shown in Figure 11.6. This can yield subsantial benefits, as shown in the “**target aware Bayesian inference**” scheme of [Rai+20].  
32 However, usually the target function  $\varphi$  is unknown or ignored, so we just try to find a “generally useful” approximation to the target.

35 One way to come up with a good proposal is to learn one, by optimizing the variational lower bound or ELBO (see Section 10.1.2). Indeed, if we fix the parameters of the generative model, we can think of importance weighted autoencoders (Section 10.5.1) as learning a good IS proposal. More details on this connection can be found in [DS18].  
39

### 40 11.5.4 Annealed importance sampling (AIS) 41

42 In this section, we describe a method known as **annealed importance sampling** [Nea01] for sampling from complex, possibly multimodal distributions. Assume we want to sample from some target distribution  $p_0(\mathbf{x}) \propto f_0(\mathbf{x})$  (where  $f_0(\mathbf{x})$  is the unnormalized version), but we cannot easily do so, because  $p_0$  is complicated in some way (e.g., high dimensional and/or multi-modal). However, suppose that there is an easier distribution which we *can* sample from, call it  $p_n(\mathbf{x}) \propto f_n(\mathbf{x})$ ; for

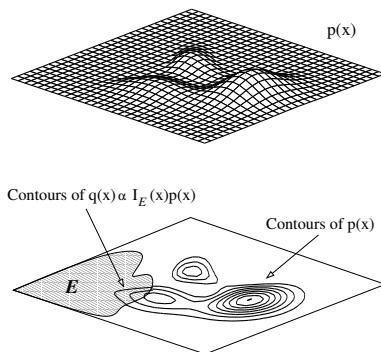


Figure 11.6: In importance sampling, we should sample from a distribution that takes into account regions where  $\pi(\mathbf{x})$  has high probability and where  $\varphi(\mathbf{x})$  is large. Here the function to be evaluated is an indicator function of a set, corresponding to a set of rare events in the tail of the distribution. From Figure 3 of [And+03]. Used with kind permission of Nando de Freitas.

example, this might be the prior. We now construct a sequence of intermediate distributions than move slowly from  $p_n$  to  $p_0$  as follows:

$$f_j(\mathbf{x}) = f_0(\mathbf{x})^{\beta_j} f_n(\mathbf{x})^{1-\beta_j} \quad (11.44)$$

where  $1 = \beta_0 > \beta_1 > \dots > \beta_n = 0$ , where  $\beta_j$  is an inverse temperature. We will sample a set of points from  $f_n$ , and then from  $f_{n-1}$ , and so on, until we eventually sample from  $f_0$ .

To sample from each  $f_j$ , suppose we can define a Markov chain  $T_j(\mathbf{x}, \mathbf{x}') = p_j(\mathbf{x}'|\mathbf{x})$ , which leaves  $p_0$  invariant (i.e.,  $\int p_j(\mathbf{x}'|\mathbf{x})p_0(\mathbf{x})d\mathbf{x} = p_0(\mathbf{x}')$ ). (See Chapter 12 for details on how to construct such chains.) Given this, we can sample  $\mathbf{x}$  from  $p_0$  as follows: sample  $\mathbf{v}_n \sim p_n$ ; sample  $\mathbf{v}_{n-1} \sim T_{n-1}(\mathbf{v}_n, \cdot)$ ; and continue in this way until we sample  $\mathbf{v}_0 \sim T_0(\mathbf{v}_1, \cdot)$ ; finally we set  $\mathbf{x} = \mathbf{v}_0$  and give it weight

$$w = \frac{f_{n-1}(\mathbf{v}_{n-1})}{f_n(\mathbf{v}_{n-1})} \frac{f_{n-2}(\mathbf{v}_{n-2})}{f_{n-1}(\mathbf{v}_{n-2})} \dots \frac{f_1(\mathbf{v}_1)}{f_2(\mathbf{v}_1)} \frac{f_0(\mathbf{v}_0)}{f_1(\mathbf{v}_0)} \quad (11.45)$$

This can be shown to be correct by viewing the algorithm as a form of importance sampling in an extended state space  $\mathbf{v} = (\mathbf{v}_0, \dots, \mathbf{v}_n)$ . Consider the following distribution on this state space:

$$p(\mathbf{v}) \propto \varphi(\mathbf{v}) = f_0(\mathbf{v}_0)\tilde{T}_0(\mathbf{v}_0, \mathbf{v}_1)\tilde{T}_2(\mathbf{v}_1, \mathbf{v}_2) \dots \tilde{T}_{n-1}(\mathbf{v}_{n-1}, \mathbf{v}_n) \quad (11.46)$$

$$\propto p(\mathbf{v}_0)p(\mathbf{v}_1|\mathbf{v}_0) \dots p(\mathbf{v}_n|\mathbf{v}_{n-1}) \quad (11.47)$$

where  $\tilde{T}_j$  is the reversal of  $T_j$ :

$$\tilde{T}_j(\mathbf{v}, \mathbf{v}') = T_j(\mathbf{v}', \mathbf{v})p_j(\mathbf{v}')/p_j(\mathbf{v}) = T_j(\mathbf{v}', \mathbf{v})f_j(\mathbf{v}')/f_j(\mathbf{v}) \quad (11.48)$$

It is clear that  $\sum_{\mathbf{v}_1, \dots, \mathbf{v}_n} \varphi(\mathbf{v}) = f_0(\mathbf{v}_0)$ , so by sampling from  $p(\mathbf{v})$ , we can effectively sample from  $p_0(\mathbf{x})$ .

We can sample on this extended state space using the above algorithm, which corresponds to the following proposal:

$$q(\mathbf{v}) \propto g(\mathbf{v}) = f_n(\mathbf{v}_n)T_{n-1}(\mathbf{v}_n, \mathbf{v}_{n-1}) \cdots T_2(\mathbf{v}_2, \mathbf{v}_1)T_0(\mathbf{v}_1, \mathbf{v}_0) \quad (11.49)$$

$$\propto p(\mathbf{v}_n)p(\mathbf{v}_{n-1}|\mathbf{v}_n) \cdots p(\mathbf{v}_1|\mathbf{v}_0) \quad (11.50)$$

One can show that the importance weights  $w = \frac{\varphi(\mathbf{v}_0, \dots, \mathbf{v}_n)}{g(\mathbf{v}_0, \dots, \mathbf{v}_n)}$  are given by Equation (11.45). Since marginals of the sampled sequences from this extended model are equivalent to samples from  $p_0(\mathbf{x})$ , we see that we are using the correct weights.

#### 11.5.4.1 Estimating normalizing constants using AIS

An important application of AIS is to evaluate a ratio of partition functions. Notice that  $Z_0 = \int f_0(\mathbf{x})d\mathbf{x} = \int \varphi(\mathbf{v})d\mathbf{v}$ , and  $Z_n = \int f_n(\mathbf{x})d\mathbf{x} = \int g(\mathbf{v})d\mathbf{v}$ . Hence

$$\frac{Z_0}{Z_n} = \frac{\int \varphi(\mathbf{v})d\mathbf{v}}{\int g(\mathbf{v})d\mathbf{v}} = \frac{\int \frac{\varphi(\mathbf{v})}{g(\mathbf{v})}g(\mathbf{v})d\mathbf{v}}{\int g(\mathbf{v})d\mathbf{v}} = \mathbb{E}_g \left[ \frac{\varphi(\mathbf{v})}{g(\mathbf{v})} \right] \approx \frac{1}{S} \sum_{s=1}^S w_s \quad (11.51)$$

where  $w_s = \varphi(\mathbf{v}_s)/g(\mathbf{v}_s)$ . If  $f_0$  is a prior and  $f_n$  is the posterior, we can estimate  $Z_n = p(\mathcal{D})$  using the above equation, provided the prior has a known normalization constant  $Z_0$ . This is generally considered the method of choice for evaluating difficult partition functions.

23

24

## 11.6 Controlling Monte Carlo variance

27 As we mentioned in Section 11.2.2, the standard error in a Monte Carlo estimate is  $O(1/\sqrt{S})$ , where  
28  $S$  is the number of (independent) samples. Consequently it may take many samples to reduce the  
29 variance to a sufficiently small value. In this section, we discuss some ways to reduce the variance of  
30 sampling methods. For more details, see e.g., [KTB11].  
31

### 11.6.1 Common random numbers

34 When performing Monte Carlo optimization, we often want to compare  $\mathbb{E}_{p(\mathbf{z})}[f(\boldsymbol{\theta}, \mathbf{z})]$  to  $\mathbb{E}_{p(\mathbf{z})}[f(\boldsymbol{\theta}', \mathbf{z})]$   
35 for different values of the parameters  $\boldsymbol{\theta}$  and  $\boldsymbol{\theta}'$ . To reduce the variance of this comparison, we can  
36 use the same random samples  $\mathbf{z}_s$  for evaluating both functions. In this way, differences in the  
37 outcome can be ascribed to differences in the parameters  $\boldsymbol{\theta}$ , rather than to the noise terms. This  
38 is called the **common random numbers** trick, and is widely used in ML (see e.g., [GBJ18;  
39 NJ00]), since it can often convert a stochastic optimization problem into a deterministic one,  
40 enabling the use of more powerful optimization methods. For more details on CRN, see e.g.,  
41 [https://en.wikipedia.org/wiki/Variance\\_reduction#Common\\_Random\\_Numbers\\_\(CRN\)](https://en.wikipedia.org/wiki/Variance_reduction#Common_Random_Numbers_(CRN)).  
42

### 11.6.2 Rao-Blackwellisation

45 In this section, we discuss a useful technique for reducing the variance of MC estimators known as  
46 **Rao-Blackwellisation**. To explain the method, suppose we have two rv's,  $X$  and  $Y$ , and we want  
47

1 to estimate  $\bar{f} = \mathbb{E}[f(X, Y)]$ . The naive approach is to use an MC approximation  
 2

3

$$4 \hat{f}_{MC} = \frac{1}{S} \sum_{s=1}^S f(X_s, Y_s) \quad (11.52)$$

5

6 where  $(X_s, Y_s) \sim p(X, Y)$ . This is an unbiased estimator of  $\bar{f}$ . However, it may have high variance.  
 7

8 Now suppose we can analytically marginalize out  $Y$ , provided we know  $X$ , i.e., we can tractably  
 9 compute  
 10

11

$$12 f_X(X_s) = \int dY p(Y|X_s) f(X_s, Y) = \mathbb{E}[f(X, Y)|X = X_s] \quad (11.53)$$

13

14 Let us define the Rao-Blackwellised estimator  
 15

16

$$17 \hat{f}_{RB} = \frac{1}{S} \sum_{s=1}^S f_X(X_s) \quad (11.54)$$

18

19 where  $X_s \sim p(X)$ . This is an unbiased estimator, since  $\mathbb{E}[\hat{f}_{RB}] = \mathbb{E}[\mathbb{E}[f(X, Y)|X]] = \bar{f}$ . However,  
 20 this estimate can have lower variance than the naive estimator. The intuitive reason is that we are  
 21 now sampling in a reduced dimensional space. Formally we can see this by using the law of iterated  
 22 variance to get  
 23

24

$$25 \mathbb{V}[\mathbb{E}[f(X, Y)|X]] = \mathbb{V}[f(X, Y)] - \mathbb{E}[\mathbb{V}[f(X, Y)]|X] \leq \mathbb{V}[f(X, Y)] \quad (11.55)$$

26

27 For some examples of this in practice, see Section 6.5.3.2, Section 13.4, and Section 12.3.8.  
 28

### 29 11.6.3 Control variates

30 Suppose we want to estimate  $\mu = \mathbb{E}[f(X)]$  using an unbiased estimator  $m(\mathcal{X}) = \frac{1}{S} \sum_{s=1}^S m(x_s)$ ,  
 31 where  $x_s \sim p(X)$  and  $\mathbb{E}[m(X)] = \mu$ . (We abuse notation slightly and use  $m$  to refer to a function of  
 32 a single random variable as well as a set of samples.) Now consider the alternative estimator  
 33

34

$$35 m^*(\mathcal{X}) = m(\mathcal{X}) + c(b(\mathcal{X}) - \mathbb{E}[b(\mathcal{X})]) \quad (11.56)$$

36

37 This is called a **control variate**, and  $b$  is called a **baseline**. (Once again we abuse notation and use  
 38  $b(\mathcal{X}) = \frac{1}{S} \sum_{s=1}^S b(x_s)$  and  $m^*(\mathcal{X}) = \frac{1}{S} \sum_{s=1}^S m^*(x_s)$ .)

39 It is easy to see that  $m^*(\mathcal{X})$  is an unbiased estimator, since  $\mathbb{E}[m^*(X)] = \mathbb{E}[m(X)] = \mu$ . However,  
 40 it can have lower variance, provided  $b$  is correlated with  $m$ . To see this, note that

41

$$42 \mathbb{V}[m^*(X)] = \mathbb{V}[m(X)] + c^2 \mathbb{V}[b(X)] + 2c \text{Cov}[m(X), b(X)] \quad (11.57)$$

43

44 By taking the derivative of  $\mathbb{V}[m^*(X)]$  wrt  $c$  and setting to 0, we find that the optimal value is  
 45

46

$$47 c^* = -\frac{\text{Cov}[m(X), b(X)]}{\mathbb{V}[b(X)]} \quad (11.58)$$

1 The corresponding variance of the new estimator is now  
2

$$\frac{3}{4} \quad \mathbb{V}[m^*(X)] = \mathbb{V}[m(X)] - \frac{\text{Cov}[m(X), b(X)]^2}{\mathbb{V}[b(X)]} = (1 - \rho_{m,b}^2) \mathbb{V}[m(X)] \leq \mathbb{V}[m(X)] \quad (11.59)$$

5

6 where  $\rho_{m,b}^2$  is the correlation of the basic estimator and the baseline function. If we can ensure  
7 this correlation is high, we can reduce the variance. Intuitively, the CV estimator is exploiting  
8 information about the errors in the estimate of a known quantity, namely  $\mathbb{E}[b(X)]$ , to reduce the  
9 errors in estimating the unknown quantity, namely  $\mu$ .

10 We give a simple worked example in Section 11.6.3.1. See Section 10.3.2 for an example of this  
11 technique applied to blackbox variational inference.

12

### 13 11.6.3.1 Example

14 We now give a simple worked example of control variates.<sup>4</sup> Consider estimating  $\mu = \mathbb{E}[f(X)]$  where  
15  $f(X) = 1/(1+X)$  and  $X \sim \text{Unif}(0, 1)$ . The exact value is  
16

$$\frac{17}{18} \quad \mu = \int_0^1 \frac{1}{1+x} dx = \ln 2 \approx 0.693 \quad (11.60)$$

19 The naive MC estimate, using  $S$  samples, is  $m(\mathcal{X}) = \frac{1}{S} \sum_{s=1}^S f(x_s)$ . Using  $S = 1500$ , we find  
20  $\mathbb{E}[m(\mathcal{X})] = 0.6935$  with standard error  $\text{se} = 0.0037$ .

22 Now let us use  $b(X) = 1 + X$  as a baseline, so  $b(\mathcal{X}) = (1/S) \sum_s (1 + x_s)$ . This has expectation  
23  $\mathbb{E}[b(X)] = \int_0^1 (1+x) dx = \frac{3}{2}$ . The control variate estimator is given by

$$\frac{24}{25} \quad m^*(\mathcal{X}) = \frac{1}{S} \sum_{s=1}^S f(x_s) + c \left( \frac{1}{S} \sum_{s=1}^S b(x_s) - \frac{3}{2} \right) \quad (11.61)$$

27 The optimal value can be estimated from the samples of  $m(x_s)$  and  $b(x_s)$ , and plugging into  
28 Equation (11.58) to get  $c^* \approx 0.4773$ . Using  $S = 1500$ , we find  $\mathbb{E}[m^*(\mathcal{X})] = 0.6941$  and  $\text{se} = 0.0007$ .

29 See also Section 11.6.4.1, where we analyze this example using antithetic sampling.

30

### 31 11.6.4 Antithetic sampling

33 In this section, we discuss **antithetic sampling**, which is a simple way to reduce variance.<sup>5</sup> Suppose  
34 we want to estimate  $\theta = \mathbb{E}[Y]$ . Let  $Y_1$  and  $Y_2$  be two samples. An unbiased estimate of  $\theta$  is given by  
35  $\hat{\theta} = (Y_1 + Y_2)/2$ . The variance of this estimate is

$$\frac{37}{38} \quad \mathbb{V}[\hat{\theta}] = \frac{\mathbb{V}[Y_1] + \mathbb{V}[Y_2] + 2\text{Cov}[Y_1, Y_2]}{4} \quad (11.62)$$

39 so the variance is reduced if  $\text{Cov}[Y_1, Y_2] < 0$ . So whenever we sample  $Y_1$ , we should set  $Y_2$  to be its  
40 “opposite”, but with the same mean.

41 For example, suppose  $Y \sim \text{Unif}(0, 1)$ . If we let  $y_1, \dots, y_n$  be iid samples from  $\text{Unif}(0, 1)$ , then we  
42 can define  $y'_i = 1 - y_i$ . The distribution of  $y'_i$  is still  $\text{Unif}(0, 1)$ , but  $\text{Cov}[y_i, y'_i] < 1$ .

43

44 4. The example is from [https://en.wikipedia.org/wiki/Control\\_variates](https://en.wikipedia.org/wiki/Control_variates), with modified notation. See [control\\_variates.ipynb](#) for some code.

45 5. Our presentation is based on [https://en.wikipedia.org/wiki/Antithetic\\_variates](https://en.wikipedia.org/wiki/Antithetic_variates). See [antithetic\\_sampling.ipynb](#) for the code.

47

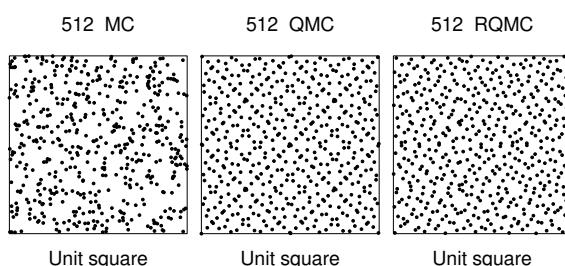


Figure 11.7: Illustration of Monte Carlo (MC), Quasi MC (QMC) from a Sobol sequence, and Randomized QMC using a scrambling method. Adapted from Figure 1 of [OR20]. Used with kind permission of Art Owen.

#### 11.6.4.1 Example

To see why this can be useful, consider the example from Section 11.6.3.1. Let  $\hat{\mu}_{\text{mc}}$  be the classic MC estimate using  $2N$  samples from  $\text{Unif}(0, 1)$ , and let  $\hat{\mu}_{\text{anti}}$  be the MC estimate using the above antithetic sampling scheme applied to  $N$  base samples from  $\text{Unif}(0, 1)$ . The exact value is  $\mu = \ln 2 \approx 0.6935$ . For the classical method, with  $N = 750$ , we find  $\mathbb{E}[\hat{\mu}_{\text{mc}}] = 0.69365$  with a standard error of 0.0037. For the antithetic method, we find  $\mathbb{E}[\hat{\mu}_{\text{anti}}] = 0.6939$  with a standard error of 0.0007, which matches the control variate method of Section 11.6.3.1.

#### 11.6.5 Quasi Monte Carlo (QMC)

**Quasi Monte Carlo** (see e.g., [Lem09; Owe13]) is an approach to numerical integration that replaces random samples with **low discrepancy sequences**, such as the **Halton sequence** (see e.g., [Owe17]) or **Sobol sequence**. Intuitively, these are **space filling** sequences of points, constructed to reduce the unwanted gaps and clusters that would arise among randomly chosen inputs. See Figure 11.7 for an example.<sup>6</sup>

More precisely, consider the problem of evaluating the following  $D$ -dimensional integral:

$$\bar{f} = \int_{[0,1]^D} f(\mathbf{x}) d\mathbf{x} \approx \hat{f}_N = \frac{1}{N} \sum_{n=1}^N f(\mathbf{x}_n) \quad (11.63)$$

Let  $\epsilon_N = |\bar{f} - \hat{f}_N|$  be the error. In standard Monte Carlo, if we draw  $N$  independent samples, then we have  $\epsilon_N \sim O\left(\frac{1}{\sqrt{N}}\right)$ . In QMC, it can be shown that  $\epsilon_N \sim O\left(\frac{(\log N)^D}{N}\right)$ . For  $N > 2^D$ , the latter is smaller than the former.

One disadvantage of QMC is that it just provides a point estimate of  $\bar{f}$ , and does not give an uncertainty estimate. By contrast, in regular MC, we can estimate the MC standard error, discussed in Section 11.2.2. **Randomized QMC** (see e.g., [L'E18]) provides a solution to this problem. The basic idea is to repeat the QMC method  $R$  times, by perturbing the sequence of  $N$  points by a

6. More details on QMC can be found at [http://roth.cs.kuleuven.be/wiki/Main\\_Page](http://roth.cs.kuleuven.be/wiki/Main_Page). For connections to Bayesian quadrature, see e.g., [DKS13; HKO22].

1 random amount. In particular, define  
2

3  $\mathbf{y}_{i,r} = \mathbf{x}_i + \mathbf{u}_r \pmod{1}$  (11.64)  
4

5 where  $\mathbf{x}_1, \dots, \mathbf{x}_N$  is a low-discrepancy sequence, and  $\mathbf{u}_r \sim \text{Unif}(0, 1)^D$  is a random perturbation.  
6 The set  $\{\mathbf{y}_j\}$  is low discrepancy, and satisfies that each  $\mathbf{y}_j \sim \text{Unif}(0, 1)^D$ , for  $j = 1 : N \times R$ . This  
7 has much lower variance than standard MC. (Typically we take  $R$  to be a power of 2.) Recently,  
8 [OR20] proved a strong law of large numbers for RQMC.

9 QMC and RQMC can be used inside of MCMC inference (see e.g., [OT05]) and variational inference  
10 (see e.g., [BWM18]). It is also commonly used to select the initial set of query points for Bayesian  
11 optimization (Section 6.8).

12 Another technique that can be used is **orthogonal Monte Carlo**, where the samples are condi-  
13 tioned to be pairwise orthogonal, but with the marginal distributions matching the original ones (see  
14 e.g., [Lin+20]).  
15

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

# 12 Markov Chain Monte Carlo inference

## 12.1 Introduction

In Chapter 11, we considered non-iterative Monte Carlo methods, including rejection sampling and importance sampling, which generate independent samples from some target distribution. The trouble with these methods is that they often do not work well in high dimensional spaces. In this chapter, we discuss a popular method for sampling from high-dimensional distributions known as **Markov chain Monte Carlo** or **MCMC**. In a survey by *SIAM News*<sup>1</sup>, MCMC was placed in the top 10 most important algorithms of the 20th century.

The basic idea behind MCMC is to construct a Markov chain (Section 2.6) on the state space  $\mathcal{X}$  whose stationary distribution is the target density  $p^*(\mathbf{x})$  of interest. (In a Bayesian context, this is usually a posterior,  $p^*(\mathbf{x}) \propto p(\mathbf{x}|\mathcal{D})$ , but MCMC can be applied to generate samples from any kind of distribution.) That is, we perform a random walk on the state space, in such a way that the fraction of time we spend in each state  $\mathbf{x}$  is proportional to  $p^*(\mathbf{x})$ . By drawing (correlated) samples  $\mathbf{x}_0, \mathbf{x}_1, \mathbf{x}_2, \dots$ , from the chain, we can perform Monte Carlo integration wrt  $p^*$ .

Note that the initial samples from the chain do not come from the stationary distribution, and should be discarded; the amount of time it takes to reach stationarity is called the **mixing time** or **burn-in time**; reducing this is one of the most important factors in making the algorithm fast, as we will see.

The MCMC algorithm has an interesting history. It was discovered by physicists working on the atomic bomb at Los Alamos during World War II, and was first published in the open literature in [Met+53] in a chemistry journal. An extension was published in the statistics literature in [Has70], but was largely unnoticed. A special case (Gibbs sampling, Section 12.3) was independently invented in [GG84] in the context of Ising models (Section 4.3.2.1). But it was not until [GS90] that the algorithm became well-known to the wider statistical community. Since then it has become wildly popular in Bayesian statistics, and is becoming increasingly popular in machine learning.

In the rest of this chapter, we give a brief introduction to MCMC methods. For more details on the theory, see e.g., [GRS96; BZ20]. For more details on the implementation side, see e.g., [Lao+20]. And for an interactive visualization of many of these algorithms in 2d, see <http://chi-feng.github.io/mcmc-demo/app.html>.

---

1. Source: <http://www.siam.org/pdf/news/637.pdf>.

1 **12.2 Metropolis Hastings algorithm**

3 In this section, we describe the simplest kinds of MCMC algorithm known as the **Metropolis**  
4 **Hastings** or **MH** algorithm.

6 **12.2.1 Basic idea**

8 The basic idea in MH is that at each step, we propose to move from the current state  $\mathbf{x}$  to a new state  
9  $\mathbf{x}'$  with probability  $q(\mathbf{x}'|\mathbf{x})$ , where  $q$  is called the **proposal distribution** (also called the **kernel**).  
10 The user is free to use any kind of proposal they want, subject to some conditions which we explain  
11 below. This makes MH quite a flexible method.

12 Having proposed a move to  $\mathbf{x}'$ , we then decide whether to **accept** this proposal, or to reject it,  
13 according to some formula, which ensures that the long-term fraction of time spent in each state is  
14 proportional to  $p^*(\mathbf{x})$ . If the proposal is accepted, the new state is  $\mathbf{x}'$ , otherwise the new state is the  
15 same as the current state,  $\mathbf{x}$  (i.e., we repeat the sample).

16 If the proposal is symmetric, so  $q(\mathbf{x}'|\mathbf{x}) = q(\mathbf{x}|\mathbf{x}')$ , the acceptance probability is given by the  
17 following formula:

$$\underline{19} \quad A = \min \left( 1, \frac{p^*(\mathbf{x}')}{p^*(\mathbf{x})} \right) \quad (12.1)$$

21 We see that if  $\mathbf{x}'$  is more probable than  $\mathbf{x}$ , we definitely move there (since  $\frac{p^*(\mathbf{x}')}{p^*(\mathbf{x})} > 1$ ), but if  $\mathbf{x}'$  is  
22 less probable, we may still move there anyway, depending on the relative probabilities. So instead of  
23 greedily moving to only more probable states, we occasionally allow “downhill” moves to less probable  
24 states. In Section 12.2.2, we prove that this procedure ensures that the fraction of time we spend in  
25 each state  $\mathbf{x}$  is equal to  $p^*(\mathbf{x})$ .

27 If the proposal is asymmetric, so  $q(\mathbf{x}'|\mathbf{x}) \neq q(\mathbf{x}|\mathbf{x}')$ , we need the **Hastings correction**, given by  
28 the following:

$$\underline{29} \quad A = \min(1, \alpha) \quad (12.2)$$

$$\underline{30} \quad \alpha = \frac{p^*(\mathbf{x}')q(\mathbf{x}|\mathbf{x}')}{p^*(\mathbf{x})q(\mathbf{x}'|\mathbf{x})} = \frac{p^*(\mathbf{x}')/q(\mathbf{x}'|\mathbf{x})}{p^*(\mathbf{x})/q(\mathbf{x}|\mathbf{x}')} \quad (12.3)$$

33 This correction is needed to compensate for the fact that the proposal distribution itself (rather than  
34 just the target distribution) might favor certain states.

35 An important reason why MH is a useful algorithm is that, when evaluating  $\alpha$ , we only need to  
36 know the target density up to a normalization constant. In particular, suppose  $p^*(\mathbf{x}) = \frac{1}{Z}\tilde{p}(\mathbf{x})$ , where  
37  $\tilde{p}(\mathbf{x})$  is an unnormalized distribution and  $Z$  is the normalization constant. Then

$$\underline{38} \quad \alpha = \frac{(\tilde{p}(\mathbf{x}')/Z) q(\mathbf{x}|\mathbf{x}')}{(\tilde{p}(\mathbf{x})/Z) q(\mathbf{x}'|\mathbf{x})} \quad (12.4)$$

41 so the  $Z$ 's cancel. Hence we can sample from  $p^*$  even if  $Z$  is unknown.

42 A proposal distribution  $q$  is valid or admissible if it “covers” the support of the target. Formally,  
43 we can write this as

$$\underline{44} \quad \text{supp}(p^*) \subseteq \cup_x \text{supp}(q(\cdot|x)) \quad (12.5)$$

46 With this, we can state the overall algorithm as in Algorithm 25.

47

---

**Algorithm 25:** Metropolis Hastings algorithm

---

```

1 Initialize  $x^0$ 
2 for  $s = 0, 1, 2, \dots$  do
3   Define  $x = x^s$ 
4   Sample  $x' \sim q(x'|x)$ 
5   Compute acceptance probability
6
7   
$$\alpha = \frac{\tilde{p}(x')q(x|x')}{\tilde{p}(x)q(x'|x)}$$

8
9   Compute  $A = \min(1, \alpha)$ 
10  Sample  $u \sim U(0, 1)$ 
11  Set new sample to
12
13  
$$x^{s+1} = \begin{cases} x' & \text{if } u \leq A \text{ (accept)} \\ x^s & \text{if } u > A \text{ (reject)} \end{cases}$$

14
15
16
17
18
19
20
21
22
```

---

**12.2.2 Why MH works**

To prove that the MH procedure generates samples from  $p^*$ , we need a bit of Markov chain theory, as discussed in Section 2.6.4.

The MH algorithm defines a Markov chain with the following transition matrix:

$$p(\mathbf{x}'|\mathbf{x}) = \begin{cases} q(\mathbf{x}'|\mathbf{x})A(\mathbf{x}'|\mathbf{x}) & \text{if } \mathbf{x}' \neq \mathbf{x} \\ q(\mathbf{x}|\mathbf{x}) + \sum_{\mathbf{x}' \neq \mathbf{x}} q(\mathbf{x}'|\mathbf{x})(1 - A(\mathbf{x}'|\mathbf{x})) & \text{otherwise} \end{cases} \quad (12.6)$$

This follows from a case analysis: if you move to  $\mathbf{x}'$  from  $\mathbf{x}$ , you must have proposed it (with probability  $q(\mathbf{x}'|\mathbf{x})$ ) and it must have been accepted (with probability  $A(\mathbf{x}'|\mathbf{x})$ ); otherwise you stay in state  $\mathbf{x}$ , either because that is what you proposed (with probability  $q(\mathbf{x}|\mathbf{x})$ ), or because you proposed something else (with probability  $q(\mathbf{x}'|\mathbf{x})$ ) but it was rejected (with probability  $1 - A(\mathbf{x}'|\mathbf{x})$ ).

Let us analyse this Markov chain. Recall that a chain satisfies **detailed balance** if

$$p(\mathbf{x}'|\mathbf{x})p^*(\mathbf{x}) = p(\mathbf{x}|\mathbf{x}')p^*(\mathbf{x}') \quad (12.7)$$

This means in the in-flow to state  $\mathbf{x}'$  from  $\mathbf{x}$  is equal to the out-flow from state  $\mathbf{x}'$  back to  $\mathbf{x}$ , and vice versa. We also showed that if a chain satisfies detailed balance, then  $p^*$  is its stationary distribution. Our goal is to show that the MH algorithm defines a transition function that satisfies detailed balance and hence that  $p^*$  is its stationary distribution. (If Equation (12.7) holds, we say that  $p^*$  is an **invariant** distribution wrt the Markov transition kernel  $q$ .)

**Theorem 12.2.1.** *If the transition matrix defined by the MH algorithm (given by Equation (12.6)) is ergodic and irreducible, then  $p^*$  is its unique limiting distribution.*

*Proof.* Consider two states  $\mathbf{x}$  and  $\mathbf{x}'$ . Either

$$p^*(\mathbf{x})q(\mathbf{x}'|\mathbf{x}) < p^*(\mathbf{x}')q(\mathbf{x}|\mathbf{x}') \quad (12.8)$$

1  
2 or

3       $p^*(\mathbf{x})q(\mathbf{x}'|\mathbf{x}) \geq p^*(\mathbf{x}')q(\mathbf{x}|\mathbf{x}')$       (12.9)  
4

5 Without loss of generality, assume that  $p^*(\mathbf{x})q(\mathbf{x}'|\mathbf{x}) > p^*(\mathbf{x}')q(\mathbf{x}|\mathbf{x}')$ . Hence  
6

7       $\alpha(\mathbf{x}'|\mathbf{x}) = \frac{p^*(\mathbf{x}')q(\mathbf{x}|\mathbf{x}')}{p^*(\mathbf{x})q(\mathbf{x}'|\mathbf{x})} < 1$       (12.10)  
8

9 Hence we have  $A(\mathbf{x}'|\mathbf{x}) = \alpha(\mathbf{x}'|\mathbf{x})$  and  $A(\mathbf{x}|\mathbf{x}') = 1$ .  
10

11 Now to move from  $\mathbf{x}$  to  $\mathbf{x}'$  we must first propose  $\mathbf{x}'$  and then accept it. Hence  
12

13       $p(\mathbf{x}'|\mathbf{x}) = q(\mathbf{x}'|\mathbf{x})A(\mathbf{x}'|\mathbf{x}) = q(\mathbf{x}'|\mathbf{x})\frac{p^*(\mathbf{x}')q(\mathbf{x}|\mathbf{x}')}{p^*(\mathbf{x})q(\mathbf{x}'|\mathbf{x})} = \frac{p^*(\mathbf{x}')}{p^*(\mathbf{x})}q(\mathbf{x}|\mathbf{x}')$       (12.11)  
14

15 Hence  
16

17       $p^*(\mathbf{x})p(\mathbf{x}'|\mathbf{x}) = p^*(\mathbf{x}')q(\mathbf{x}|\mathbf{x}')$       (12.12)

18 The backwards probability is  
19

20       $p(\mathbf{x}|\mathbf{x}') = q(\mathbf{x}|\mathbf{x}')A(\mathbf{x}|\mathbf{x}') = q(\mathbf{x}|\mathbf{x}')$       (12.13)  
21

22 since  $A(\mathbf{x}|\mathbf{x}') = 1$ . Inserting this into Equation (12.12) we get  
23

24       $p^*(\mathbf{x})p(\mathbf{x}'|\mathbf{x}) = p^*(\mathbf{x}')p(\mathbf{x}|\mathbf{x}')$       (12.14)  
25

26 so detailed balance holds wrt  $p^*$ . Hence, from Theorem 2.6.3,  $p^*$  is a stationary distribution.  
27 Furthermore, from Theorem 2.6.2, this distribution is unique, since the chain is ergodic and irreducible.  
28 □

### 29 12.2.3 Proposal distributions

30 In this section, we discuss some common proposal distributions. Note, however, that good proposal  
31 design is often intimately dependent on the form of the target distribution (most often the posterior).  
32

#### 33 12.2.3.1 Independence sampler

34 If we use a proposal of the form  $q(\mathbf{x}'|\mathbf{x}) = q(\mathbf{x}')$ , where the new state is independent of the old  
35 state, we get a method known as the **independence sampler**, which is similar to importance  
36 sampling (Section 11.5). The function  $q(\mathbf{x}')$  can be any suitable distribution, such as a Gaussian.  
37 This has non-zero probability density on the entire state space, and hence is a valid proposal for any  
38 unconstrained continuous state space.  
39

#### 40 12.2.3.2 Random walk Metropolis (RWM) algorithm

41 The **random walk Metropolis** algorithm corresponds to MH with the following proposal distribu-  
42 tion:  
43

44       $q(\mathbf{x}'|\mathbf{x}) = \mathcal{N}(\mathbf{x}'|\mathbf{x}, \tau^2 \mathbf{I})$       (12.15)  
45

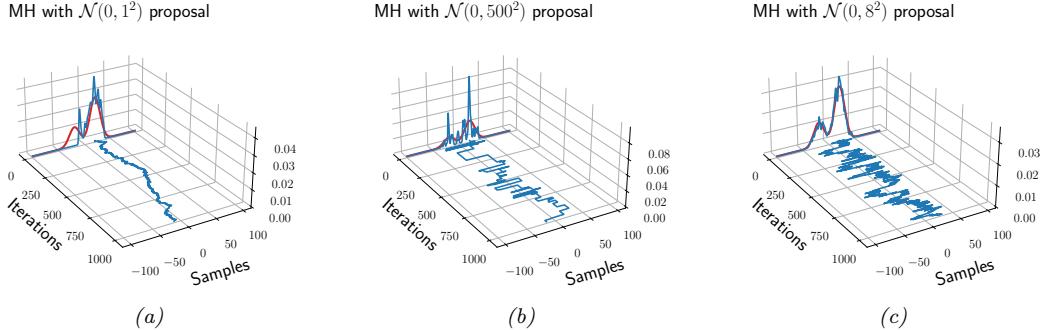


Figure 12.1: An example of the Metropolis Hastings algorithm for sampling from a mixture of two 1D Gaussians ( $\mu = (-20, 20)$ ,  $\pi = (0.3, 0.7)$ ,  $\Sigma = (100, 100)$ ), using a Gaussian proposal with standard deviation of  $\tau \in \{1, 8, 500\}$ . (a) When  $\tau = 1$ , the chain gets trapped near the starting state and fails to sample from the mode at  $\mu = -20$ . (b) When  $\tau = 500$ , the chain is very “sticky”, so its effective sample size is low (as reflected by the rough histogram approximation at the end). (c) Using a variance of  $\tau = 8$  is just right and leads to a good approximation of the true distribution (shown in red). Compare to Figure 12.4. Generated by `mcmc_gmm_demo.ipynb`.

Here  $\tau$  is a scale factor chosen to facilitate rapid mixing. [RR01b] prove that, if the posterior is Gaussian, the asymptotically optimal value is to use  $\tau^2 = 2.38^2/D$ , where  $D$  is the dimensionality of  $\mathbf{x}$ ; this results in an acceptance rate of 0.234, which (in this case) is the optimal tradeoff between exploring widely enough to cover the distribution without being rejected too often. (See [Béd08] for a more recent account of optimal acceptance rates for random walk Metropolis methods.)

Figure 12.1 shows an example where we use RWM to sample from a mixture of two 1D Gaussians. This is a somewhat tricky target distribution, since it consists of two somewhat separated modes. It is very important to set the variance of the proposal  $\tau^2$  correctly: If the variance is too low, the chain will only explore one of the modes, as shown in Figure 12.1(a), but if the variance is too large, most of the moves will be rejected, and the chain will be very **sticky**, i.e., it will stay in the same state for a long time. This is evident from the long stretches of repeated values in Figure 12.1(b). If we set the proposal’s variance just right, we get the trace in Figure 12.1(c), where the samples clearly explore the support of the target distribution.

### 12.2.3.3 Composing proposals

If there are several proposals that might be useful, one can combine them using a **mixture proposal**, which is a convex combination of base proposals:

$$q(\mathbf{x}'|\mathbf{x}) = \sum_{k=1}^K w_k q_k(\mathbf{x}'|\mathbf{x}) \quad (12.16)$$

where  $w_k$  are the mixing weights that sum to one. As long as each  $q_k$  is an individually valid proposal, and each  $w_k > 0$ , then the overall mixture proposal will also be valid. In particular, if each proposal is reversible, so it satisfies detailed balance (Section 2.6.4.4), then so does the mixture.

1 It is also possible to compose individual proposals by chaining them together to get  
 2

3

$$4 q(\mathbf{x}'|\mathbf{x}) = \sum_{\mathbf{x}_1} \cdots \sum_{\mathbf{x}_{K-1}} q_1(\mathbf{x}_1|\mathbf{x}) q_2(\mathbf{x}_2|\mathbf{x}_1) \cdots q_K(\mathbf{x}'|\mathbf{x}_{K-1}) \quad (12.17)$$

5

6 A common example is where each base proposal only updates a subset of the variables (see e.g.,  
 7 Section 12.3).  
 8

9

#### 10 12.2.3.4 Data-driven MCMC

11

12 In the case where the target distribution is a posterior,  $p^*(\mathbf{x}) = p(\mathbf{x}|\mathcal{D})$ , it is helpful to condition the  
 13 proposal not just on the previous hidden state, but also the visible data, i.e., to use  $q(\mathbf{x}'|\mathbf{x}, \mathcal{D})$ . This  
 14 is called **data-driven MCMC** (see e.g., [TZ02; Jih+12]).

15 One way to create such a proposal is to train a recognition network to propose states using  
 16  $q(\mathbf{x}'|\mathbf{x}, \mathcal{D}) = f(\mathbf{x})$ . If the state space is high-dimensional, it might be hard to predict all the hidden  
 17 components, so we can alternatively train individual “experts” to predict specific pieces of the hidden  
 18 state. For example, in the context of estimating the 3d pose of a person from an image, we might  
 19 combine a face detector with a limb detector. We can then use a mixture proposal of the form

20

$$21 q(\mathbf{x}'|\mathbf{x}, \mathcal{D}) = \pi_0 q_0(\mathbf{x}'|\mathbf{x}) + \sum_k \pi_k q_k(x'_k|f_k(\mathcal{D})) \quad (12.18)$$

22

23 where  $q_0$  is a standard data-independent proposal (e.g., random walk), and  $q_k$  updates the  $k$ 'th  
 24 component of the state space.  
 25

26 The overall procedure is a form of **generate and test**: the discriminative proposals  $q(\mathbf{x}'|\mathbf{x}, \mathcal{D})$   
 27 generate new hypotheses, which are then “tested” by computing the posterior ratio  $\frac{p(\mathbf{x}'|\mathcal{D})}{p(\mathbf{x}|\mathcal{D})}$ , to see if  
 28 the new hypothesis is better or worse. (See also Section 13.3, where we discuss learning proposal  
 29 distributions for particle filters.)

30

#### 31 12.2.3.5 Adaptive MCMC

32

33 One can change the parameters of the proposal as the algorithm is running to increase efficiency.  
 34 This is called **adaptive MCMC**. This allows one to start with a broad covariance (say), allowing  
 35 large moves through the space until a mode is found, followed by a narrowing of the covariance to  
 36 ensure careful exploration of the region around the mode.

37 However, one must be careful not to violate the Markov property; thus the parameters of the  
 38 proposal should not depend on the entire history of the chain. It turns out that a sufficient condition  
 39 to ensure this is that the adaption is “faded out” gradually over time. See e.g., [AT08] for details.

40

#### 41 12.2.4 Initialization

42

43 It is necessary to start MCMC in an initial state that has non-zero probability. A natural approach  
 44 is to first use an optimizer to find a local mode. However, at such points the gradients of the log  
 45 joint are zero, which can cause problems for some gradient-based MCMC methods, such as HMC  
 46 (Section 12.5), so it can be better to start “close” to a MAP estimate (see e.g., [HFM17, Sec 7.]).

1 **12.3 Gibbs sampling**

2 The major problems with MH are the need to choose the proposal distribution, and the fact that the  
3 acceptance rate may be low. In this section, we describe an MH method that exploits conditional  
4 independence properties of a graphical model to automatically create a good proposal, with acceptance  
5 probability 1. This method is known as **Gibbs sampling**.<sup>2</sup> (In physics, this method is known as  
6 **Glauber dynamics** or the **heat bath** method.) This is the MCMC analog of coordinate descent.<sup>3</sup>

7  
8 **12.3.1 Basic idea**

9 The idea behind Gibbs sampling is to sample each variable in turn, conditioned on the values of all  
10 the other variables in the distribution. For example, if we have  $D = 3$  variables, we use

- 11 •  $x_1^{s+1} \sim p(x_1|x_2^s, x_3^s)$
- 12 •  $x_2^{s+1} \sim p(x_2|x_1^{s+1}, x_3^s)$
- 13 •  $x_3^{s+1} \sim p(x_3|x_1^{s+1}, x_2^{s+1})$

14 This readily generalizes to  $D$  variables. (Note that if  $x_i$  is a known variable, we do not sample it, but  
15 it may be used as input to the another conditional distribution.)

16 The expression  $p(x_i|\mathbf{x}_{-i})$  is called the **full conditional** for variable  $i$ . In general,  $x_i$  may only  
17 depend on some of the other variables. If we represent  $p(\mathbf{x})$  as a graphical model, we can infer  
18 the dependencies by looking at  $i$ 's Markov blanket, which are its neighbors in the graph (see  
19 Section 4.2.4.3), so we can write

$$\underline{26} \quad x_i^{s+1} \sim p(x_i|\mathbf{x}_{-i}) = p(x_i|\mathbf{x}_{\text{mb}(i)}) \quad (12.19)$$

27 (Compare to the equation for mean field variational inference in Equation (10.28).)

28 We can sample some of the nodes in parallel, without affecting correctness. In particular, suppose  
29 we can create a **coloring** of the (moralized) undirected graph, such that no two neighboring nodes  
30 have the same color. (In general, computing an optimal coloring is NP-complete, but we can use  
31 efficient heuristics such as those in [Kub04].) Then we can sample all the nodes of the same color in  
32 parallel, and cycle through the colors sequentially [Gon+11].

33  
34 **12.3.2 Gibbs sampling is a special case of MH**

35 It turns out that Gibbs sampling is a special case of MH where we use a sequence of proposals of the  
36 form

$$\underline{39} \quad q_i(\mathbf{x}'|\mathbf{x}) = p(x'_i|\mathbf{x}_{-i})\mathbb{I}(\mathbf{x}'_{-i} = \mathbf{x}_{-i}) \quad (12.20)$$

40 That is, we move to a new state where  $x_i$  is sampled from its full conditional, but  $\mathbf{x}_{-i}$  is left  
41 unchanged.

42 2. Josiah Willard Gibbs, 1839–1903, was an American physicist.

43 3. Several software libraries exist for applying Gibbs sampling to general graphical models, including [Nimble](#), which is  
44 a C++ library with an R wrapper, and which replaces older programs such as BUGS and JAGS.

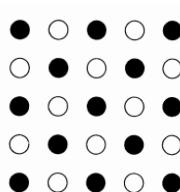


Figure 12.2: Illustration of checkerboard pattern for a 2d MRF. This allows for parallel updates.

We now prove that the acceptance rate of each such proposal is 100%, so the overall algorithm also has an acceptance rate of 100%. We have

$$\alpha = \frac{p(\mathbf{x}')q_i(\mathbf{x}|\mathbf{x}')}{p(\mathbf{x})q_i(\mathbf{x}'|\mathbf{x})} = \frac{p(x'_i|\mathbf{x}'_{-i})p(\mathbf{x}'_{-i})p(x_i|\mathbf{x}'_{-i})}{p(x_i|\mathbf{x}_{-i})p(\mathbf{x}_{-i})p(x'_i|\mathbf{x}_{-i})} \quad (12.21)$$

$$= \frac{p(x'_i|\mathbf{x}_{-i})p(\mathbf{x}_{-i})p(x_i|\mathbf{x}_{-i})}{p(x_i|\mathbf{x}_{-i})p(\mathbf{x}_{-i})p(x'_i|\mathbf{x}_{-i})} = 1 \quad (12.22)$$

where we exploited the fact that  $\mathbf{x}'_{-i} = \mathbf{x}_{-i}$ .

The fact that the acceptance rate is 100% does not necessarily mean that Gibbs will converge rapidly, since it only updates one coordinate at a time (see Section 12.3.7). However, if we can group together correlated variables, then we can sample them as a group, which can significantly help mixing.

### 12.3.3 Example: Gibbs sampling for Ising models

In Section 4.3.2.1, we discuss Ising models and Potts models, which are pairwise MRFs with a 2d grid structure. The joint distribution has the form

$$p(\mathbf{x}) = \frac{1}{Z} \prod_{i \sim j} \psi_{ij}(x_i, x_j | \boldsymbol{\theta}) \quad (12.23)$$

where  $i \sim j$  means  $i$  and  $j$  are neighbors in the graph.

To apply Gibbs sampling to such a model, we just need to iteratively sample from each full conditional:

$$p(x_i | \mathbf{x}_{-i}) \propto \prod_{j \in \text{nbr}(i)} \psi_{ij}(x_i, x_j) \quad (12.24)$$

Note that although Gibbs sampling is a sequential algorithm, we can sometimes exploit conditional independence properties to perform parallel updates [RS97a]. In the case of a 2d grid, we can color code nodes using a checkerboard pattern shown in Figure 12.2. This has the property that the black nodes are conditionally independent of each other given the white nodes, and vice versa. Hence we can sample all the black nodes in parallel (as a single group), and then sample all the white nodes, etc.

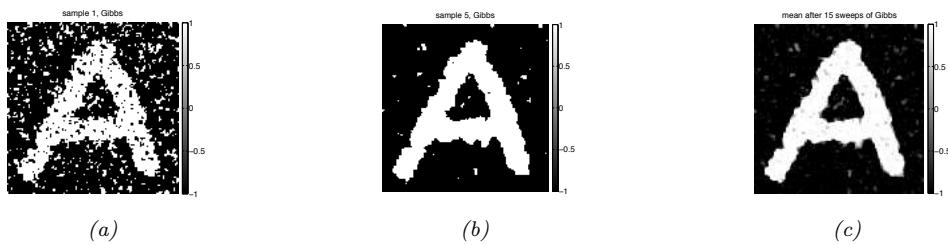


Figure 12.3: Example of image denoising using Gibbs sampling. We use an Ising prior with  $J = 1$  and a Gaussian noise model with  $\sigma = 2$ . (a) Sample from the posterior after one sweep over the image. (b) Sample after 5 sweeps. (c) Posterior mean, computed by averaging over 15 sweeps. Compare with Figure 10.3 which shows the results of mean field inference. Generated by [ising\\_image\\_denoise\\_demo.ipynb](#).

To perform the sampling, we need to compute the full conditional in Equation (12.24). In the case of an Ising model with edge potentials  $\psi(x_i, x_j) = \exp(Jx_i x_j)$ , where  $x_i \in \{-1, +1\}$ , the full conditional becomes

$$p(x_i = +1 | \mathbf{x}_{-i}) = \frac{\prod_{j \in \text{nbr}(i)} \psi_{ij}(x_i = +1, x_j)}{\prod_{j \in \text{nbr}(i)} \psi(x_i = +1, x_j) + \prod_{j \in \text{nbr}(i)} \psi(x_i = -1, x_j)} \quad (12.25)$$

$$= \frac{\exp[J \sum_{j \in \text{nbr}(i)} x_j]}{\exp[J \sum_{j \in \text{nbr}(i)} x_j] + \exp[-J \sum_{j \in \text{nbr}(i)} x_j]} \quad (12.26)$$

$$= \frac{\exp[J\eta_i]}{\exp[J\eta_i] + \exp[-J\eta_i]} = \sigma(2J\eta_i) \quad (12.27)$$

where  $J$  is the coupling strength,  $\eta_i \triangleq \sum_{j \in \text{nbr}(i)} x_j$  and  $\sigma(u) = 1/(1 + e^{-u})$  is the sigmoid function. (If we use  $x_i \in \{0, 1\}$ , this becomes  $p(x_i = +1 | \mathbf{x}_{-i}) = \sigma(J\eta_i)$ .) It is easy to see that  $\eta_i = x_i(a_i - d_i)$ , where  $a_i$  is the number of neighbors that agree with (have the same sign as) node  $i$ , and  $d_i$  is the number of neighbors who disagree. If this number is equal, the ‘‘forces’’ on  $x_i$  cancel out, so the full conditional is uniform. Some samples from this model are shown in Figure 4.17.

One application of Ising models is as a prior for binary image denoising problems. In particular, suppose  $\mathbf{y}$  is a noisy version of  $\mathbf{x}$ , and we wish to compute the posterior  $p(\mathbf{x} | \mathbf{y}) \propto p(\mathbf{x})p(\mathbf{y} | \mathbf{x})$ , where  $p(\mathbf{x})$  is an Ising prior, and  $p(\mathbf{y} | \mathbf{x}) = \prod_i p(y_i | x_i)$  is a per-site likelihood term. Suppose this is a Gaussian. Let  $\psi_i(x_i) = \mathcal{N}(y_i | x_i, \sigma^2)$  be the corresponding ‘‘local evidence’’ term. The full conditional becomes

$$p(x_i = +1 | \mathbf{x}_{-i}, \mathbf{y}) = \frac{\exp[J\eta_i]\psi_i(+1)}{\exp[J\eta_i]\psi_i(+1) + \exp[-J\eta_i]\psi_i(-1)} \quad (12.28)$$

$$= \sigma\left(2J\eta_i - \log \frac{\psi_i(+1)}{\psi_i(-1)}\right) \quad (12.29)$$

Now the probability of  $x_i$  entering each state is determined both by compatibility with its neighbors (the Ising prior) and compatibility with the data (the local likelihood term).

See Figure 12.3 for an example of this algorithm applied to a simple image denoising problem. The results are similar to the mean field results in Figure 10.3.

---

1 **12.3.4 Example: Gibbs sampling for Potts models**

3 We can extend Section 12.3.3 to the Potts models as follows. Recall that the model has the following  
4 form:

6  $p(\mathbf{x}) = \frac{1}{Z} \exp(-\mathcal{E}(\mathbf{x}))$  (12.30)

8  $\mathcal{E}(\mathbf{x}) = -J \sum_{i \sim j} \mathbb{I}(x_i = x_j)$  (12.31)

11 For a node  $i$  with neighbors  $\text{nbr}(i)$ , the full conditional is thus given by

13  $p(x_i = k | \mathbf{x}_{-i}) = \frac{\exp(J \sum_{n \in \text{nbr}(i)} \mathbb{I}(x_n = k))}{\sum_{k'} \exp(J \sum_{n \in \text{nbr}(i)} \mathbb{I}(x_n = k'))}$  (12.32)

16 So if  $J > 0$ , a node  $i$  is more likely to enter a state  $k$  if most of its neighbors are already in state  $k$ ,  
17 corresponding to an attractive MRF. If  $J < 0$ , a node  $i$  is more likely to enter a different state from  
18 its neighbors, corresponding to a repulsive MRF. See Figure 4.18 for some samples from this model  
19 created using this method.

20

21 **12.3.5 Example: Gibbs sampling for GMMs**

23 In this section, we consider sampling from a Bayesian Gaussian mixture model of the form

25  $p(z = k, \mathbf{x} | \boldsymbol{\theta}) = \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \Sigma_k)$  (12.33)

27  $p(\boldsymbol{\theta}) = \text{Dir}(\boldsymbol{\pi} | \boldsymbol{\alpha}) \prod_{k=1}^K \mathcal{N}(\boldsymbol{\mu}_k | \mathbf{m}_0, \mathbf{V}_0) \text{IW}(\boldsymbol{\Sigma}_k, \mathbf{S}_0, \nu_0)$  (12.34)

30 **12.3.5.1 Known parameters**

32 Suppose, initially, that the parameters  $\boldsymbol{\theta}$  are known. We can easily draw independent samples from  
33  $p(\mathbf{x} | \boldsymbol{\theta})$  by using ancestral sampling: first sample  $z$  and then  $\mathbf{x}$ . However, for illustrative purposes, we  
34 will use Gibbs sampling to draw correlated samples. The full conditional for  $p(\mathbf{x} | z = k, \boldsymbol{\theta})$  is just  
35  $\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \Sigma_k)$ , and the full conditional for  $p(z = k | \mathbf{x})$  is given by Bayes rule:

37  $p(z = k | \mathbf{x}, \boldsymbol{\theta}) = \frac{\pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \Sigma_k)}{\sum_{k'} \pi_{k'} \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_{k'}, \Sigma_{k'})}$  (12.35)

40 An example of this procedure, applied to a mixture of two 1D Gaussians with means at  $-20$  and  
41  $+20$ , is shown in Figure 12.4. We see that the samples are auto correlated, meaning that if we are  
42 in state 1, we will likely stay in that state for a while, and generate values near  $\mu_1$ ; then we will  
43 stochastically jump to state 2, and stay near there for a while, etc. (See Section 12.6.3 for a way to  
44 measure this.) By contrast, independent samples from the joint would not be correlated at all.

45 In Section 12.3.5.2, we modify this example to sample the parameters of the GMM from their  
46 posterior,  $p(\boldsymbol{\theta} | \mathcal{D})$ , instead of sampling from  $p(\mathcal{D} | \boldsymbol{\theta})$ .

47

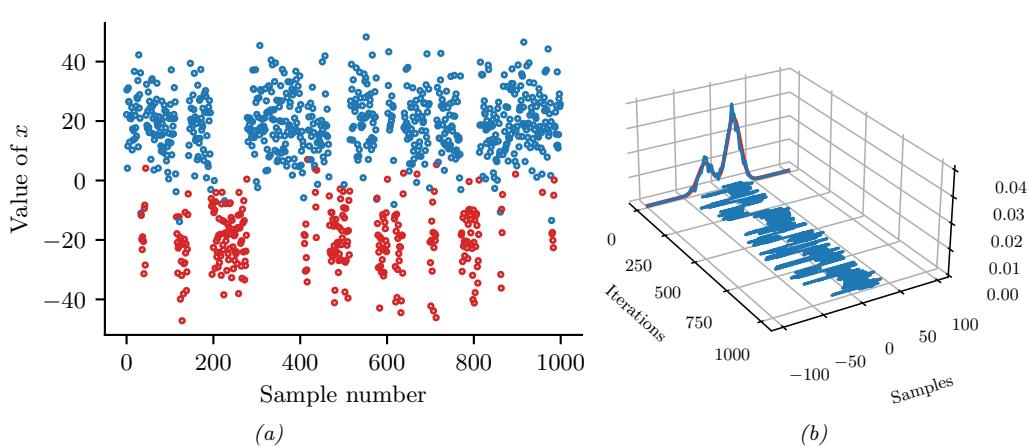


Figure 12.4: (a) Some samples from a mixture of two 1d Gaussians generated using Gibbs sampling. Color denotes the value of  $z$ , vertical location denotes the value of  $x$ . Horizontal axis represents time (sample number). (b) Traceplot of  $x$  over time, and the resulting empirical distribution is shown in blue. The true distribution is shown in red. Compare to Figure 12.1. Generated by [mcmc\\_gmm.ipynb](#).

### 12.3.5.2 Unknown parameters

Now suppose the parameters are unknown, so we want to fit the model to data. If we use a conditionally conjugate factored prior, then the full joint distribution is given by

$$p(\mathbf{x}, \mathbf{z}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = p(\mathbf{x}|\mathbf{z}, \boldsymbol{\mu}, \boldsymbol{\Sigma})p(\mathbf{z}|\boldsymbol{\pi})p(\boldsymbol{\pi}) \prod_{k=1}^K p(\boldsymbol{\mu}_k)p(\boldsymbol{\Sigma}_k) \quad (12.36)$$

$$= \left( \prod_{i=1}^N \prod_{k=1}^K (\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k))^{\mathbb{I}(z_i=k)} \right) \times \quad (12.37)$$

$$\text{Dir}(\boldsymbol{\pi}|\boldsymbol{\alpha}) \prod_{k=1}^K \mathcal{N}(\boldsymbol{\mu}_k | \mathbf{m}_0, \mathbf{V}_0) \text{IW}(\boldsymbol{\Sigma}_k | \mathbf{S}_0, \nu_0) \quad (12.38)$$

We use the same prior for each mixture component.

The full conditionals are as follows. For the discrete indicators, we have

$$p(z_i = k | \mathbf{x}_i, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) \propto \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (12.39)$$

For the mixing weights, we have (using results from Section 3.2.2)

$$p(\boldsymbol{\pi} | \mathbf{z}) = \text{Dir}(\{\alpha_k + \sum_{i=1}^N \mathbb{I}(z_i = k)\}_{k=1}^K) \quad (12.40)$$

1 For the means, we have (using results from Section 3.3.1)

$$\underline{3} \quad p(\boldsymbol{\mu}_k | \boldsymbol{\Sigma}_k, \mathbf{z}, \mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_k | \mathbf{m}_k, \mathbf{V}_k) \quad (12.41)$$

$$\underline{4} \quad \mathbf{V}_k^{-1} = \mathbf{V}_0^{-1} + N_k \boldsymbol{\Sigma}_k^{-1} \quad (12.42)$$

$$\underline{5} \quad \mathbf{m}_k = \mathbf{V}_k (\boldsymbol{\Sigma}_k^{-1} N_k \bar{\mathbf{x}}_k + \mathbf{V}_0^{-1} \mathbf{m}_0) \quad (12.43)$$

$$\underline{6} \quad N_k \triangleq \sum_{i=1}^N \mathbb{I}(z_i = k) \quad (12.44)$$

$$\underline{7} \quad \bar{\mathbf{x}}_k \triangleq \frac{\sum_{i=1}^N \mathbb{I}(z_i = k) \mathbf{x}_i}{N_k} \quad (12.45)$$

8 For the covariances, we have (using results from Section 3.3.2)

$$\underline{9} \quad p(\boldsymbol{\Sigma}_k | \boldsymbol{\mu}_k, \mathbf{z}, \mathbf{x}) = \text{IW}(\boldsymbol{\Sigma}_k | \mathbf{S}_k, \nu_k) \quad (12.46)$$

$$\underline{10} \quad \mathbf{S}_k = \mathbf{S}_0 + \sum_{i=1}^N \mathbb{I}(z_i = k) (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \quad (12.47)$$

$$\underline{11} \quad \nu_k = \nu_0 + N_k \quad (12.48)$$

### 12 12.3.6 Metropolis within Gibbs

23 When implementing Gibbs sampling, we have to sample from the full conditionals. If the distributions  
24 are conjugate, we can compute the full conditional in closed form, but in the general case, we will  
25 need to devise special algorithms to sample from the full conditionals.

26 One approach is to use the MH algorithm; this is called **Metropolis within Gibbs**. In particular,  
27 to sample from  $x_i^{s+1} \sim p(x_i | \mathbf{x}_{1:i-1}^{s+1}, \mathbf{x}_{i+1:D}^s)$ , we proceed in 3 steps:

29 1. Propose  $x'_i \sim q(x'_i | x_i^s)$

30 2. Compute the acceptance probability  $A_i = \min(1, \alpha_i)$  where

$$\underline{32} \quad \alpha_i = \frac{p(\mathbf{x}_{1:i-1}^{s+1}, x'_i, \mathbf{x}_{i+1:D}^s) / q(x'_i | x_i^s)}{p(\mathbf{x}_{1:i-1}^s, x_i^s, \mathbf{x}_{i+1:D}^s) / q(x_i^s | x'_i)} \quad (12.49)$$

36 3. Sample  $u \sim U(0, 1)$  and set  $x_i^{s+1} = x'_i$  if  $u < A_i$ , and set  $x_i^{s+1} = x_i^s$  otherwise.

### 38 12.3.7 Blocked Gibbs sampling

40 Gibbs sampling can be quite slow, since it only updates one variable at a time (so-called **single**  
41 **site updating**). If the variables are highly correlated, the chain will move slowly through the state  
42 space. This is illustrated in Figure 12.5, where we illustrate sampling from a 2d Gaussian. The ellipse  
43 represents the covariance matrix. The size of the moves taken by Gibbs sampling is controlled by the  
44 variance of the conditional distributions. If the variance is  $\ell$  along some coordinate direction, but the  
45 support of the distribution is  $L$  along this dimension, then we need  $O((L/\ell)^2)$  steps to obtain an  
46 independent sample.

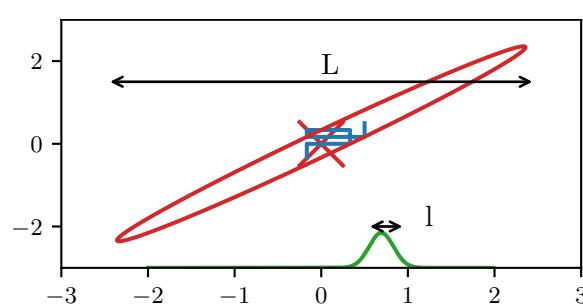


Figure 12.5: Illustration of potentially slow sampling when using Gibbs sampling for a skewed 2D Gaussian. Adapted from Figure 11.11 of [Bis06]. Generated by [gibbs\\_gauss\\_demo.ipynb](#).

In some cases we can efficiently sample groups of variables at a time. This is called **blocked Gibbs sampling** [JKK95; WY02], and can make much bigger moves through the state space.

As an example, suppose we want to perform Bayesian inference for a state-space model, such as an HMM, i.e., we want to sample from

$$p(\boldsymbol{\theta}, \mathbf{z} | \mathbf{x}) \propto p(\boldsymbol{\theta}) \prod_{t=1}^T p(\mathbf{x}_t | \mathbf{z}_t, \boldsymbol{\theta}) p(\mathbf{z}_t | \mathbf{z}_{t-1}, \boldsymbol{\theta}) \quad (12.50)$$

We can use blocked Gibbs sampling, where we alternate between sampling from  $p(\boldsymbol{\theta} | \mathbf{z}, \mathbf{x})$  and  $p(\mathbf{z} | \mathbf{x}, \boldsymbol{\theta})$ ; The former is easy to do (assuming conjugate priors), since all variables in the model are observed (see Section 29.8.4.1). The latter can be done using the forwards-filtering backwards-sampling (Section 8.2.8).

### 12.3.8 Collapsed Gibbs sampling

We can sometimes gain even greater speedups by analytically integrating out some of the unknown quantities. This is called a **collapsed Gibbs sampler**, and it tends to be more efficient, since it is sampling in a lower dimensional space. This can result in lower variance, as discussed in Section 11.6.2.

As an example, consider a GMM with a fully conjugate prior. This can be represented as a DPGM as shown in Figure 12.6a. Since the prior is conjugate, we can analytically integrate out the model parameters  $\boldsymbol{\mu}_k$ ,  $\boldsymbol{\Sigma}_k$  and  $\boldsymbol{\pi}$ , so the only remaining hidden variables are the discrete indicator variables  $\mathbf{z}$ . However, once we integrate out  $\boldsymbol{\pi}$ , all the  $z_i$  nodes become inter-dependent. Similarly, once we integrate out  $\boldsymbol{\theta}_k = (\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ , all the  $\mathbf{x}_i$  nodes become inter-dependent, as shown in Figure 12.6b. Nevertheless, we can easily compute the full conditionals, and hence implement a Gibbs sampler, as we explain below. In particular, the full conditional for the latent indicators is given by

$$p(z_i = k | \mathbf{z}_{-i}, \mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta}) \propto p(z_i = k | \mathbf{z}_{-i}, \boldsymbol{\alpha}, \boldsymbol{\beta}) p(\mathbf{x}|z_i = k, \mathbf{z}_{-i}, \boldsymbol{\alpha}, \boldsymbol{\beta}) \quad (12.51)$$

$$\propto p(z_i = k | \mathbf{z}_{-i}, \boldsymbol{\alpha}) p(\mathbf{x}_i | \mathbf{x}_{-i}, z_i = k, \mathbf{z}_{-i}, \boldsymbol{\beta}) \quad (12.52)$$

$$\propto p(z_i = k | \mathbf{z}_{-i}, \boldsymbol{\alpha}) p(\mathbf{x}_i | \mathbf{x}_{-i}, z_i = k, \mathbf{z}_{-i}, \boldsymbol{\beta}) \quad (12.53)$$

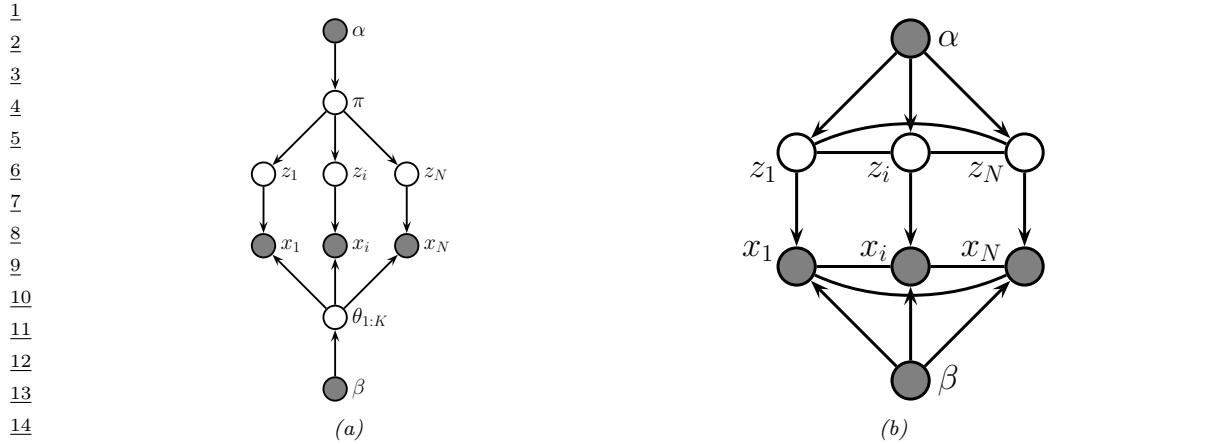


Figure 12.6: (a) A mixture model represented as an “unrolled” DPGM. (b) After integrating out the continuous latent parameters.

where  $\beta = (\mathbf{m}_0, \mathbf{V}_0, \mathbf{S}_0, \nu_0)$  are the hyper-parameters for the class-conditional densities. We now discuss how to compute these terms.

Suppose we use a symmetric prior of the form  $\pi \sim \text{Dir}(\alpha)$ , where  $\alpha_k = \alpha/K$ , for the mixing weights. Then we can obtain the first term in Equation (12.53), from Equation (3.29), where

$$p(z_1, \dots, z_N | \alpha) = \frac{\Gamma(\alpha)}{\Gamma(N + \alpha)} \prod_{k=1}^K \frac{\Gamma(N_k + \alpha/K)}{\Gamma(\alpha/K)} \quad (12.54)$$

Hence

$$p(z_i = k | \mathbf{z}_{-i}, \alpha) = \frac{p(\mathbf{z}_{1:N} | \alpha)}{p(\mathbf{z}_{-i} | \alpha)} = \frac{\frac{1}{\Gamma(N + \alpha)}}{\frac{1}{\Gamma(N + \alpha - 1)}} \times \frac{\Gamma(N_k + \alpha/K)}{\Gamma(N_{k,-i} + \alpha/K)} \quad (12.55)$$

$$= \frac{\Gamma(N + \alpha - 1)}{\Gamma(N + \alpha)} \frac{\Gamma(N_{k,-i} + 1 + \alpha/K)}{\Gamma(N_{k,-i} + \alpha/K)} = \frac{N_{k,-i} + \alpha}{N + \alpha - 1} \quad (12.56)$$

where  $N_{k,-i} \triangleq \sum_{n \neq i} \mathbb{I}(z_n = k) = N_k - 1$ , and where we exploited the fact that  $\Gamma(x + 1) = x\Gamma(x)$ .

To obtain the second term in Equation (12.53), which is the posterior predictive distribution for  $\mathbf{x}_i$  given all the other data and all the assignments, we use the fact that

$$p(\mathbf{x}_i | \mathbf{x}_{-i}, \mathbf{z}_{-i}, z_i = k, \beta) = p(\mathbf{x}_i | \mathcal{D}_{-i,k}, \beta) \quad (12.57)$$

where  $\mathcal{D}_{-i,k} = \{\mathbf{x}_j : z_j = k, j \neq i\}$  is all the data assigned to cluster  $k$  except for  $\mathbf{x}_i$ . If we use a conjugate prior for  $\theta_k$ , we can compute  $p(\mathbf{x}_i | \mathcal{D}_{-i,k}, \beta)$  in closed form. Furthermore, we can efficiently update these predictive likelihoods by caching the sufficient statistics for each cluster. To compute the above expression, we remove  $\mathbf{x}_i$ ’s statistics from its current cluster (namely  $z_i$ ), and then evaluate  $\mathbf{x}_i$  under each cluster’s posterior predictive distribution. Once we have picked a new cluster, we add  $\mathbf{x}_i$ ’s statistics to this new cluster.

47

Some pseudo-code for one step of the algorithm is shown in Algorithm 26, based on [Sud06, p94]. (We update the nodes in random order to improve the mixing time, as suggested in [RS97b].) We can initialize the sample by sequentially sampling from  $p(z_i|\mathbf{z}_{1:i-1}, \mathbf{x}_{1:i})$ . In the case of GMMs, both the naive sampler and collapsed sampler take  $O(NKD)$  time per step.

---

**Algorithm 26:** Collapsed Gibbs sampler for a mixture model

---

```

1 for each  $i = 1 : N$  in random order do
2   Remove  $\mathbf{x}_i$ 's sufficient statistics from old cluster  $z_i$ 
3   for each  $k = 1 : K$  do
4      $\lfloor$  Compute  $p_k(\mathbf{x}_i|\boldsymbol{\beta}) = p(\mathbf{x}_i|\{\mathbf{x}_j : z_j = k, j \neq i\}, \boldsymbol{\beta})$ 
5   Compute  $p(z_i = k|\mathbf{z}_{-i}, \alpha) \propto (N_{k,-i} + \alpha/K)p_k(\mathbf{x}_i)$ 
6   Sample  $z_i \sim p(z_i|\cdot)$ 
7   Add  $\mathbf{x}_i$ 's sufficient statistics to new cluster  $z_i$ 

```

---

The primary advantage of using the collapsed sampler is that it extends to the case where we have an “infinite” number of mixture components, as in the Dirichlet process mixture model of Section 31.3.2.

## 12.4 Auxiliary variable MCMC

Sometimes we can dramatically improve the efficiency of sampling by introducing **auxiliary variables**, in order to reduce correlation between the original variables. If the original variables are denoted by  $\mathbf{x}$ , and the auxiliary variables by  $\mathbf{v}$ , then the augmented distribution becomes  $p(\mathbf{x}, \mathbf{v})$ . We assume it is easier to sample from this than the marginal distribution  $p(\mathbf{x})$ . If so, we can draw joint samples  $(\mathbf{x}^s, \mathbf{v}^s) \sim p(\mathbf{x}, \mathbf{v})$ , and then just “throw away” the  $\mathbf{v}^s$ , and the result will be samples from the desired marginal,  $\mathbf{x}^s \sim \sum_{\mathbf{v}} p(\mathbf{x}, \mathbf{v})$ . We give some examples of this below.

### 12.4.1 Slice sampling

Consider sampling from a univariate, but multimodal, distribution  $p(x) = \tilde{p}(x)/Z_p$ , where  $\tilde{p}(x)$  is unnormalized, and  $Z_p = \int \tilde{p}(x)dx$ . We can sometimes improve the ability to make large moves by adding a uniform auxiliary variable  $v$ . We define the joint distribution as follows:

$$\hat{p}(x, v) = \begin{cases} 1/Z_p & \text{if } 0 \leq v \leq \tilde{p}(x) \\ 0 & \text{otherwise} \end{cases} \quad (12.58)$$

The marginal distribution over  $x$  is given by

$$\int \hat{p}(x, v)dv = \int_0^{\tilde{p}(x)} \frac{1}{Z_p} dv = \frac{\tilde{p}(x)}{Z_p} = p(x) \quad (12.59)$$

so we can sample from  $p(x)$  by sampling from  $\hat{p}(x, v)$  and then ignoring  $v$ . To do this, we will use a technique called **slice sampling** [Nea03].

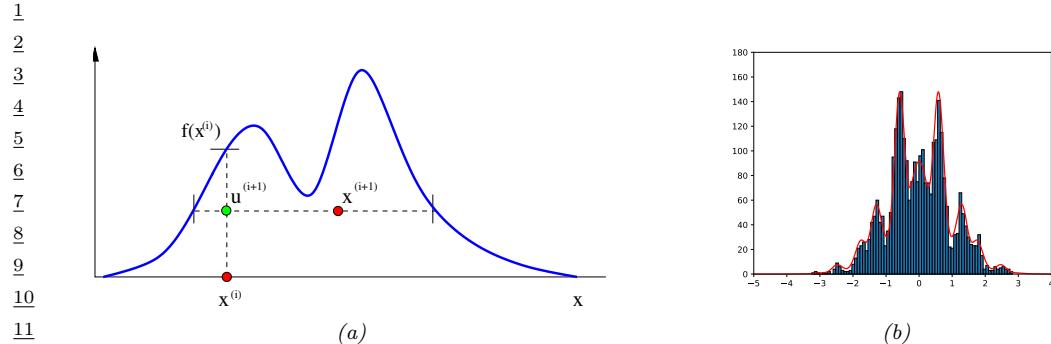


Figure 12.7: Slice sampling. (a) Illustration of one step of the algorithm in 1d. Given a previous sample  $x^i$ , we sample  $u^{i+1}$  uniformly on  $[0, f(x^i)]$ , where  $f = \tilde{p}$  is the (unnormalized) target density. We then sample  $x^{i+1}$  along the slice where  $f(x) \geq u^{i+1}$ . From Figure 15 of [And+03]. Used with kind permission of Nando de Freitas.  
(b) Output of slice sampling applied to a 1d distribution. Generated by [slice\\_sampling\\_demo\\_1d.ipynb](#).

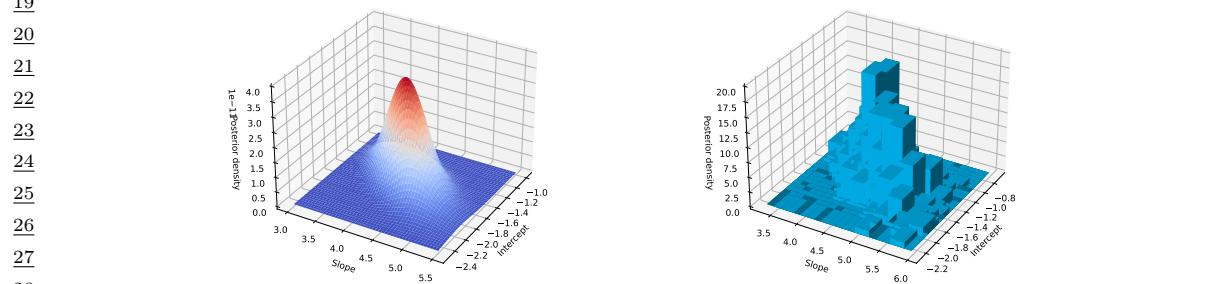


Figure 12.8: Posterior for binomial regression for 1d data. Left: Slice sampling approximation. Right: Grid approximation. Generated by [slice\\_sampling\\_demo\\_2d.ipynb](#).

This works as follows. Given previous sample  $x^i$ , we sample  $v^{i+1}$  from

$$p(v|x^i) = U_{[0,\tilde{p}(x^i)]}(v) \quad (12.60)$$

This amounts to uniformly picking a point on the vertical line between 0 and  $\tilde{p}(x^i)$ . We use this to construct a “slice” of the density at or above this height, by computing  $A^{i+1} = \{x : \tilde{p}(x) \geq v^{i+1}\}$ . We then sample  $x^{i+1}$  uniformly from this set. See Figure 12.7(a) for an illustration.

To compute the level set  $A$ , we can use an iterative search procedure called **stepping out**, in which we start with an interval  $x_{min} \leq x \leq x_{max}$  around the current point  $x^i$  of some width, and then we keep extending it until the endpoints fall outside the slice. We can then use rejection sampling to sample from the interval. For the details, see [Nea03].

To apply the method to multivariate distributions, we sample one extra auxiliary variable for each dimension. Thus we perform 2D sampling operations to draw a single joint sample, where

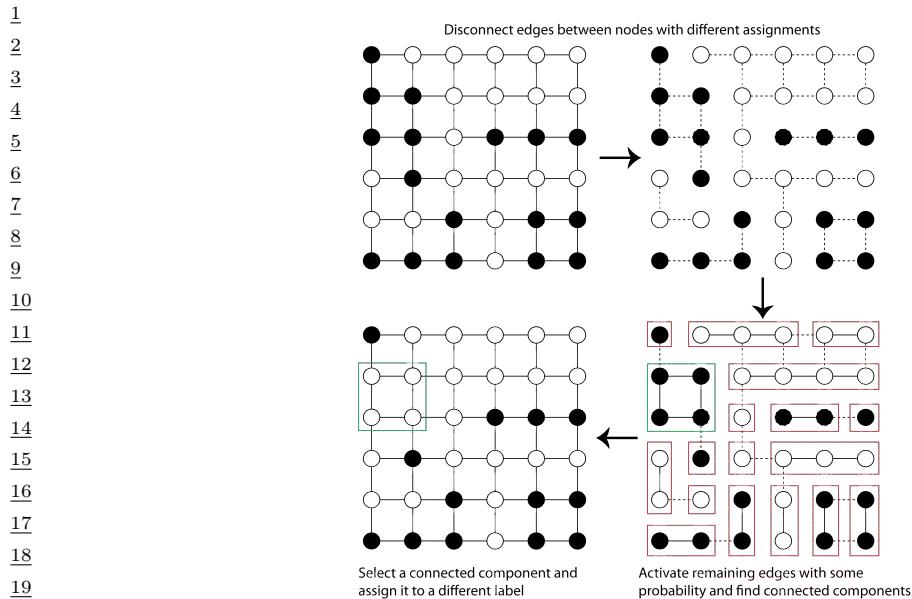


Figure 12.9: Illustration of the Swendsen Wang algorithm on a 2d grid. Used with kind permission of Kevin Tang.

$D$  is the number of random variables. The advantage of this over Gibbs sampling applied to the original (non-augmented) distribution is that it only needs access to the unnormalized joint, not the full-conditionals.

Figure 12.7(b) illustrates the algorithm in action on a synthetic 1d problem. Figure 12.8 illustrates its behavior on a slightly harder problem, namely binomial logistic regression. The model has the form  $y_i \sim \text{Bin}(n_i, \text{logit}(\beta_1 + \beta_2 x_i))$ . We use a vague Gaussian prior for the  $\beta_j$ 's. On the left we show the slice sampling approximation to the posterior, and on the right we shpw a grid-based approximation, as a simple deteterministic proxy for the true posterior. We see a close correpondence.

### 12.4.2 Swendsen Wang

Consider an Ising model of the following form:  $p(\mathbf{x}) = \frac{1}{Z} \prod_e \Psi(\mathbf{x}_e)$ , where  $\mathbf{x}_e = (x_i, x_j)$  for edge  $e = (i, j)$ ,  $x_i \in \{+1, -1\}$ , and the edge potential is defined by  $\begin{pmatrix} e^J & e^{-J} \\ e^{-J} & e^J \end{pmatrix}$ , where  $J$  is the edge strength. In Section 12.3.3, we discussed how to apply Gibbs sampling to this model. However, this can be slow when  $J$  is large in absolute value, because neighboring states can be highly correlated. The **Swendsen Wang** algorithm [SW87b] is an auxiliary variable MCMC sampler which mixes much faster, at least for the case of attractive or ferromagnetic models, with  $J > 0$ .

Suppose we introduce auxiliary binary variables, one per edge.<sup>4</sup> These are called **bond variables**, and will be denoted by  $\mathbf{v}$ . We then define an extended model  $p(\mathbf{x}, \mathbf{v})$  of the form  $p(\mathbf{x}, \mathbf{v}) =$

4. Our presentation of the method is based on notes by David Mackay, available from <http://www.inference.phy.cam.ac.uk/mackay/itila/swendsen.pdf>.

1  $\frac{1}{Z'} \prod_e \Psi(\mathbf{x}_e, v_e)$ , where  $v_e \in \{0, 1\}$ , and we define the new edge potentials as follows:

$$\Psi(\mathbf{x}_e, v_e = 0) = \begin{pmatrix} e^{-J} & e^{-J} \\ e^{-J} & e^{-J} \end{pmatrix}, \quad \Psi(\mathbf{x}_e, v_e = 1) = \begin{pmatrix} e^J - e^{-J} & 0 \\ 0 & e^J - e^{-J} \end{pmatrix} \quad (12.61)$$

2 It is clear that  $\sum_{v_e=0}^1 \Psi(\mathbf{x}_e, v_e) = \Psi(\mathbf{x}_e)$ , and hence that  $\sum_{\mathbf{v}} p(\mathbf{x}, \mathbf{v}) = p(\mathbf{x})$ , as required.

3 Fortunately, it is easy to apply Gibbs sampling to this extended model. The full conditional  $p(\mathbf{v}|\mathbf{x})$   
4 factorizes over the edges, since the bond variables are conditionally independent given the node  
5 variables. Furthermore, the full conditional  $p(v_e|\mathbf{x}_e)$  is simple to compute: if the nodes on either end  
6 of the edge are in the same state ( $x_i = x_j$ ), we set the bond  $v_e$  to 1 with probability  $p = 1 - e^{-2J}$ ,  
7 otherwise we set it to 0. In Figure 12.9 (top right), the bonds that could be turned on (because their  
8 corresponding nodes are in the same state) are represented by dotted edges. In Figure 12.9 (bottom  
9 right), the bonds that are randomly turned on are represented by solid edges.

10 To sample  $p(\mathbf{x}|\mathbf{v})$ , we proceed as follows. Find the connected components defined by the graph  
11 induced by the bonds that are turned on. (Note that a connected component may consist of a  
12 singleton node.) Pick one of these components uniformly at random. All the nodes in each such  
13 component must have the same state. Pick a state  $\pm 1$  uniformly at random, and set all the variables  
14 in this component to adopt this new state. This is illustrated in Figure 12.9 (bottom right), where  
15 the green square denotes the selected connected component; we set all the nodes within this square  
16 to white, to get the bottom left configuration.

17 It should be intuitively clear that Swendsen Wang makes much larger moves through the state space  
18 than Gibbs sampling. The gains are exponentially large for certain settings of the edge parameter.  
19 More precisely, let the edge strength be parameterized by  $J/T$ , where  $T > 0$  is a computational  
20 temperature. For large  $T$ , the nodes are roughly independent, so both methods work equally well.  
21 However, as  $T$  approaches a **critical temperature**  $T_c$ , the typical states of the system have very  
22 long correlation lengths, and Gibbs sampling takes a very long time to generate independent samples.  
23 As the temperature continues to drop, the typical states are either all on or all off. The frequency  
24 with which Gibbs sampling moves between these two modes is exponentially small. By contrast, SW  
25 mixes rapidly at all temperatures.

26 Unfortunately, if any of the edge weights are negative,  $J < 0$ , the system is **frustrated**, and there  
27 are exponentially many modes, even at low temperature. SW does not work very well in this setting,  
28 since it tries to force many neighboring variables to have the same state. In fact, sampling from these  
29 kinds of frustrated systems is provably computationally hard for any algorithm [JS93; JS96].

30

## 31 12.5 Hamiltonian Monte Carlo (HMC)

32 Many MCMC algorithms perform poorly in high dimensional spaces, because they rely on a form  
33 of random search based on local perturbations. In this section, we discuss a method known as  
34 **Hamiltonian Monte Carlo** or **HMC**, that leverages gradient information to guide the local moves.  
35 This is an auxiliary variable method (Section 12.4) derived from physics [Dua+87; Nea93; Mac03;  
36 Nea10; Bet17].<sup>5</sup> In particular, the method builds on **Hamiltonian mechanics**, which we describe  
37 below.

38

39 5. The method was originally called **hybrid MC** [Dua+87]. It was introduced to the statistics community in [Nea93],  
40 and was renamed to Hamiltonian MC in [Mac03].

41

1 **12.5.1 Hamiltonian mechanics**

2 Consider a particle rolling around an energy landscape. We can characterize the motion of the particle  
3 in terms of its position  $\boldsymbol{\theta} \in \mathbb{R}^D$  (often denoted by  $\mathbf{q}$ ) and its momentum  $\mathbf{v} \in \mathbb{R}^D$  (often denoted by  
4  $\mathbf{p}$ ). The set of possible values for  $(\boldsymbol{\theta}, \mathbf{v})$  is called the **phase space**. We define the **Hamiltonian**  
5 function for each point in phase space as follows:

6 
$$\mathcal{H}(\boldsymbol{\theta}, \mathbf{v}) \triangleq \mathcal{E}(\boldsymbol{\theta}) + \mathcal{K}(\mathbf{v}) \quad (12.62)$$

7 where  $\mathcal{E}(\boldsymbol{\theta})$  is the **potential energy**,  $\mathcal{K}(\mathbf{v})$  is the **kinetic energy**, and the Hamiltonian is the total  
8 energy. In a physical setting, the potential energy is due to the pull of gravity, and the momentum is  
9 due to the motion of the particle. In a statistical setting, we often take the potential energy to be

10 
$$\mathcal{E}(\boldsymbol{\theta}) = -\log \tilde{p}(\boldsymbol{\theta}) \quad (12.63)$$

11 where  $\tilde{p}(\boldsymbol{\theta})$  is a possibly unnormalized distribution, such as  $p(\boldsymbol{\theta}, \mathcal{D})$ , and the kinetic energy to be

12 
$$\mathcal{K}(\mathbf{v}) = \frac{1}{2} \mathbf{v}^\top \boldsymbol{\Sigma}^{-1} \mathbf{v} \quad (12.64)$$

13 where  $\boldsymbol{\Sigma}$  is a positive definite matrix, known as the **inverse mass matrix**.

14 Stable orbits are defined by trajectories in phase space that have a constant energy. The trajectory  
15 of a particle within an energy level set can be obtained by solving the following continuous time  
16 differential equations, known as **Hamilton's equations**:

17 
$$\begin{aligned} \frac{d\boldsymbol{\theta}}{dt} &= \frac{\partial \mathcal{H}}{\partial \mathbf{v}} = \frac{\partial \mathcal{K}}{\partial \mathbf{v}} \\ \frac{d\mathbf{v}}{dt} &= -\frac{\partial \mathcal{H}}{\partial \boldsymbol{\theta}} = -\frac{\partial \mathcal{E}}{\partial \boldsymbol{\theta}} \end{aligned} \quad (12.65)$$

18 To see why energy is conserved, note that

19 
$$\frac{d\mathcal{H}}{dt} = \sum_{i=1}^D \left[ \frac{\partial \mathcal{H}}{\partial \boldsymbol{\theta}_i} \frac{d\boldsymbol{\theta}_i}{dt} + \frac{\partial \mathcal{H}}{\partial \mathbf{v}_i} \frac{d\mathbf{v}_i}{dt} \right] = \sum_{i=1}^D \left[ \frac{\partial \mathcal{H}}{\partial \boldsymbol{\theta}_i} \frac{\partial \mathcal{H}}{\partial \mathbf{v}_i} - \frac{\partial \mathcal{H}}{\partial \boldsymbol{\theta}_i} \frac{\partial \mathcal{H}}{\partial \mathbf{v}_i} \right] = 0 \quad (12.66)$$

20 Intuitively, we can understand this result as follows: a satellite in orbit around a planet will “want”  
21 to continue in a straight line due to its momentum, but will get pulled in towards the planet due  
22 to gravity, and if these forces cancel, the orbit is stable. If the satellite starts spiraling towards the  
23 planet, its kinetic energy will increase but its potential energy will decrease.

24 Note that the mapping from  $(\boldsymbol{\theta}(t), \mathbf{v}(t))$  to  $(\boldsymbol{\theta}(t+s), \mathbf{v}(t+s))$  for some time increment  $s$  is invertible  
25 for small enough time steps. Furthermore, this mapping is volume preserving, so has a Jacobian  
26 determinant of 1. (See e.g., [BZ20, p287] for a proof.) These facts will be important later when we  
27 turn this system into an MCMC algorithm.

28 **12.5.2 Integrating Hamilton's equations**

29 In this section, we discuss how to simulate Hamilton's equations in discrete time.

1    **12.5.2.1 Euler's method**

3    The simplest way to model the time evolution is to update the position and momentum simultaneously  
4    by a small amount, known as the step size  $\eta$ :

6     $\mathbf{v}_{t+1} = \mathbf{v}_t + \eta \frac{d\mathbf{v}}{dt}(\boldsymbol{\theta}_t, \mathbf{v}_t) = \mathbf{v}(t) - \eta \frac{\partial \mathcal{E}(\boldsymbol{\theta}_t)}{\partial \boldsymbol{\theta}}$     (12.67)

8     $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \eta \frac{d\boldsymbol{\theta}}{dt}(\boldsymbol{\theta}_t, \mathbf{v}_t) = \boldsymbol{\theta}_t + \eta \frac{\partial \mathcal{K}(\mathbf{v}_t)}{\partial \mathbf{v}}$     (12.68)

10   If the kinetic energy has the form in Equation (12.64) then the second expression simplifies to  
11

12    $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \eta \boldsymbol{\Sigma}^{-1} \mathbf{v}_{t+1}$     (12.69)

14   This is known as **Euler's method**.

16   **12.5.2.2 Modified Euler's method**

18   The **Modified Euler's method** is slightly more accurate, and works as follows: first update the  
19   momentum, and then update the position using the new momentum:

21    $\mathbf{v}_{t+1} = \mathbf{v}_t + \eta \frac{d\mathbf{v}}{dt}(\boldsymbol{\theta}_t, \mathbf{v}_t) = \mathbf{v}_t - \eta \frac{\partial \mathcal{E}(\boldsymbol{\theta}_t)}{\partial \boldsymbol{\theta}}$     (12.70)

23    $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \eta \frac{d\boldsymbol{\theta}}{dt}(\boldsymbol{\theta}_t, \mathbf{v}_{t+1}) = \boldsymbol{\theta}_t + \eta \frac{\partial \mathcal{K}(\mathbf{v}_{t+1})}{\partial \mathbf{v}}$     (12.71)

25   Unfortunately, the asymmetry of this method can cause some theoretical problems (see e.g., [BZ20,  
26   p287]) which we resolve below.

27

28   **12.5.2.3 Leapfrog integrator**

30   In this section, we discuss the **leapfrog integrator** which is a symmetrized version of the modified  
31   Euler method. We first perform a “half” update of the momentum, then a full update of the position,  
32   and then finally another “half” update of the momentum:

34    $\mathbf{v}_{t+1/2} = \mathbf{v}_t - \frac{\eta}{2} \frac{\partial \mathcal{E}(\boldsymbol{\theta}_t)}{\partial \boldsymbol{\theta}}$     (12.72)

36    $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \eta \frac{\partial \mathcal{K}(\mathbf{v}_{t+1/2})}{\partial \mathbf{v}}$     (12.73)

38    $\mathbf{v}_{t+1} = \mathbf{v}_{t+1/2} - \frac{\eta}{2} \frac{\partial \mathcal{E}(\boldsymbol{\theta}_{t+1})}{\partial \boldsymbol{\theta}}$     (12.74)

40   If we perform multiple leapfrog steps, it is equivalent to performing a half step update of  $\mathbf{v}$  at the  
41   beginning and end of the trajectory, and alternating between full step updates of  $\boldsymbol{\theta}$  and  $\mathbf{v}$  in between.  
42

43   **12.5.2.4 Higher order integrators**

45   It is possible to define higher order integrators that are more accurate, but take more steps. For  
46   details, see [BRSS18].

47

---

1    **12.5.3 The HMC algorithm**

2    We now describe how to use Hamiltonian dynamics to define an MCMC sampler in the expanded  
3    state space  $(\boldsymbol{\theta}, \mathbf{v})$ . The target distribution has the form  
4

5

$$p(\boldsymbol{\theta}, \mathbf{v}) = \frac{1}{Z} \exp[-\mathcal{H}(\boldsymbol{\theta}, \mathbf{v})] = \frac{1}{Z} \exp \left[ -\mathcal{E}(\boldsymbol{\theta}) - \frac{1}{2} \mathbf{v}^\top \boldsymbol{\Sigma} \mathbf{v} \right] \quad (12.75)$$

6

7    The marginal distribution over the latent variables of interest has the form  
8

9

$$p(\boldsymbol{\theta}) = \int p(\boldsymbol{\theta}, \mathbf{v}) d\mathbf{v} = \frac{1}{Z_q} e^{-\mathcal{E}(\boldsymbol{\theta})} \int \frac{1}{Z_p} e^{-\frac{1}{2} \mathbf{v}^\top \boldsymbol{\Sigma} \mathbf{v}} d\mathbf{v} = \frac{1}{Z_q} e^{-\mathcal{E}(\boldsymbol{\theta})} \quad (12.76)$$

10

11   Suppose the previous state of the Markov chain is  $(\boldsymbol{\theta}_{t-1}, \mathbf{v}_{t-1})$ . To sample the next state, we  
12   proceed as follows. We set the initial position to  $\boldsymbol{\theta}'_0 = \boldsymbol{\theta}_{t-1}$ , and sample a new random momentum,  
13    $\mathbf{v}'_0 \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$ . We then initialize a random trajectory in the phase space, starting at  $(\boldsymbol{\theta}'_0, \mathbf{v}'_0)$ , and  
14   followed for  $L$  leapfrog steps, until we get to the final proposed state  $(\boldsymbol{\theta}^*, \mathbf{v}^*) = (\boldsymbol{\theta}'_L, \mathbf{v}'_L)$ . If we have  
15   simulated Hamiltonian mechanics correctly, the energy should be the same at the start and end  
16   of this process; if not, we say the HMC has **diverged**, and we reject the sample. If the energy is  
17   constant, we compute the MH acceptance probability  
18

19

$$\alpha = \min \left( 1, \frac{p(\boldsymbol{\theta}^*, \mathbf{v}^*)}{p(\boldsymbol{\theta}_{t-1}, \mathbf{v}_{t-1})} \right) = \min (1, \exp[-\mathcal{H}(\boldsymbol{\theta}^*, \mathbf{v}^*) + \mathcal{H}(\boldsymbol{\theta}_{t-1}, \mathbf{v}_{t-1})]) \quad (12.77)$$

20

21   (The transition probabilities cancel since the proposal is reversible.) Finally, we accept the proposal  
22   by setting  $(\boldsymbol{\theta}_t, \mathbf{v}_t) = (\boldsymbol{\theta}^*, \mathbf{v}^*)$  with probability  $\alpha$ , otherwise we set  $(\boldsymbol{\theta}_t, \mathbf{v}_t) = (\boldsymbol{\theta}_{t-1}, \mathbf{v}_{t-1})$ . (In practice  
23   we don't need to keep the momentum term, it is only used inside of the leapfrog algorithm.) See  
24   Algorithm 27 for the pseudocode.<sup>6</sup>

25

---

26   **Algorithm 27:** Hamiltonian Monte Carlo

---

27

28   1 **for**  $t = 1 : T$  **do**  
29     2 Generate random momentum  $\mathbf{v}_{t-1} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$   
30     3 Set  $(\boldsymbol{\theta}'_0, \mathbf{v}'_0) = (\boldsymbol{\theta}_{t-1}, \mathbf{v}_{t-1})$   
31     4 Half step for momentum:  $\mathbf{v}'_{\frac{1}{2}} = \mathbf{v}'_0 - \frac{\eta}{2} \nabla \mathcal{E}(\boldsymbol{\theta}'_0)$   
32     5 **for**  $l = 1 : L - 1$  **do**  
33       6      $\boldsymbol{\theta}'_l = \boldsymbol{\theta}'_{l-1} + \eta \boldsymbol{\Sigma}^{-1} \mathbf{v}'_{l-1/2}$   
34       7      $\mathbf{v}'_{l+1/2} = \mathbf{v}'_{l-1/2} - \eta \nabla \mathcal{E}(\boldsymbol{\theta}'_l)$   
35     8 Full step for location:  $\boldsymbol{\theta}'_L = \boldsymbol{\theta}'_{L-1} + \eta \boldsymbol{\Sigma}^{-1} \mathbf{v}'_{L-1/2}$   
36     9 Half step for momentum:  $\mathbf{v}'_L = \mathbf{v}'_{L-1/2} - \frac{\eta}{2} \nabla \mathcal{E}(\boldsymbol{\theta}'_L)$   
37     10 Compute proposal  $(\boldsymbol{\theta}^*, \mathbf{v}^*) = (\boldsymbol{\theta}'_L, \mathbf{v}'_L)$   
38     11 Compute  $\alpha = \min (1, \exp[-\mathcal{H}(\boldsymbol{\theta}^*, \mathbf{v}^*) + \mathcal{H}(\boldsymbol{\theta}_{t-1}, \mathbf{v}_{t-1})])$   
39     12 Set  $\boldsymbol{\theta}_t = \boldsymbol{\theta}^*$  with probability  $\alpha$ , otherwise  $\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1}$ .

---

40

41   6. There are many high-quality implementations of HMC. For example, [BlackJAX](#) in JAX.

42

We need to sample a new momentum at each iteration to satisfy ergodicity. To see why, recall that  $\mathcal{H}(\boldsymbol{\theta}, \mathbf{v})$  stays approximately constant as we move through phase space. If  $\mathcal{H}(\boldsymbol{\theta}, \mathbf{v}) = \mathcal{E}(\boldsymbol{\theta}) + \frac{1}{2}\mathbf{v}^\top \Sigma \mathbf{v}$ , then clearly  $\mathcal{E}(\boldsymbol{\theta}) \leq \mathcal{H}(\boldsymbol{\theta}, \mathbf{v}) = h$  for all locations  $\boldsymbol{\theta}$  along the trajectory. Thus the sampler cannot reach states where  $\mathcal{E}(\boldsymbol{\theta}) > h$ . To ensure the sampler explores the full space, we must pick a random momentum at the start of each iteration.

#### 12.5.4 Tuning HMC

We need to specify three hyperparameters for HMC: the number of leapfrog steps  $L$ , the step size  $\eta$ , and the covariance  $\Sigma$ .

We want to choose the number of leapfrog steps  $L$  to be large enough that the algorithm explores the level set of constant energy, but without doubling back on itself, which would waste computation, due to correlated samples. Fortunately, there is an algorithm, known as the **No-U-Turn Sampler** or **NUTS** algorithm [HG14], which can adaptively choose  $L$  for us.

##### 12.5.4.2 Choosing the step size

When  $\Sigma = \mathbf{I}$ , the ideal step size  $\eta$  should be roughly equal to the width of  $\mathcal{E}(\boldsymbol{\theta})$  in the most constrained direction of the local energy landscape. For a locally quadratic potential, this corresponds to the square root of the smallest marginal standard deviation of the local covariance matrix. (If we think of the energy surface as a valley, this corresponds to the direction with the steepest sides.) A step size much larger than this will cause moves that are likely to be rejected because they move to places which increase the potential energy too much. On the other hand, if the step size is too low, the proposal distribution will not move much from the starting position, and the algorithm will be very slow.

In [BZ20, Sec 9.5.4] they recommend the following heuristic for picking  $\eta$ : set  $\Sigma = \mathbf{I}$  and  $L = 1$ , and then vary  $\eta$  until the acceptance rates are in the range of 40%–80%. Of course, different step sizes might be needed in different parts of the state space. In this case, we can use learning rate schedules from the optimization literature, such as cyclical schedules [Zha+20d].

##### 12.5.4.3 Choosing the covariance (inverse mass) matrix

To allow for larger step sizes, we can use a smarter choice for  $\Sigma$ , also called the **inverse mass matrix**. One way to estimate a fixed  $\Sigma$  is to run HMC with  $\Sigma = \mathbf{I}$  for a **warmup** period, until the chain is “burned in” (see Section 12.6); then we run for a few more steps, so we can compute the empirical covariance matrix using  $\Sigma = \mathbb{E}[(\boldsymbol{\theta} - \bar{\boldsymbol{\theta}})(\boldsymbol{\theta} - \bar{\boldsymbol{\theta}})^\top]$ . In [Hof+19] they propose a method called **NeuTra HMC** algorithm which “neutralizes” bad geometry by learning an inverse autoregressive flow model (Section 23.2.4.3) in order to map the warped distribution to an isotropic Gaussian. This is often an order of magnitude faster than vanilla HMC.

---

### 12.5.5 Riemann Manifold HMC

If we let the covariance matrix change as we move position, so  $\Sigma$  is a function of  $\theta$ , the method is known as **Riemann Manifold HMC** or **RM-HMC** [GC11; Bet13], since the moves follow a curved manifold, rather than the flat manifold induced by a constant  $\Sigma$ .

A natural choice for the covariance matrix is to use the Hessian at the current location, to capture the local geometry:

$$\Sigma(\theta) = \nabla^2 \mathcal{E}(\theta) \quad (12.78)$$

Since this is not always positive definite, an alternative, that can be used for some problems, is to use the Fisher information matrix (Section 2.4), given by

$$\Sigma(\theta) = -\mathbb{E}_{p(x|\theta)} [\nabla^2 \log p(x|\theta)] \quad (12.79)$$

Once we have computed  $\Sigma(\theta)$ , we can compute the kinetic energy as follows:

$$\mathcal{K}(\theta, v) = \frac{1}{2} \log((2\pi)^D |\Sigma(\theta)|) + \frac{1}{2} v^\top \Sigma(\theta) v \quad (12.80)$$

Unfortunately the Hamiltonian updates of  $\theta$  and  $v$  are no longer separable, which makes the RM-HMC algorithm more complex to implement, so it is not widely used.

### 12.5.6 Langevin Monte Carlo (MALA)

A special case of HMC occurs when we take  $L = 1$  leapfrog steps. This is known as **Langevin Monte Carlo (LMC)**, or **Metropolis Adjusted Langevin Algorithm (MALA)** [RT96]. This gives rise to the simplified algorithm shown in Algorithm 28.

---

#### Algorithm 28: Langevin Monte Carlo

---

```

1 for t = 1 : T do
2   Generate random momentum  $v_{t-1} \sim \mathcal{N}(\mathbf{0}, \Sigma)$ 
3    $\theta^* = \theta_{t-1} - \frac{\eta^2}{2} \Sigma^{-1} \nabla \mathcal{E}(\theta_{t-1}) + \eta \Sigma^{-1} v_{t-1}$ 
4    $v^* = v_{t-1} - \frac{\eta}{2} \nabla \mathcal{E}(\theta_{t-1}) - \frac{\eta}{2} \nabla \mathcal{E}(\theta^*)$ 
5   Compute  $\alpha = \min(1, \exp[-\mathcal{H}(\theta^*, v^*)] / \exp[-\mathcal{H}(\theta_{t-1}, v_{t-1})])$ 
6   Set  $\theta_t = \theta^*$  with probability  $\alpha$ , otherwise  $\theta_t = \theta_{t-1}$ .

```

---

A further simplification is to eliminate the MH acceptance step. In this case, the update becomes

$$\theta_t = \theta_{t-1} - \frac{\eta^2}{2} \Sigma^{-1} \nabla \mathcal{E}(\theta_{t-1}) + \eta \Sigma^{-1} v_{t-1} \quad (12.81)$$

$$= \theta_{t-1} - \frac{\eta^2}{2} \Sigma^{-1} \nabla \mathcal{E}(\theta_{t-1}) + \eta \sqrt{\Sigma^{-1}} \epsilon_{t-1} \quad (12.82)$$

where  $v_{t-1} \sim \mathcal{N}(\mathbf{0}, \Sigma)$  and  $\epsilon_{t-1} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . This is just like gradient descent with added noise. If we set  $\Sigma$  to be the Fisher information matrix, this becomes natural gradient descent (Section 6.4) with

1 added noise. If we approximate the gradient with a stochastic gradient, we get a method known as  
2 **SGLD**, or **Stochastic Gradient Langevin Descent** (see Section 12.7.1 for details).  
3

4 Now suppose  $\Sigma = \mathbf{I}$ , and we set  $\eta = \sqrt{2}$ . In continuous time, we get the following stochastic  
5 differential equation (SDE), known as **Langevin diffusion**:

$$\underline{6} \quad d\boldsymbol{\theta}_t = -\nabla \mathcal{E}(\boldsymbol{\theta}_t) dt + \sqrt{2} dB_t \quad (12.83)$$

7 where  $B_t$  represents  $D$ -dimensional **Brownian motion**. If we use this to generate the samples, the  
8 method is known as the **Unadjusted Langevin Algorithm** or **ULA** [Par81; RT96].  
9

10

### 11 12.5.7 Connection between SGD and Langevin sampling

12 In this section, we discuss a deep connection between stochastic gradient descent (SGD) and Langevin  
13 sampling, following the presentation of [BZ20, Sec 10.2.3].  
14

15 Consider the minimization of the additive loss

$$\underline{16} \quad \mathcal{L}(\boldsymbol{\theta}) = \sum_{n=1}^N \mathcal{L}_n(\boldsymbol{\theta}) \quad (12.84)$$

17 For example, we may define  $\mathcal{L}_n(\boldsymbol{\theta}) = -\log p(y_n | \mathbf{x}_n, \boldsymbol{\theta})$ . We will use a minibatch approximation to  
18 the gradients:  
19

$$\underline{20} \quad \nabla_B \mathcal{L}(\boldsymbol{\theta}) = \frac{1}{B} \sum_{n \in \mathcal{S}} \nabla \mathcal{L}_n(\boldsymbol{\theta}) \quad (12.85)$$

21 where  $\mathcal{S} = \{i_1, \dots, i_B\}$  is a randomly chosen set of indices of size  $B$ . For simplicity of analysis, we  
22 assume the indices are chosen with replacement from  $\{1, \dots, N\}$ .  
23

24 Let us define the (scaled) error (due to minibatching) in the estimated gradient by  
25

$$\underline{26} \quad \mathbf{v}_t \triangleq \sqrt{\eta} (\nabla \mathcal{L}(\boldsymbol{\theta}_t) - \nabla_B \mathcal{L}(\boldsymbol{\theta}_t)) \quad (12.86)$$

27 This is called the **diffusion term**. Then we can rewrite the SGD update as  
28

$$\underline{29} \quad \boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \nabla_B \mathcal{L}(\boldsymbol{\theta}_t) = \boldsymbol{\theta}_t - \eta \nabla \mathcal{L}(\boldsymbol{\theta}_t) + \sqrt{\eta} \mathbf{v}_t \quad (12.87)$$

30 The diffusion term  $\mathbf{v}_t$  has mean 0, since  
31

$$\underline{32} \quad \mathbb{E} [\nabla_B \mathcal{L}(\boldsymbol{\theta})] = \frac{1}{B} \sum_{j=1}^B \mathbb{E} [\nabla \mathcal{L}_{i_j}(\boldsymbol{\theta})] = \frac{1}{B} \sum_{j=1}^B \nabla \mathcal{L}(\boldsymbol{\theta}) = \nabla \mathcal{L}(\boldsymbol{\theta}) \quad (12.88)$$

33 To compute the variance of the diffusion term, note that  
34

$$\underline{35} \quad \mathbb{V} [\nabla_B \mathcal{L}(\boldsymbol{\theta})] = \frac{1}{B^2} \sum_{j=1}^B \mathbb{V} [\nabla \mathcal{L}_{i_j}(\boldsymbol{\theta})] \quad (12.89)$$

36

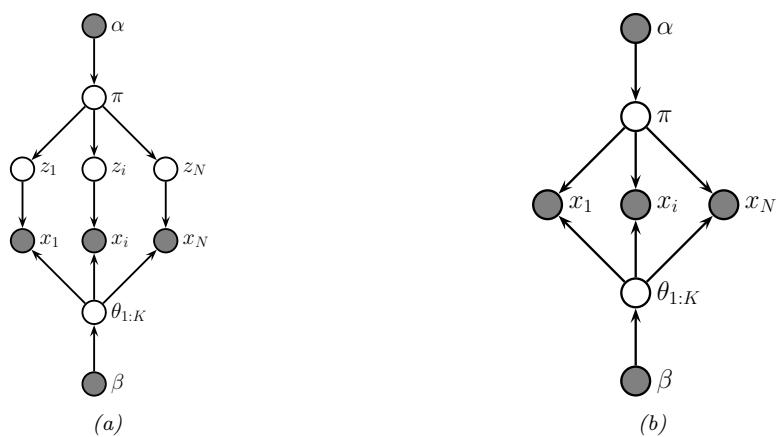


Figure 12.10: (a) A mixture model. (b) After integrating out the discrete latent variables.

where

$$\mathbb{V} [\nabla \mathcal{L}_{i_j}(\boldsymbol{\theta})] = \mathbb{E} [\nabla \mathcal{L}_{i_j}(\boldsymbol{\theta}) \nabla \mathcal{L}_{i_j}(\boldsymbol{\theta})^\top] - \mathbb{E} [\nabla \mathcal{L}_{i_j}(\boldsymbol{\theta})] \mathbb{E} [\nabla \mathcal{L}_{i_j}(\boldsymbol{\theta})^\top] \quad (12.90)$$

$$= \left( \frac{1}{N} \sum_{n=1}^N \nabla \mathcal{L}_n(\boldsymbol{\theta}) \nabla \mathcal{L}_n(\boldsymbol{\theta})^\top \right) - \nabla \mathcal{L}(\boldsymbol{\theta}) \nabla \mathcal{L}(\boldsymbol{\theta})^\top \triangleq \mathbf{D}(\boldsymbol{\theta}) \quad (12.91)$$

where  $\mathbf{D}(\boldsymbol{\theta})$  is called the **diffusion matrix**. Hence  $\mathbb{V}[\mathbf{v}_t] = \frac{\eta}{B} \mathbf{D}(\boldsymbol{\theta}_t)$ .

[LTW15] prove that the following continuous time stochastic differential equation is a first-order approximation of minibatch SGD (assuming the loss function is Lipschitz continuous):

$$d\boldsymbol{\theta}(t) = -\nabla \mathcal{L}(\boldsymbol{\theta}(t)) dt + \sqrt{\frac{\eta}{B} \mathbf{D}(\boldsymbol{\theta}_t)} d\mathbf{B}(t) \quad (12.92)$$

where  $\mathbf{B}(t)$  is Brownian motion. Thus the noise from minibatching causes SGD to act like a Langevin sampler. (See [Hu+17] for more information.)

The scale factor for the noise,  $\tau = \frac{\eta}{B}$ , plays the role of **temperature**. Thus we see that using a smaller batch size is like using a larger temperature; the added noise ensures that SGD avoids going into narrow ravines, and instead spends most of its time in flat minima which have better generalization performance [Kes+17]. See Section 17.4.1 for more discussion of this point.

#### 12.5.8 Applying HMC to constrained parameters

To apply HMC, we require that all the latent quantities be continuous (real-valued) and have unconstrained support, i.e.,  $\boldsymbol{\theta} \in \mathbb{R}^D$ , so discrete latent variables need to be marginalized out (although some recent work, such as [NDL20; Zho20], relaxes this requirement).

As an example of how this can be done, consider a GMM. We can easily write the likelihood

1 without discrete latents as follows:

2

$$\underline{4} \quad p(\mathbf{x}_n | \boldsymbol{\theta}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (12.93)$$

5

6 The corresponding “collapsed” model is shown in Figure 12.10(b). (Note that this is the opposite of  
7 Section 12.3.8, where we integrated out the continuous parameters in order to apply Gibbs sampling  
8 to the discrete latents.) We can apply similar techniques to other discrete latent variable models. For  
9 example, to apply HMC to HMMs, we can use the forwards algorithm (Section 8.2.2) to efficiently  
10 compute  $p(\mathbf{x}_n | \boldsymbol{\theta}) = \sum_{\mathbf{z}_{1:T}} p(\mathbf{x}_n, \mathbf{z}_{n,1:T} | \boldsymbol{\theta})$ .

11 In addition to marginalizing out any discrete latent variables, we need to ensure the remaining  
12 continuous latent variables are unconstrained. This often requires performing a change of variables  
13 using a bijector. For example, instead of sampling the discrete probability vector from the probability  
14 simplex  $\boldsymbol{\pi} \in \mathbb{S}^K$ , we should sample the logits  $\boldsymbol{\eta} \in \mathbb{R}^K$ . After sampling, we can transform back, since  
15 bijectors are invertible. (For a practical example, see [change\\_of\\_variable\\_hmc.ipynb](#).)

16

17

### 18 12.5.9 Speeding up HMC

19 Although HMC uses gradient information to explore the typical set, sometimes the geometry of the  
20 typical set can be difficult to sample from. See Section 12.5.4.3 for ways to estimate the mass matrix,  
21 which can help with such difficult cases.

22 Another issue is the cost of evaluating the target distribution,  $\mathcal{E}(\boldsymbol{\theta}) = -\log \tilde{p}(\boldsymbol{\theta})$ . For many ML  
23 applications, this has the form  $\log \tilde{p}(\boldsymbol{\theta}) = \log p_0(\boldsymbol{\theta}) + \sum_{n=1}^N \log p(\boldsymbol{\theta}_n | \boldsymbol{\theta})$ . This takes  $O(N)$  time to  
24 compute. We can speed this up by using stochastic gradient methods; see Section 12.7 for details.

25

26

## 27 12.6 MCMC convergence

28

29 We start MCMC from an arbitrary initial state. As we explained in Section 2.6.4, the samples will be  
30 coming from the chain’s stationary distribution only when the chain has “forgotten” where it started  
31 from. The amount of time it takes to enter the stationary distribution is called the mixing time (see  
32 Section 12.6.1 for details). Samples collected before the chain has reached its stationary distribution  
33 do not come from  $p^*$ , and are usually thrown away. The initial period, whose samples will be ignored,  
34 is called the **burn-in phase**.

35 For example, consider a uniform distribution on the integers  $\{0, 1, \dots, 20\}$ . Suppose we sample  
36 from this using a symmetric random walk. In Figure 12.11, we show two runs of the algorithm. On  
37 the left, we start in state 10; on the right, we start in state 17. Even in this small problem it takes  
38 over 200 steps until the chain has “forgotten” where it started from. Proposal distributions that  
39 make larger changes can converge faster. For example, [BD92; Man] prove that it takes about 7 riffle  
40 shuffles to properly mix a deck of 52 cards (i.e., to ensure the distribution is uniform).

41 In Section 12.6.1 we discuss how to compute the mixing time theoretically. In practice, this can be  
42 very hard [BBM10] (this is one of the fundamental weaknesses of MCMC), so in Section 12.6.2, we  
43 discuss practical heuristics.

44

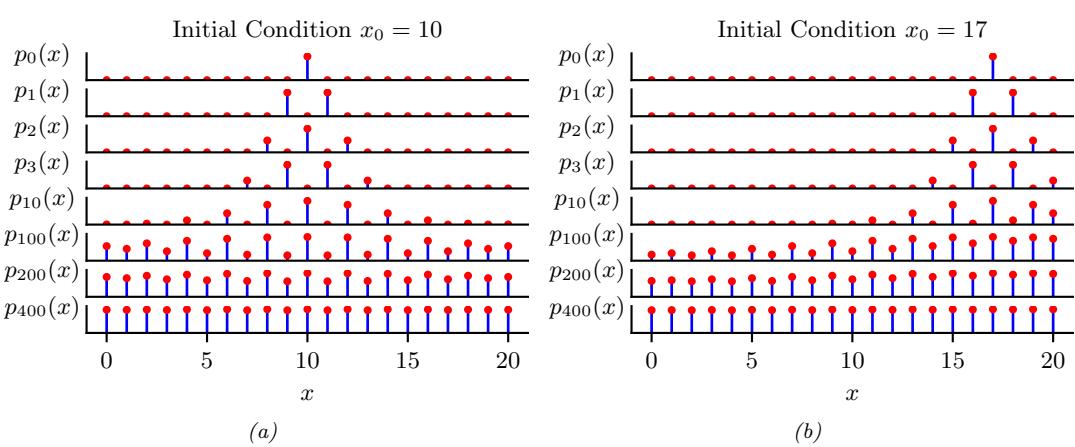


Figure 12.11: Illustration of convergence to the uniform distribution over  $\{0, 1, \dots, 20\}$  using a symmetric random walk starting from (left) state 10, and (right) state 17. Adapted from Figures 29.14 and 29.15 of [Mac03]. Generated by `random_walk_integers.ipynb`.

### 12.6.1 Mixing rates of Markov chains

The amount of time it takes for a Markov chain to converge to the stationary distribution, and forget its initial state, is called the **mixing time**. More formally, we say that the mixing time from state  $x_0$  is the minimal time such that, for any constant  $\epsilon > 0$ , we have that

$$\tau_\epsilon(x_0) \triangleq \min\{t : \|\delta_{x_0}(x)T^t - p^*\|_1 \leq \epsilon\} \quad (12.94)$$

where  $\delta_{x_0}(x)$  is a distribution with all its mass in state  $x_0$ ,  $T$  is the transition matrix of the chain (which depends on the target  $p^*$  and the proposal  $q$ ), and  $\delta_{x_0}(x)T^t$  is the distribution after  $t$  steps. The mixing time of the chain is defined as

$$\tau_\epsilon \triangleq \max_{x_0} \tau_\epsilon(x_0) \quad (12.95)$$

This is the maximum amount of time it takes for the chain's distribution to get  $\epsilon$  close to  $p^*$  from any starting state.

The mixing time is determined by the eigengap  $\gamma = \lambda_1 - \lambda_2$ , which is the difference between the first and second eigenvalues of the transition matrix. For a finite state chain, one can show  $\tau_\epsilon = O(\frac{1}{\gamma} \log \frac{n}{\epsilon})$ , where  $n$  is the number of states.

We can also study the problem by examining the geometry of the state space. For example, consider the chain in Figure 12.12. We see that the state space consists of two “islands”, each of which is connected via a narrow “bottleneck”. (If they were completely disconnected, the chain would not be ergodic, and there would no longer be a unique stationary distribution, as discussed in Section 2.6.4.3.) We define the **conductance**  $\phi$  of a chain as the minimum probability, over all subsets  $S$  of states, of transitioning from that set to its complement:

$$\phi \triangleq \min_{S: 0 \leq p^*(S) \leq 0.5} \frac{\sum_{x \in S, x' \in S^c} T(x \rightarrow x')}{p^*(S)}, \quad (12.96)$$

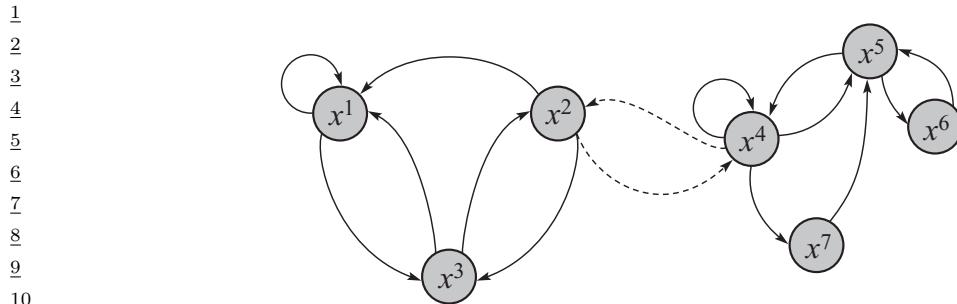


Figure 12.12: A Markov chain with low conductance. The dotted arcs represent transitions with very low probability. From Figure 12.6 of [KF09a]. Used with kind permission of Daphne Koller.

One can show that  $\tau_\epsilon \leq O\left(\frac{1}{\phi^2} \log \frac{n}{\epsilon}\right)$ . Hence chains with low conductance have high mixing time. For example, distributions with well-separated modes usually have high mixing time. Simple MCMC methods, such as MH and Gibbs, often do not work well in such cases, and more advanced algorithms, such as parallel tempering, are necessary (see e.g., [ED05; Kat+06; BZ20]).

### 12.6.2 Practical convergence diagnostics

Computing the mixing time of a chain is in general quite difficult, since the transition matrix is usually very hard to compute. Furthermore, diagnosing convergence is computationally intractable in general [BBM10]. Nevertheless, various heuristics have been proposed — see e.g., [Gey92; CC96; BR98; Veh+19]. We discuss some of the current recommended approaches below, following [Veh+19].

#### 12.6.2.1 Trace plots

One of the simplest approaches to assessing if the method has converged is to run multiple chains from very different **overdispersed** starting points, and to plot the samples of some quantity of interest, such as the value of a certain component of the state vector, or some event such as the value taking on an extreme value. This is called a **trace plot**. If the chain has mixed, it should have “forgotten” where it started from, so the trace plots should converge to the same distribution, and thus overlap with each other.

To illustrate this, we will consider a very simple, but enlightening, example from [McE20, Sec 9.5]. The model is a univariate Gaussian,  $y_i \sim \mathcal{N}(\alpha, \sigma)$ , with just 2 observations,  $y_1 = -1$  and  $y_2 = +1$ . We first consider a very diffuse prior,  $\alpha \sim \mathcal{N}(0, 1000)$  and  $\sigma \sim \text{Expon}(0.0001)$ , both of which allow for very large values of  $\alpha$  and  $\sigma$ . We fit the model using HMC using 3 chains and 500 samples. The result is shown in Figure 12.13. On the right, we show the trace plots for  $\alpha$  and  $\sigma$  for 3 different chains. We see that they do not overlap much with each other. In addition, the numerous black vertical lines at the bottom of the plot indicate that HMC had many divergences.

The problem is caused by the overly diffuse priors, which do not get overwhelmed by the likelihood because we only have 2 data points. Thus the posterior is also diffuse. We can fix this by using slightly stronger priors, that keep the parameters close to more sensible values. For example, suppose we use  $\alpha \sim \mathcal{N}(1, 10)$  and  $\sigma \sim \text{Expon}(1)$ . Now we get the results in Figure 12.14. On the right we see

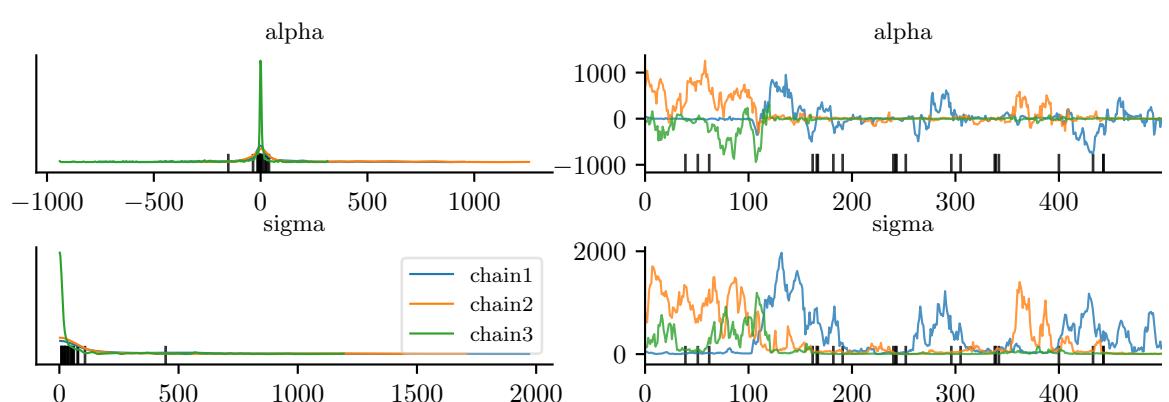


Figure 12.13: Marginals (left) and trace plot (right) for the univariate Gaussian using the diffuse prior. Black vertical lines indicate HMC divergences. Adapted from Figures 9.9–9.10 of [McE20]. Generated by `mcmc_traceplots_unigauss.ipynb`.

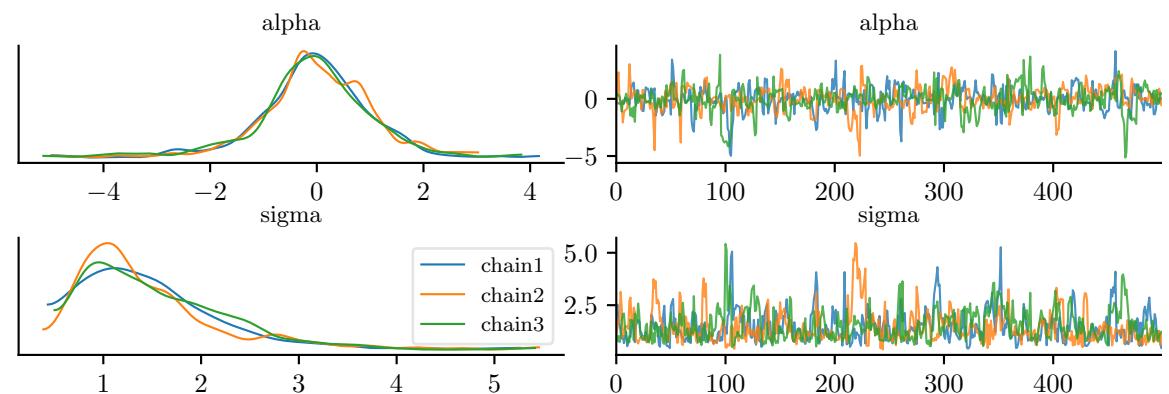
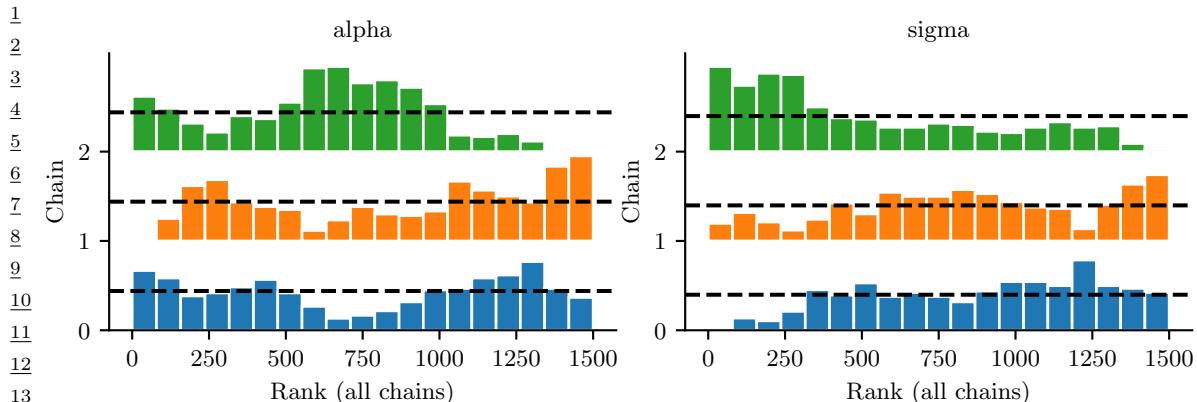


Figure 12.14: Marginals (left) and trace plot (right) for the univariate Gaussian using the sensible prior. Adapted from Figures 9.9–9.10 of [McE20]. Generated by `mcmc_traceplots_unigauss.ipynb`.

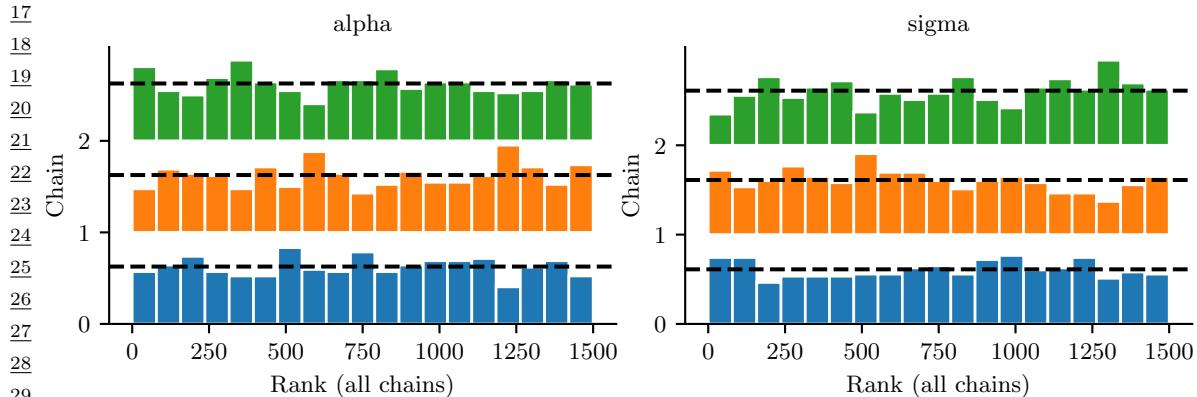
that the traceplots overlap. On the left, we see that the marginal distributions from each chain have support over a reasonable interval, and have a peak at the “right” place (the MLE for  $\alpha$  is 0, and for  $\sigma$  is 1). And we don’t see any divergence warnings (vertical black markers in the plot).

Since trace plots of converging chains correspond to overlapping lines, it can be hard to distinguish success from failure. An alternative plot, known as a **trace rank plot**, was recently proposed in [Veh+19]. (In [McE20], this is called a **trankplot**, a term we borrow.) The idea is to compute the rank of each sample based on all the samples from all the chains, after burnin. We then plot a histogram of the ranks for each chain separately. If the chains have converged, the distribution over ranks should be uniform, since there should be no preference for high or low scoring samples amongst the chains.

The trankplot for the model with the diffuse prior is shown in Figure 12.15. (The x-axis is from 1



*Figure 12.15: Trace rank plot for the univariate Gaussian using the diffuse prior. Adapted from Figures 9.9–9.10 of [McE20]. Generated by `mcmc_traceplots_unigauss.ipynb`.*



*Figure 12.16: Trace rank plot for the univariate Gaussian using the sensible prior. Adapted from Figures 9.9–9.10 of [McE20]. Generated by `mcmc_traceplots_unigauss.ipynb`.*

to the total number of samples, which in this example is 1500, since we use 3 chains and draw 500 samples from each.) We can see that the different chains are clearly not mixing. The traceplot for the model with the sensible prior is shown in Figure 12.16; this looks much better.

### 12.6.2.2 Estimated potential scale reduction (EPSR)

In this section, we discuss a way to assess convergence more quantitatively. The basic idea is this: if one or more chains has not mixed well, then the variance of all the chains combined together will be higher than the variance of the individual chains. So we will compare the variance of the quantity of interest computed between and within chains.

More precisely, suppose we have  $M$  chains, and we draw  $N$  samples from each. Let  $x_{nm}$  denote the quantity of interest derived from the  $n$ 'th sample from the  $m$ 'th chain. We compute the between

1 and within-sequence variances as follows:  
 2

$$3 \quad B = \frac{N}{M-1} \sum_{m=1}^M (\bar{x}_{\cdot m} - \bar{x}_{\cdot \cdot})^2, \text{ where } \bar{x}_{\cdot m} = \frac{1}{N} \sum_{n=1}^N x_{nm}, \bar{x}_{\cdot \cdot} = \frac{1}{M} \sum_{m=1}^M \bar{x}_{\cdot m} \quad (12.97)$$

$$4 \quad W = \frac{1}{M} \sum_{m=1}^M s_m^2, \text{ where } s_m^2 = \frac{1}{N-1} \sum_{n=1}^N (x_{nm} - \bar{x}_{\cdot m})^2 \quad (12.98)$$

5 The formula for  $s_m^2$  is the usual unbiased estimate for the variance from a set of  $N$  samples;  $W$  is  
 6 just the average of this. The formula for  $B$  is similar, but scaled up by  $N$  since it is based on the  
 7 variance of  $\bar{x}_{\cdot m}$ , which are averaged over  $N$  values.  
 8

9 Next we compute the following average variance:  
 10

$$11 \quad \hat{V}^+ \triangleq \frac{N-1}{N} W + \frac{1}{N} B \quad (12.99)$$

12 Finally, we compute the following quantity, known as the **estimated potential scale reduction**  
 13 or **R-hat**:  
 14

$$15 \quad \hat{R} \triangleq \sqrt{\frac{\hat{V}^+}{W}} \quad (12.100)$$

16 In [Veh+19], they recommend checking if  $\hat{R} < 1.01$  before declaring convergence.  
 17

18 For example, consider the  $\hat{R}$  values for various samplers for our univariate GMM example. In  
 19 particular, consider the 3 MH samplers in Figure 12.1, and the Gibbs sampler in Figure 12.4. The  $\hat{R}$   
 20 values are 1.493, 1.039, 1.005 and 1.007. So this diagnostic has correctly identified that the first two  
 21 samplers are unreliable, which evident from the figure.  
 22

23 In practice, it is recommended to use a slightly different quantity, known as **split- $\hat{R}$** . This can  
 24 be computed by splitting each chain into the first and second halves, thus doubling the number of  
 25 chains  $M$  (but halving the number of samples  $N$  from each), before computing  $\hat{R}$ . This can detect  
 26 non-stationarity within a single chain.  
 27

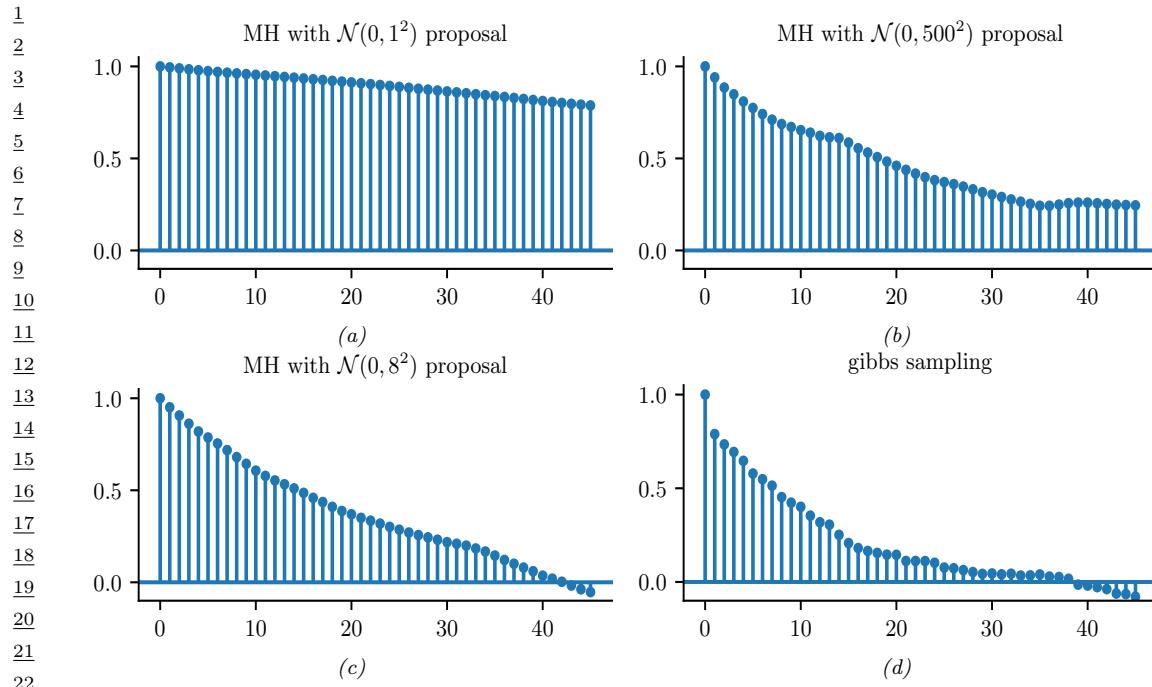
### 32 12.6.3 Effective sample size

33 Although MCMC lets us draw samples from a target distribution (assuming it has converged), the  
 34 samples are not independent, so we may need to draw a lot of them to get a reliable estimate. In  
 35 this section, we discuss how to compute the **effective sample size** or **ESS** from a set of (possibly  
 36 correlated) samples.  
 37

38 To start, suppose we draw  $N$  *independent* samples from the target distribution, and let  $\hat{x} =$   
 39  $\frac{1}{N} \sum_{n=1}^N x_n$  be our empirical estimate of the mean of the quantity of interest. The variance of this  
 40 estimate is given by  
 41

$$42 \quad \mathbb{V}[\hat{x}] = \frac{1}{N^2} \mathbb{V} \left[ \sum_{n=1}^N x_n \right] = \frac{1}{N^2} \sum_{n=1}^N \mathbb{V}[x_n] = \frac{1}{N} \sigma^2 \quad (12.101)$$

43 where  $\sigma^2 = \mathbb{V}[X]$ . If the samples are correlated, the variance of the estimate will be higher, as we  
 44 show below.  
 45



23 *Figure 12.17: Autocorrelation functions for various MCMC samplers for the mixture of two 1D Gaussians.*  
 24 (a-c) These are the MH samplers in Figure 12.1. (d) This is the Gibbs sampler in Figure 12.4. Generated by  
 25 `mcmc_gmm_demo.ipynb`.

26  
 27  
 28 Recall that for  $N$  (not necessarily independent) random variables we have  
 29

$$30 \quad \mathbb{V} \left[ \sum_{n=1}^N x_n \right] = \sum_{i=1}^N \sum_{j=1}^N \text{Cov}[x_i, x_j] = \sum_{i=1}^N \mathbb{V}[x_i] + 2 \sum_{1 \leq i < j \leq N} \text{Cov}[x_i, x_j] \quad (12.102)$$

33  
 34 Let  $\bar{x} = \frac{1}{N} \sum_{n=1}^N x_n$  be our estimate based on these correlated samples. The variance of this estimate  
 35 is given by

$$36 \quad \mathbb{V}[\bar{x}] = \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N \text{Cov}[x_i, x_j] \quad (12.103)$$

40 We now rewrite this in a more convenient form. First recall that the correlation of  $x_i$  and  $x_j$  is  
 41 given by

$$43 \quad \text{corr}[x_i, x_j] = \frac{\text{Cov}[x_i, x_j]}{\sqrt{\mathbb{V}[x_i] \mathbb{V}[x_j]}} \quad (12.104)$$

46 Since we assume we are drawing samples from the target distribution, we have  $\mathbb{V}[x_i] = \sigma^2$ , where  $\sigma^2$   
 47

1  
2 is the true variance. Hence

3  
4  $\mathbb{V}[\bar{x}] = \frac{\sigma^2}{N^2} \sum_{i=1}^N \sum_{j=1}^N \text{corr}[x_i, x_j]$  (12.105)  
5

6  
7 For a fixed  $i$ , we can think of  $\text{corr}[x_i, x_j]$  as a function of  $j$ . This will usually decay as  $j$  gets further  
8 from  $i$ . As  $N \rightarrow \infty$  we can approximate the sum of correlations by

9  
10  $\sum_{j=1}^N \text{corr}[x_i, x_j] \rightarrow \sum_{\ell=-\infty}^{\infty} \text{corr}[x_i, x_{i+\ell}] = 1 + 2 \sum_{\ell=1}^{\infty} \text{corr}[x_i, x_{i+\ell}]$  (12.106)  
11

12  
13 since  $\text{corr}[x_i, x_i] = 1$  and  $\text{corr}[x_i, x_{i-\ell}] = \text{corr}[x_i, x_{i+\ell}]$  for lag  $\ell > 0$ . Since we assume the samples  
14 are coming from a stationary distribution, the index  $i$  does not matter. This we can define the  
15 **autocorrelation time** as

16  
17  $\rho = 1 + 2 \sum_{\ell=1}^{\infty} \rho(\ell)$  (12.107)  
18

19  
20 where  $\rho(\ell)$  is the **autocorrelation function** (ACF), defined as

21  
22  $\rho(\ell) \triangleq \text{corr}[x_0, x_\ell]$  (12.108)

23 The ACF can be approximated efficiently by convolving the signal  $\mathbf{x}$  with itself. In Figure 12.17,  
24 we plot the ACF for our four samplers for the GMM. We see that the ACF of the Gibbs sampler  
25 (bottom right) dies off to 0 much more rapidly than the MH samplers. Intuitively this indicates that  
26 each Gibbs sample is “worth” more than each MH sample. We quantify this below.

27 From Equation (12.105), we can compute the variance of our estimate in terms of the ACF as  
28 follows:  $\mathbb{V}[\bar{x}] = \frac{\sigma^2}{N^2} \sum_{i=1}^N \rho = \frac{\sigma^2}{N} \rho$ . By contrast, the variance of the estimate from independent  
29 samples is  $\mathbb{V}[\hat{x}] = \frac{\sigma^2}{N}$ . So we see that the variance is a factor  $\rho$  larger when there is correlation. We  
30 therefore define the **effective sample size** of our set of samples to be

31  
32  $N_{\text{eff}} \triangleq \frac{N}{\rho} = \frac{N}{1 + 2 \sum_{\ell=1}^{\infty} \rho(\ell)}$  (12.109)  
33

34  
35 In practice, we truncate the sum at lag  $L$ , which is the last integer at which  $\rho(L)$  is positive. Also, if  
36 we run  $M$  chains, the numerator should be  $NM$ , so we get the following estimate:

37  
38  $\hat{N}_{\text{eff}} = \frac{NM}{1 + 2 \sum_{\ell=1}^L \hat{\rho}(\ell)}$  (12.110)  
39

40  
41 In [Veh+19], they propose various extensions of the above estimator, such as using rank statistics,  
42 to make the estimate more robust.

#### 43 12.6.4 Improving speed of convergence 44

45 There are many possible things you could try if the  $\hat{R}$  value is too large, and/or the effective sample  
46 size is too low. Here is a brief list:

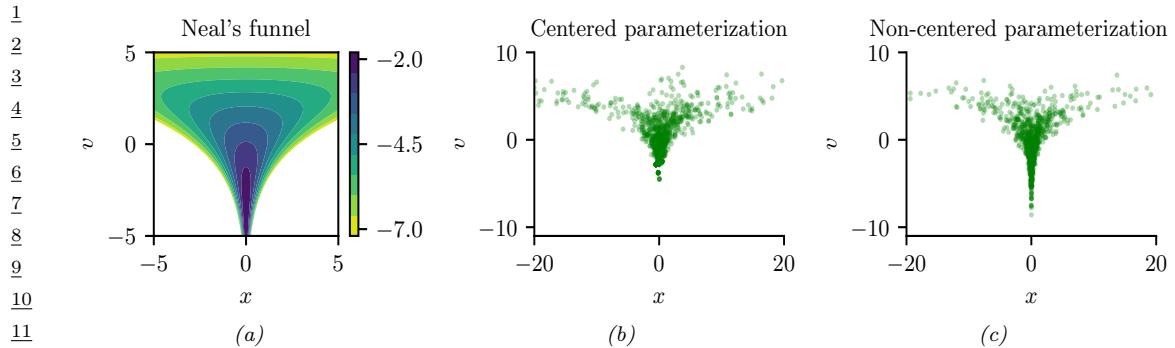


Figure 12.18: Neal’s funnel. (a) Joint density. (b) HMC samples from centered representation. (c) HMC samples from non-centered representation. Generated by [neals\\_funnel.ipynb](#).

- Try using a non-centered parameterization (see Section 12.6.5).
- Try sampling variables in groups or blocks (see Section 12.3.7).
- Try using Rao-Blackwellisation, i.e., analytically integrating out some of the variables (see Section 12.3.8).
- Try adding auxiliary variables (see Section 12.4).
- Try using adaptive proposal distributions (see Section 12.2.3.5).

More details can be found in [Rob+18].

## 12.6.5 Non-centered parameterizations and Neal’s funnel

A common problem that arises when applying sampling to hierarchical Bayesian models is when a set of parameters at one level of the model have a tight dependence on parameters at the level above. We show some practical examples of this in the hierarchical Gaussian 8-schools example in Section 3.7.2.2 and the hierarchical radon regression example in Section 15.5.3.2. Here, we focus on the following simple toy model that captures the essence of the problem:

$$\nu \sim \mathcal{N}(0, 3) \tag{12.111}$$

$$x \sim \mathcal{N}(0, \exp(\nu)) \tag{12.112}$$

The corresponding joint density  $p(x, \nu)$  is shown in Figure 12.18a. This is known **Neal’s funnel**, named after [Nea03]. It is hard for a sampler to “descend” in the narrow “neck” of the distribution, corresponding to areas where the variance  $\nu$  is small [BG13].

Fortunately, we can represent this model in an equivalent way that makes it easier to sample from, providing we use a **non-centered parameterization** [PR03]. This has the form

$$\nu \sim \mathcal{N}(0, 3) \tag{12.113}$$

$$z \sim \mathcal{N}(0, 1) \tag{12.114}$$

$$x = z \exp(\nu) \tag{12.115}$$

This is easier to sample from, since  $p(z, \nu)$  is a product of 2 independent Gaussians, and we can derive  $x$  deterministically from these Gaussian samples. The advantage of this reparameterization is shown in Figure 12.18. A method to automatically derive such reparameterizations is discussed in [GMH20].

## 12.7 Stochastic gradient MCMC

Consider an unnormalized target distribution of the following form:

$$\pi(\boldsymbol{\theta}) \propto p(\boldsymbol{\theta}, \mathcal{D}) = p_0(\boldsymbol{\theta}) \prod_{n=1}^N p(\mathbf{x}_n | \boldsymbol{\theta}) \quad (12.116)$$

where  $\mathcal{D} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$ . Alternatively we can define the target distribution in terms of an energy function (negative log joint) as follows:

$$p(\boldsymbol{\theta}, \mathcal{D}) \propto \exp(-\mathcal{E}(\boldsymbol{\theta})) \quad (12.117)$$

The energy function can be decomposed over data samples:

$$\mathcal{E}(\boldsymbol{\theta}) = \sum_{n=1}^N \mathcal{E}_n(\boldsymbol{\theta}) \quad (12.118)$$

$$\mathcal{E}_n(\boldsymbol{\theta}) = -\log p(\mathbf{x}_n | \boldsymbol{\theta}) - \frac{1}{N} \log p_0(\boldsymbol{\theta}) \quad (12.119)$$

Evaluating the full energy (e.g., to compute an acceptance probability in the Metropolis Hastings algorithm, or to compute the gradient in HMC) takes  $O(N)$  time, which does not scale to large data. In this section, we discuss some solutions to this problem.

### 12.7.1 Stochastic Gradient Langevin Dynamics (SGLD)

Recall from Equation (12.83) that the **Langevin diffusion** SDE has the following form

$$d\boldsymbol{\theta}_t = -\nabla \mathcal{E}(\boldsymbol{\theta}_t) dt + \sqrt{2} d\mathbf{W}_t \quad (12.120)$$

where  $d\mathbf{W}_t$  is a Wiener noise (also called Brownian noise) process. In discrete time, we can use the following Euler approximation:

$$\boldsymbol{\theta}_{t+1} \approx \boldsymbol{\theta}_t - \eta_t \nabla \mathcal{E}(\boldsymbol{\theta}_t) + \sqrt{2\eta_t} \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (12.121)$$

Computing the gradient  $\mathbf{g}(\boldsymbol{\theta}_t) = \nabla \mathcal{E}(\boldsymbol{\theta}_t)$  at each step takes  $O(N)$  time. We can compute an unbiased minibatch approximation to the gradient term in  $O(B)$  time using

$$\hat{\mathbf{g}}(\boldsymbol{\theta}_t) = \frac{N}{B} \sum_{n \in \mathcal{B}_t} \nabla \mathcal{E}_n(\boldsymbol{\theta}_t) = -\frac{N}{B} \left( \sum_{n \in \mathcal{B}_t} \nabla \log p(\mathbf{x}_n | \boldsymbol{\theta}_t) + \frac{B}{N} \nabla \log p_0(\boldsymbol{\theta}_t) \right) \quad (12.122)$$

1 where  $\mathcal{B}_t$  is the minibatch at step  $t$ . This gives rise to the following approximate update:  
 2

$$3 \quad \boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_t \hat{\mathbf{g}}(\boldsymbol{\theta}_t) + \sqrt{2\eta_t} \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (12.123)$$

4 This is called **stochastic gradient Langevin dynamics** or **SGLD** [Wel11]. The resulting update  
 5 step is identical to SGD, except for the addition of a Gaussian noise term. (See [Neg+21] for some  
 6 recent analysis of this method; they also suggest setting  $\eta_t \propto N^{-2/3}$ .)  
 7

### 8    12.7.2 Preconditionining

9 As in SGD, we can get better results (especially for models such as neural networks) if we use  
 10 preconditioning to scale the gradient updates. In [PT13], they use the Fisher information matrix  
 11 (FIM) as the preconditioner; this method is known as **Stochastic Gradient Riemannian Langevin**  
 12 **Dynamics** or **SGRLD**.  
 13

14 Unfortunately, computing the FIM is often hard. In [Li+16], they propose to use the same kind of  
 15 diagonal approximation as used by RMSprop; this is called **preconditioned SGLD**. An alternative  
 16 is to use an Adam-like preconditioner, as proposed in [KSL21]. This is called **SGLD-Adam**. For  
 17 more details, see [CSN21].  
 18

### 19    12.7.3 Reducing the variance of the gradient estimate

20 The variance of the noise introduced by minibatching can be quite large, which can hurt the  
 21 performance of methods such as SGLD [BDM18]. In [Bak+17], they propose to reduce the variance  
 22 of this estimate by using a **control variate** estimator; this method is therefore called **SGLD-CV**.  
 23 Specifically they use the following gradient approximation:  
 24

$$25 \quad \hat{\nabla}_{cv} \mathcal{E}(\boldsymbol{\theta}_t) = \nabla \mathcal{E}(\hat{\boldsymbol{\theta}}) + \frac{N}{B} \sum_{n \in \mathcal{S}_t} (\nabla \mathcal{E}_n(\boldsymbol{\theta}_t) - \nabla \mathcal{E}_n(\hat{\boldsymbol{\theta}})) \quad (12.124)$$

26 Here  $\hat{\boldsymbol{\theta}}$  is any fixed value, but it is often taken to be an approximate MAP estimate (e.g., based on  
 27 one epoch of SGD). The reason Equation (12.124) is valid is because the terms we add and subtract  
 28 are equal in expectation, and hence we get an unbiased estimate:  
 29

$$30 \quad \mathbb{E} [\hat{\nabla}_{cv} \mathcal{E}(\boldsymbol{\theta}_t)] = \nabla \mathcal{E}(\hat{\boldsymbol{\theta}}) + \mathbb{E} \left[ \frac{N}{B} \sum_{n \in \mathcal{S}_t} (\nabla \mathcal{E}_n(\boldsymbol{\theta}_t) - \nabla \mathcal{E}_n(\hat{\boldsymbol{\theta}})) \right] \quad (12.125)$$

$$31 \quad = \nabla \mathcal{E}(\hat{\boldsymbol{\theta}}) + \nabla \mathcal{E}(\boldsymbol{\theta}_t) - \nabla \mathcal{E}(\hat{\boldsymbol{\theta}}) = \nabla \mathcal{E}(\boldsymbol{\theta}_t) \quad (12.126)$$

32 Note that the first term,  $\nabla \mathcal{E}(\hat{\boldsymbol{\theta}}) = \sum_{n=1}^N \nabla \mathcal{E}_n(\hat{\boldsymbol{\theta}})$ , requires a single pass over the entire dataset, but  
 33 only has to be computed once (e.g., while estimating  $\hat{\boldsymbol{\theta}}$ ).  
 34

35 One disadvantage of SGLD-CV is that the reference point  $\hat{\boldsymbol{\theta}}$  has to be precomputed, and is then  
 36 fixed. An alternative is to update the reference point online, by performing periodic full batch  
 37 estimates. This is called **SVRG-LD** [Dub+16; Cha+18], where SVRG stands for stochastic variance  
 38 reduced gradient, and LD stands for Langevin Dynamics. If we use  $\tilde{\boldsymbol{\theta}}_t$  to denote the most recent  
 39 snapshot (reference point), the corresponding gradient estimate is given by  
 40

$$41 \quad \hat{\nabla}_{svrg} \mathcal{E}(\boldsymbol{\theta}_t) = \nabla \mathcal{E}(\tilde{\boldsymbol{\theta}}_t) + \frac{N}{B} \sum_{n \in \mathcal{S}_t} (\nabla \mathcal{E}_n(\boldsymbol{\theta}_t) - \nabla \mathcal{E}_n(\tilde{\boldsymbol{\theta}}_t)) \quad (12.127)$$

We recompute the snapshot every  $\tau$  steps (known as the epoch length). See Algorithm 29 for the pseudo-code.

---

**Algorithm 29:** SVRG Langevin Descent

---

```

1 Initialize  $\theta_0$ 
2 for  $t = 1 : T$  do
3   if  $t \bmod \tau = 0$  then
4      $\tilde{\theta} = \theta_t$ 
5      $\tilde{g} = \sum_{n=1}^N \mathcal{E}_n(\tilde{\theta})$ 
6     Sample minibatch  $\mathcal{B}_t \in \{1, \dots, N\}$ 
7      $g_t = \tilde{g} + \frac{N}{B} \sum_{n \in \mathcal{B}_t} (\nabla \mathcal{E}_n(\theta_t) - \nabla \mathcal{E}_n(\tilde{\theta}))$ 
8      $\theta_{t+1} = \theta_t - \eta_t g_t + \sqrt{2\eta_t} \mathcal{N}(\mathbf{0}, \mathbf{I})$ 

```

---

The disadvantage of SVRG is that it needs to perform a full pass over the data every  $\tau$  steps. An alternative approach, called **SAGA-LD** [Dub+16; Cha+18] (which stands for stochastic averaged gradient acceleration), avoids this by storing all  $N$  gradient vectors, and then doing incremental updates. Unfortunately the memory requirements of this algorithm usually make it impractical.

#### 12.7.4 SG-HMC

We discussed Hamiltonian Monte Carlo (HMC) in Section 12.5, which uses auxiliary momentum variables to improve performance over Langevin MC. In this section, we discuss a way to speed it up by approximating the gradients using minibatches. This is called called **SG-HMC** [CFG14; ZG21], where SG stands for “stochastic gradient”.

Recall that the leapfrog updates have the following form:

$$\mathbf{v}_{t+1/2} = \mathbf{v}_t - \frac{\eta}{2} \nabla \mathcal{E}(\theta_t) \quad (12.128)$$

$$\theta_{t+1} = \theta_t + \eta \mathbf{v}_{t+1/2} = \theta_t + \eta \mathbf{v}_t - \frac{\eta}{2} \nabla \mathcal{E}(\theta_t) \quad (12.129)$$

$$\mathbf{v}_{t+1} = \mathbf{v}_{t+1/2} - \frac{\eta}{2} \nabla \mathcal{E}(\theta_{t+1}) = \mathbf{v}_t - \frac{\eta}{2} \nabla \mathcal{E}(\theta_t) - \frac{\eta}{2} \nabla \mathcal{E}(\theta_{t+1}) \quad (12.130)$$

We can replace the full batch gradient with a stochastic approximation, to get

$$\theta_{t+1} = \theta_t + \eta \mathbf{v}_t - \frac{\eta^2}{2} \mathbf{g}(\theta_t, \xi_t) \quad (12.131)$$

$$\mathbf{v}_{t+1} = \mathbf{v}_t - \frac{\eta}{2} \mathbf{g}(\theta_t, \xi_t) - \frac{\eta}{2} \mathbf{g}(\theta_{t+1}, \xi_{t+1/2}) \quad (12.132)$$

where  $\xi_t$  and  $\xi_{t+1/2}$  are independent sources of randomness (e.g., batch indices). In [ZG21], they show that this algorithm (even without the MH rejection step) provides a good approximation to the posterior (in the sense of having small Wasserstein-2 distance) for the case where the energy function is strongly convex. Furthermore, performance can be considerably improved if we use the variance reduction methods discussed in Section 12.7.3.

---

1 **12.7.5 Underdamped Langevin Dynamics**

3 The **underdamped Langevin dynamics (ULD)** has the form of the following SDE [CDC15;  
4 LMS16; Che+18a; Che+18d]:  
5

6  $d\theta_t = v_t dt$   
7  $dv_t = -g(\theta_t)dt - \gamma v_t dt + \sqrt{2\gamma} dW_t$  (12.133)  
8

9 where  $g(\theta_t) = \nabla \mathcal{E}(\theta_t)$  is the gradient or **force** acting on the particle,  $\gamma > 0$  is the **friction** parameter,  
10 and  $dW_t$  is Wiener noise.

11 Equation (12.133) is like the Langevin dynamics of Equation (12.83) but with an added momentum  
12 term  $v_t$ . We can solve the dynamics using various integration methods. It can be shown (see e.g.,  
13 [LMS16]) that these methods are accurate to second order, whereas solving standard (overdamped)  
14 Langevin is only accurate to first order, and thus will require more sampling steps to achieve a given  
15 accuracy.

16

17 **12.8 Reversible jump (trans-dimensional) MCMC**  
18

19 Suppose we have a set of models with different numbers of parameters, e.g., mixture models in which  
20 the number of mixture components is unknown. Let the model be denoted by  $m$ , and let its unknowns  
21 (e.g., parameters) be denoted by  $x_m \in \mathcal{X}_m$  (e.g.,  $\mathcal{X}_m = \mathbb{R}^{n_m}$ , where  $n_m$  is the dimensionality of  
22 model  $m$ ). Sampling in spaces of differing dimensionality is called **trans-dimensional MCMC**. We  
23 could sample the model indicator  $m \in \{1, \dots, M\}$  and sample all the parameters from the product  
24 space  $\prod_{m=1}^M \mathcal{X}_m$ , but this is very inefficient, and only works if  $M$  is finite. It is more parsimonious  
25 to sample in the union space  $\mathcal{X} = \cup_{m=1}^M \{m\} \times \mathcal{X}_m$ , where we only worry about parameters for the  
26 currently active model.

27 The difficulty with this approach arises when we move between models of different dimensionality.  
28 The trouble is that when we compute the MH acceptance ratio, we are comparing densities defined  
29 on spaces of different dimensionality, which is not well defined. For example, comparing densities on  
30 two points of a sphere makes sense, but comparing a density on a sphere to a density on a circle  
31 does not, as there is a dimensional mismatch in the two concepts. The solution, proposed by [Gre98]  
32 and known as **reversible jump MCMC** or **RJMCMC**, is to augment the low dimensional space  
33 with extra random variables so that the two spaces have a common measure. This is illustrated in  
34 Figure 12.19.

35 We give a sketch of the algorithm below. For more details, see e.g., [Gre03; HG12].  
36

37 **12.8.1 Basic idea**  
38

39 To explain the method in more detail, we follow the presentation of [And+03]. To ensure a common  
40 measure, we need to define a way to extend each pair of subspaces  $\mathcal{X}_m$  and  $\mathcal{X}_n$  to  $\mathcal{X}_{m,n} = \mathcal{X}_m \times \mathcal{U}_{m,n}$   
41 and  $\mathcal{X}_{n,m} = \mathcal{X}_n \times \mathcal{U}_{n,m}$ . We also need to define a deterministic, differentiable and invertible mapping

42  $(x_m, u_{m,n}) = f_{n \rightarrow m}(x_n, u_{n,m}) = (f_{n \rightarrow m}^x(x_n, u_{n,m}), f_{n \rightarrow m}^u(x_n, u_{n,m}))$  (12.134)  
43

44 Invertibility means that

45  $f_{m \rightarrow n}(f_{n \rightarrow m}(x_n, u_{n,m})) = (x_n, u_{n,m})$  (12.135)  
46

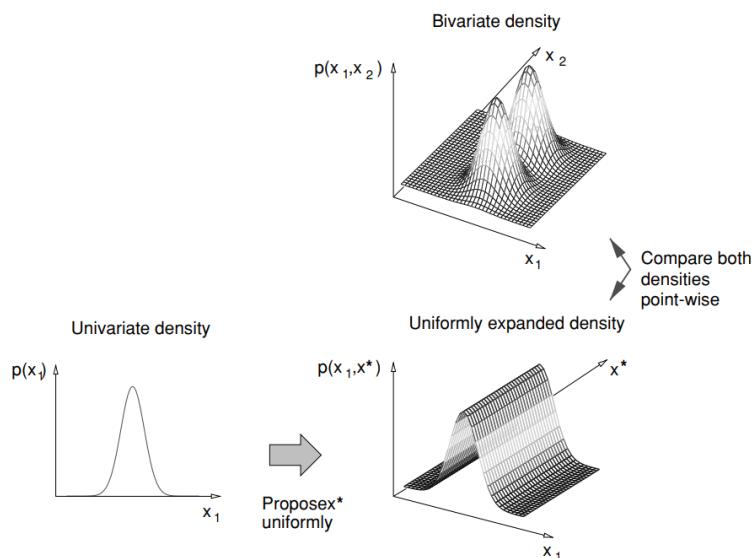


Figure 12.19: To compare a 1d model against a 2d model, we first have to map the 1d model to 2d space so the two have a common measure. Note that we assume the ridge has finite support, so it is integrable. From Figure 17 of [And+03]. Used with kind permission of Nando de Freitas.

Finally, we need to define proposals  $q_{n \rightarrow m}(\mathbf{u}_{n,m}|n, \mathbf{x}_n)$  and  $q_{m \rightarrow n}(\mathbf{u}_{m,n}|m, \mathbf{x}_m)$ .

Suppose we are in state  $(n, \mathbf{x}_n)$ . We move to  $(m, \mathbf{x}_m)$  by generating  $\mathbf{u}_{n,m} \sim q_{n \rightarrow m}(\cdot|n, \mathbf{x}_n)$ , and then computing  $(\mathbf{x}_m, \mathbf{u}_{m,n}) = f_{n \rightarrow m}(\mathbf{x}_n, \mathbf{u}_{n,m})$ . We then accept the move with probability

$$A_{n \rightarrow m} = \min \left\{ 1, \frac{p(m, \mathbf{x}_m^*)}{p(n, \mathbf{x}_n)} \times \frac{q(n|m)}{q(m|n)} \times \frac{q_{m \rightarrow n}(\mathbf{u}_{m,n}|m, \mathbf{x}_m^*)}{q_{n \rightarrow m}(\mathbf{u}_{n,m}|n, \mathbf{x}_n)} \times |\det \mathbf{J}_{f_{m \rightarrow n}}| \right\} \quad (12.136)$$

where  $\mathbf{x}_m^* = f_{n \rightarrow m}^x(\mathbf{x}_n, \mathbf{u}_{n,m})$ ,  $\mathbf{J}_{f_{m \rightarrow n}}$  is the Jacobian of the transformation

$$\mathbf{J}_{f_{m \rightarrow n}} = \frac{\partial f_{n \rightarrow m}(\mathbf{x}_m, \mathbf{u}_{m,n})}{\partial (\mathbf{x}_m, \mathbf{u}_{m,n})} \quad (12.137)$$

and  $|\det \mathbf{J}|$  is the absolute value of the determinant of the Jacobian.

### 12.8.2 Example

Let us consider an example from [AFD01]. They consider an RBF network for nonlinear regression of the form

$$f(\mathbf{x}) = \sum_{j=1}^k a_j \mathcal{K}(\|\mathbf{x} - \boldsymbol{\mu}_j\|) + \boldsymbol{\beta}^\top \mathbf{x} + \beta_0 + \epsilon \quad (12.138)$$

where  $\mathcal{K}()$  is some kernel function (e.g., a Gaussian),  $k$  is the number of such basis functions, and  $\epsilon$  is a Gaussian noise term. If  $k = 0$ , the model corresponds to linear regression.

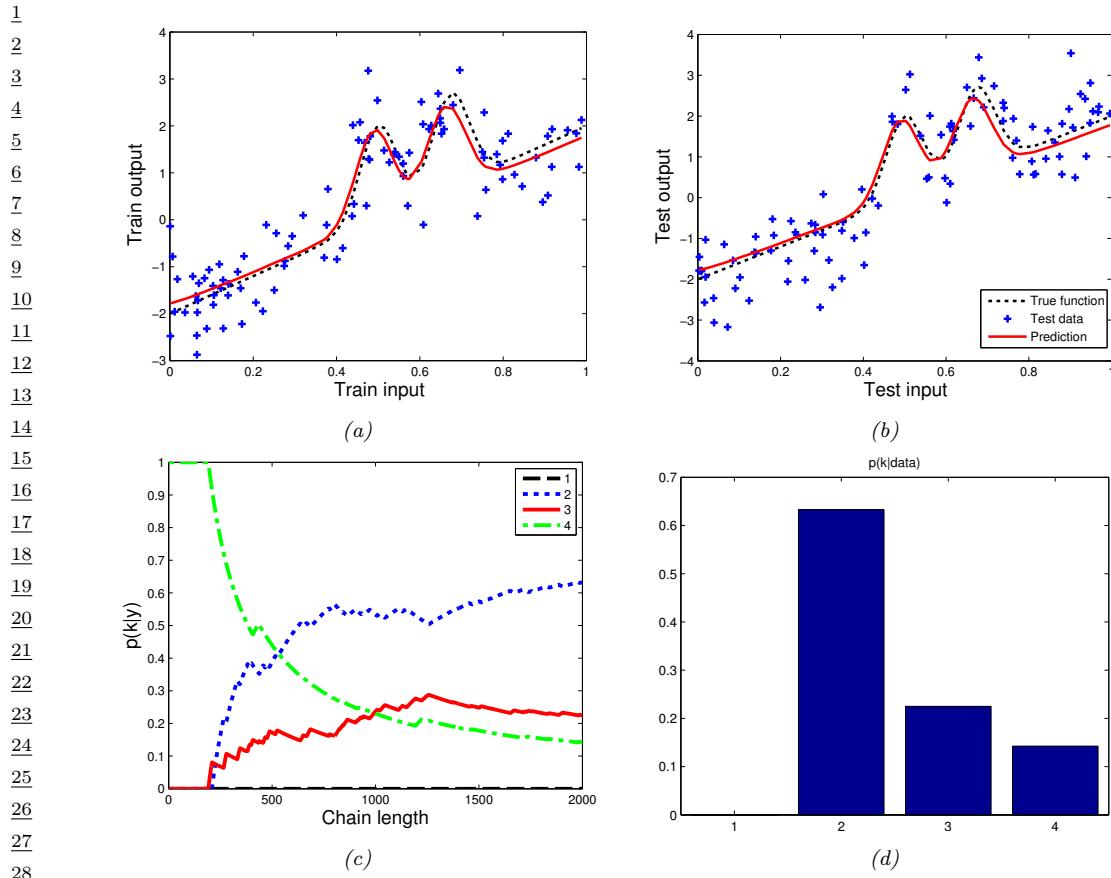


Figure 12.20: Fitting an RBF network to some 1d data using RJMCMC. (a) Prediction on train set. (b) Prediction on test set. (c) Plot of  $p(k|D)$  vs iteration. (d) Final posterior  $p(k|D)$ . Adapted from Figure 4 of [AFD01]. Generated by [rjmcme\\_rbf](#), written by Nando de Freitas.

They fit this model to the data in Figure 12.20(a). The predictions on the test set are shown in Figure 12.20(b). Estimates of  $p(k|D)$ , the (distribution over the) number of basis functions, are shown in Figure 12.20(c) as a function of the iteration number; the posterior at the final iteration is shown in Figure 12.20(d). There is clearly the most posterior support for  $k = 2$ , which makes sense given the two “bumps” in the data.

To generate these results, they consider several kinds of proposal. One of them is to split a current basis function  $\mu$  into two new ones using

$$\mu_1 = \mu - u_{n,n+1}\alpha, \quad \mu_2 = \mu + u_{n,n+1}\alpha \tag{12.139}$$

where  $\alpha$  is a parameter of the proposal, and  $u_{n,m}$  is sampled from some distribution (e.g., uniform).

To ensure reversibility, they define a corresponding merge move

$$\mu = \frac{\mu_1 + \mu_2}{2} \tag{12.140}$$

where  $\mu_1$  is chosen at random, and  $\mu_2$  is its nearest neighbor. To ensure these moves are reversible, we require  $\|\mu_1 - \mu_2\| < 2\beta$ .

The acceptance ratio for the split move is given by

$$A_{split} = \min \left\{ 1, \frac{p(k+1, \mu_{k+1})}{p(k, \mu_k)} \times \frac{1/(k+1)}{1/k} \times \frac{1}{p(u_{n,m})} \times |\det \mathbf{J}_{split}| \right\} \quad (12.141)$$

where  $1/k$  is the probability of choosing one of the  $k$  bases uniformly at random. The Jacobian is

$$\mathbf{J}_{split} = \frac{\partial(\mu_1, \mu_2)}{\partial(\mu, u_{n,m})} = \det \begin{pmatrix} 1 & 1 \\ -\beta & \beta \end{pmatrix} \quad (12.142)$$

so  $|\det \mathbf{J}_{split}| = 2\beta$ . The acceptance ratio for the merge move is given by

$$A_{merge} = \min \left\{ 1, \frac{p(k-1, \mu_{k-1})}{p(k, \mu_k)} \times \frac{1/(k-1)}{1/k} \times |\det \mathbf{J}_{merge}| \right\} \quad (12.143)$$

where  $|\det \mathbf{J}_{merge}| = 1/(2\beta)$ .

---

**Algorithm 30:** Generic reversible jump MCMC (single step)

---

- 1 Sample  $u \sim U(0, 1)$
  - 2 If  $u \leq b_k$
  - 3 then birth move
  - 4 else if  $u \leq (b_k + d_k)$  then death move
  - 5 else if  $u \leq (b_k + d_k + s_k)$  then split move
  - 6 else if  $u \leq (b_k + d_k + s_k + m_k)$  then merge move
  - 7 else update parameters
- 

The overall pseudo-code for the algorithm, assuming the current model has index  $k$ , is given in Algorithm 30. Here  $b_k$  is the probability of a birth move,  $d_k$  is the probability of a death move,  $s_k$  is the probability of a split move, and  $m_k$  is the probability of a merge move. If we don't make a dimension-changing move, we just update the parameters of the current model using random walk MH.

### 12.8.3 Discussion

RJMCMC algorithms can be quite tricky to implement. If, however, the continuous parameters can be integrated out (resulting in a method called collapsed RJMCMC), much of the difficulty goes away, since we are just left with a discrete state space, where there is no need to worry about change of measure. For example, if we fix the centers  $\mu_j$  in Equation (12.138) (e.g., using samples from the data, or using K-means clustering), we are left with a linear model, where we can integrate out the parameters. All that is left to do is sample which of these fixed basis functions to include in the model, which is a discrete variable selection problem. See e.g., [Den+02] for details.

In Chapter 31, we discuss Bayesian nonparametric models, which allow for an infinite number of different models. Surprisingly, such models are often easier to deal with computationally (as well as more realistic, statistically) than working with a finite set of different models.

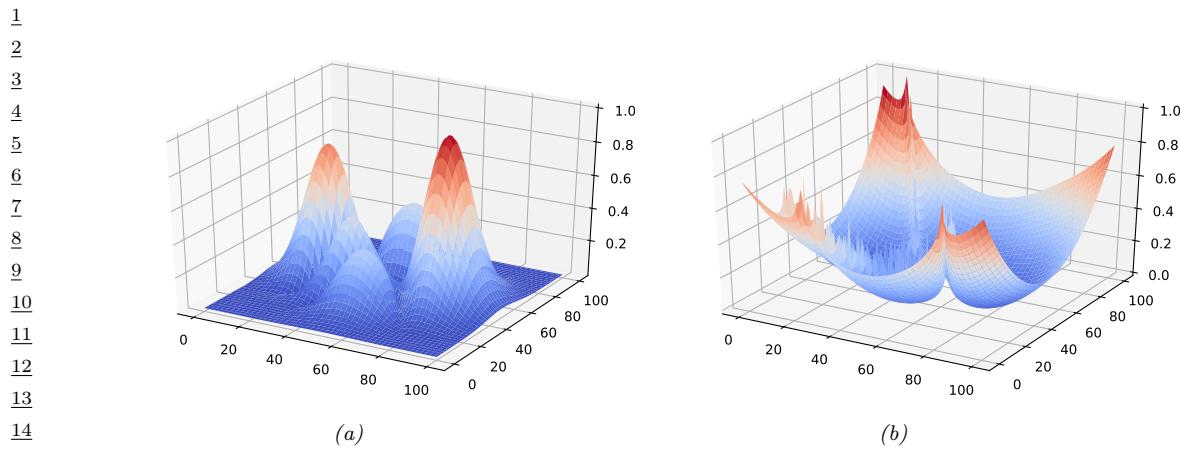


Figure 12.21: (a) A peaky distribution. (b) Corresponding energy function. Generated by [simulated\\_annealing\\_2d\\_demo.ipynb](#).

## 12.9 Annealing methods

Many distributions are multimodal and hence hard to sample from. However, by analogy to the way metals are heated up and then cooled down in order to make the molecules align, we can imagine using a computational temperature parameter to “smooth out” a distribution, gradually cooling it to recover the original “bumpy” distribution. We first explain this idea in more detail in the context of an algorithm for MAP estimation. We then discuss extensions to the sampling case.

### 12.9.1 Simulated annealing

In this section, we discuss the **simulated annealing** algorithm [KJV83; LA87], which is a variant of the Metropolis Hastings algorithm which is designed to find the global optimum of blackbox function. (Other approaches to blackbox optimization are discussed in Section 6.9.)

Annealing is a physical process of heating a solid until thermal stresses are released, then cooling it very slowly until the crystals are perfectly arranged, achieving a minimum energy state. Depending on how fast or slow the temperature is cooled, the results will have worse or better the quality. We can apply this approach to probability distributions, to control the number of modes (low energy states) that they have, by defining

$$p_T(\mathbf{x}) = \exp(-\mathcal{E}(\mathbf{x})/T) \quad (12.144)$$

where  $T$  is the temperature, which is reduced over time. As an example, consider the **peaks** function:

$$p(x, y) \propto |3(1-x)^2 e^{-x^2-(y+1)^2} - 10(\frac{x}{5} - x^3 - y^5) e^{-x^2-y^2} - \frac{1}{3} e^{-(x+1)^2-y^2}| \quad (12.145)$$

This is plotted in Figure 12.21a. The corresponding energy is in Figure 12.21b. We plot annealed versions of this distribution in Figure 12.22. At high temperatures,  $T \gg 1$ , the surface is approximately

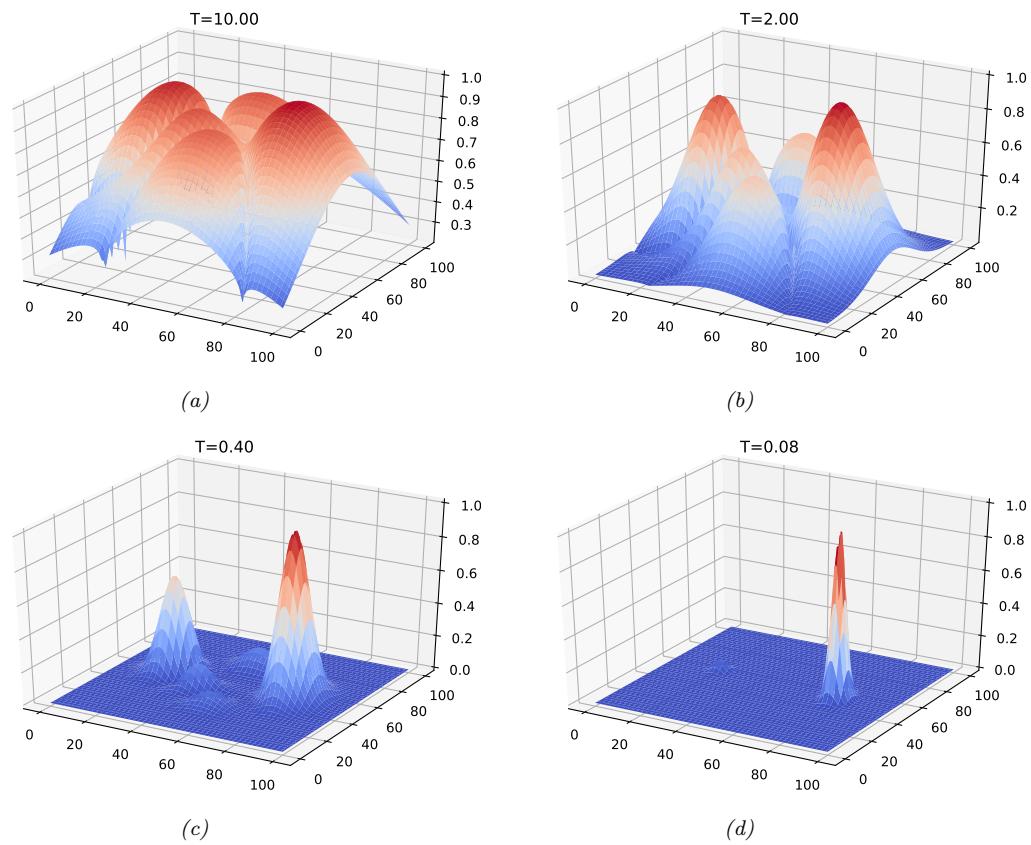


Figure 12.22: Annealed version of the distribution in Figure 12.21a at different temperatures. Generated by [simulated\\_annealing\\_2d\\_demo.ipynb](#).

flat, and hence it is easy to move around (i.e., to avoid local optima). As the temperature cools, the largest peaks become larger, and the smallest peaks disappear. By cooling slowly enough, it is possible to “track” the largest peak, and thus find the global optimum (minimum energy state). This is an example of a **continuation method**.

In more detail, at each step, we sample a new state according to some proposal distribution  $\mathbf{x}' \sim q(\cdot | \mathbf{x}_t)$ . For real-valued parameters, this is often simply a random walk proposal centered on the current iterate,  $\mathbf{x}' = \mathbf{x}_t + \boldsymbol{\epsilon}_{t+1}$ , where  $\boldsymbol{\epsilon}_{t+1} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$ . (The matrix  $\boldsymbol{\Sigma}$  is often diagonal, and may be updated over time using the method in [Cor+87].) Having proposed a new state, we compute the acceptance probability

$$\alpha_{t+1} = \exp(-(\mathcal{E}(\mathbf{x}') - \mathcal{E}(\mathbf{x}_t))/T_t) \quad (12.146)$$

where  $T_t$  is the temperature of the system. We then accept the new state (i.e., set  $\mathbf{x}_{t+1} = \mathbf{x}'$ ) with probability  $\min(1, \alpha_{t+1})$ , otherwise we stay in the current state (i.e., set  $\mathbf{x}_{t+1} = \mathbf{x}_t$ ). This means that

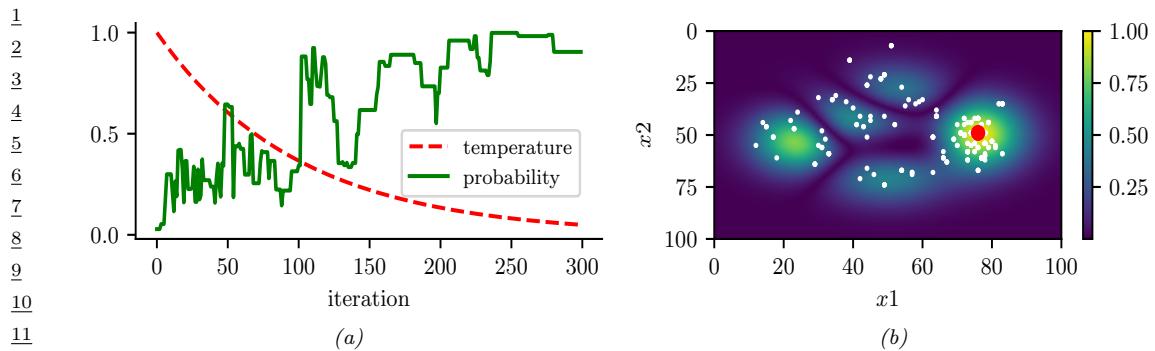


Figure 12.23: Simulated annealing applied to the distribution in Figure 12.21a. (a) Temperature vs iteration and probability of each visited point vs iteration. (b) Visited samples, superimposed on the target distribution. The big red dot is the highest probability point found. Generated by [simulated\\_annealing\\_2d\\_demo.ipynb](#).

if the new state has lower energy (is more probable), we will definitely accept it, but if it has higher energy (is less probable), we might still accept, depending on the current temperature. Thus the algorithm allows “downhill” moves in probability space (uphill in energy space), but less frequently as the temperature drops.

The rate at which the temperature changes over time is called the **cooling schedule**. It has been shown [Haj88] that if one cools according to a logarithmic schedule,  $T_t \propto 1/\log(t+1)$ , then the method is guaranteed to find the global optimum under certain assumptions. However, this schedule is often too slow. In practice it is common to use an **exponential cooling schedule** of the form  $T_{t+1} = \gamma T_t$ , where  $\gamma \in (0, 1]$  is the cooling rate. Cooling too quickly means one can get stuck in a local maximum, but cooling too slowly just wastes time. The best cooling schedule is difficult to determine; this is one of the main drawbacks of simulated annealing.

In Figure 12.23a, we show a cooling schedule using  $\gamma = 0.9$ . If we combine this with a Gaussian random walk proposal with  $\sigma = 10$  to the peaky distribution in Figure 12.21a, we get the results shown in Figure 12.23 and Figure 12.23b. We see that the algorithm concentrates its samples near the global optimum (the peak on the middle right).

### 12.9.2 Parallel tempering

Another way to combine MCMC and annealing is to run multiple chains in parallel at different temperatures, and allow one chain to sample from another chain at a neighboring temperature. In this way, the high temperature chain can make long distance moves through the state space, and have this influence lower temperature chains. This is known as **parallel tempering**. See e.g., [ED05; Kat+06] for details.

# 13 Sequential Monte Carlo inference

## 13.1 Introduction

In this chapter, we discuss **sequential Monte Carlo** or **SMC** algorithms, which can be used to sample from a sequence of related probability distributions. SMC is most commonly used to solve filtering in state-space models (SSM, Chapter 29), but it can also be applied to other problems, such as sampling from a static (but possibly multi-modal) distribution, or for sampling rare events from some process.

Our presentation is based on the excellent tutorial [NLS19], and differs from traditional presentations, such as [Aru+02], by emphasizing the fact that we are sampling sequences of related variables, not just computing the filtering distribution of an SSM. This more general perspective will let us tackle static estimation problems, as we will see. For another good introduction to SMC, see [DJ11]. For a more formal (measure theoretic) treatment of SMC, using the **Feynman-Kac** formalism, see [CP20b].

### 13.1.1 Problem statement

In SMC, the goal is to sample from a sequence of related distributions of the form

$$\pi_t(\mathbf{z}_{1:t}) = \frac{1}{Z_t} \tilde{\gamma}_t(\mathbf{z}_{1:t}) \tag{13.1}$$

for  $t = 1 : T$ , where  $\tilde{\gamma}_t$  is the unnormalized **target distribution**,  $\pi_t$  is the normalized version, and  $\mathbf{z}_{1:t}$  are the random variables of interest. In some applications (e.g., filtering in an SSM), we care about each intermediate marginal distribution,  $\pi_t(\mathbf{z}_t)$ , for  $t = 1 : T$ ; this is called **particle filtering**. (The word “particle” just means “sample”.) In other applications, we only care about the final distribution,  $\pi_T(\mathbf{z}_T)$ , and the intermediate steps are introduced just for computational reasons; this is called an **SMC sampler**. We briefly review both of these below, and go into more detail in later sections.

### 13.1.2 Particle filtering for state-space models

An important application of SMC is to sequential (online) inference (state estimation) in SSMs. As an example, consider a Markovian state-space model with the following joint distribution:

$$\pi_T(\mathbf{z}_{1:T}) \propto p(\mathbf{z}_{1:T}, \mathbf{y}_{1:T}) = p(\mathbf{z}_1)p(\mathbf{y}_1|\mathbf{z}_1) \prod_{t=1}^T p(\mathbf{z}_t|\mathbf{z}_{t-1})p(\mathbf{y}_t|\mathbf{z}_t) \tag{13.2}$$

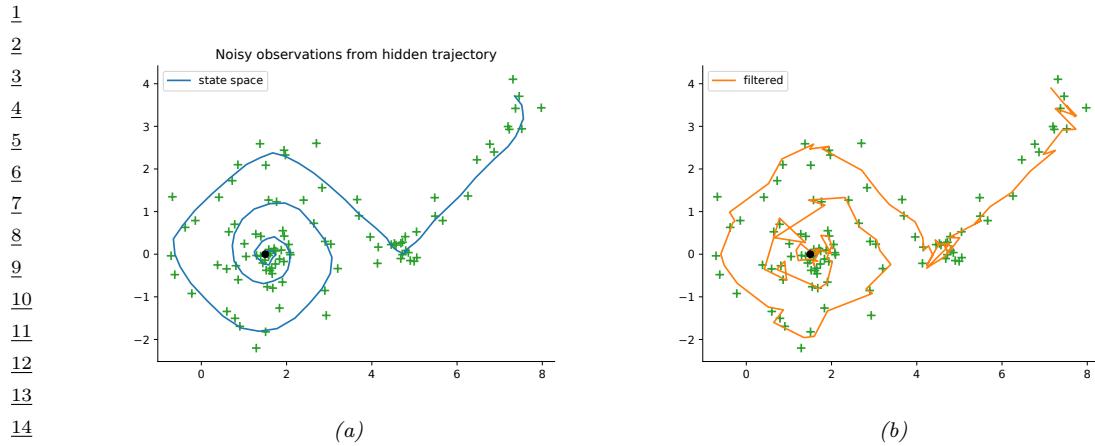


Figure 13.1: Illustration of particle filtering (using the dynamical prior as the proposal) applied to a 2d nonlinear dynamical system. (a) True underlying state and observed data. (b) PF estimate of the posterior mean. Generated by `bootstrap_filter_spiral.ipynb`.

A common choice is to define the unnormalized target distribution at step  $t$  to be

$$\tilde{\gamma}_t(\mathbf{z}_{1:t}) = p(\mathbf{z}_{1:t}, \mathbf{y}_{1:t}) = p(\mathbf{z}_1)p(\mathbf{y}_1|\mathbf{z}_1) \prod_{s=1}^t p(\mathbf{z}_s|\mathbf{z}_{s-1})p(\mathbf{y}_s|\mathbf{z}_s) \quad (13.3)$$

Note that this is a distribution over an (ever growing) sequence of latent variables. However, we often only care about the most recent marginal of this distribution, in which case we just need to compute  $\tilde{\gamma}_t(\mathbf{z}_t)$ , which avoids having to store the full history.

For example, consider the following 2d nonlinear tracking problem (the same one as in Section 8.5.5.1):

$$\begin{aligned} p(\mathbf{z}_t|\mathbf{z}_{t-1}) &= \mathcal{N}(\mathbf{z}_t|f(\mathbf{z}_{t-1}), q\mathbf{I}) \\ p(\mathbf{y}_t|\mathbf{z}_t) &= \mathcal{N}(\mathbf{y}_t|\mathbf{z}_t, r\mathbf{I}) \\ f(\mathbf{z}) &= (z_1 + \Delta \sin(z_2), z_2 + \Delta \cos(z_1)) \end{aligned} \quad (13.4)$$

where  $\Delta$  is the step size of the underlying continuous system,  $q$  is the variance of the system noise, and  $r$  is the variance of the observation noise. (We treat  $\Delta$ ,  $q$  and  $r$  as fixed constants; see Supplementary Section 13.1.3 for a discussion of joint state and parameter estimation.) The true underlying state trajectory, and the corresponding noisy measurements, are shown in Figure 13.1a. The posterior mean estimate of the state, computed using 2000 samples in a simple form of SMC called the bootstrap filter (Section 13.2.3.1), is shown in Figure 13.1b.

Particle filtering can also be applied to **non-Markovian models**, where  $\mathbf{z}_t$  may depend on all the past hidden states,  $\mathbf{z}_{1:t-1}$ , and  $\mathbf{y}_t$  depends on the current  $\mathbf{z}_t$  and possibly also all the past hidden states,  $\mathbf{z}_{1:t-1}$ , and optionally the past observations,  $\mathbf{y}_{1:t-1}$ . In this case, the unnormalized target

1 distribution at step  $t$  is  
2

$$\tilde{\gamma}_t(\mathbf{z}_{1:t}) = p(\mathbf{z}_1)p(\mathbf{y}_1|\mathbf{z}_1) \prod_{s=1}^t p(\mathbf{z}_s|\mathbf{z}_{1:s-1})p(\mathbf{y}_s|\mathbf{z}_{1:s}) \quad (13.5)$$

3  
4 For example, consider a 1d Gaussian sequence model where the dynamics are first-order Markov, but  
5 the observations depend on the entire past sequence (this is example 1.2.1 from [NLS19]):  
6

$$\begin{aligned} p(z_t|\mathbf{z}_{1:t-1}) &= \mathcal{N}(z_t|\phi z_{t-1}, q) \\ p(y_t|\mathbf{z}_{1:t}) &= \mathcal{N}(y_t| \sum_{s=1}^t \beta^{t-s} z_s, r) \end{aligned} \quad (13.6)$$

7 If we set  $\beta = 0$ , we get  $p(y_t|\mathbf{z}_{1:t}) = \mathcal{N}(y_t|z_t, r)$  (where we define  $0^0 = 1$ ), so the model becomes a  
8 linear-Gaussian SSM. As  $\beta$  gets larger, the dependence on the past increases, making the inference  
9 problem harder. (We will revisit this example below.)  
10

### 11 13.1.3 SMC samplers for static parameter estimation

12 Now consider the problem of parameter estimation from a fixed dataset,  $\mathcal{D} = \{\mathbf{y}_n : n = 1 : N_{\mathcal{D}}\}$ . We suppose the observations are conditionally iid, so the posterior has the form  $p(\mathbf{z}|\mathcal{D}) \propto$   
13  $p(\mathbf{z}) \prod_{n=1}^{N_{\mathcal{D}}} p(\mathbf{y}_n|\mathbf{z})$ , where  $\mathbf{z}$  is the unknown parameter. It is not immediately obvious how to  
14 approximate  $p(\mathbf{z}|\mathcal{D})$  using SMC, since we just have one distribution. However, we can convert  
15 this into a sequential inference problem in several different ways. One approach, known as **data**  
16 **tempering**, defines the (marginal) target distribution at step  $t$  as  $\tilde{\gamma}_t(\mathbf{z}_t) = p(\mathbf{z}_t)p(\mathbf{y}_{1:t}|\mathbf{z}_t)$ . In  
17 this case, the number of time steps  $T$  is the same as the number of data samples,  $N_{\mathcal{D}}$ . Another  
18 approach, known as **likelihood tempering**, defines the (marginal) target distribution at step  $t$  as  
19  $\tilde{\gamma}_t(\mathbf{z}_t) = p(\mathbf{z}_t)p(\mathcal{D}|\mathbf{z}_t)^{\tau_t}$ , where  $0 = \tau_t < \dots < \tau_T = 1$  is a temperature parameter. In this case, the  
20 number of steps  $T$  depends on how quickly we anneal the distribution from the initial prior  $p(\mathbf{z}_1)$  to  
21 the final target  $p(\mathbf{z}_T)p(\mathcal{D}|\mathbf{z}_T)$ .

22 Once we have defined the marginal target distributions  $\tilde{\gamma}_t(\mathbf{z}_t)$ , we need a way to expand this to a joint  
23 target distribution over a *sequence* of variables,  $\tilde{\gamma}_t(\mathbf{z}_{1:t})$ , so the distributions become connected to each  
24 other. We explain how to do this in Section 13.6. We can then treat the model as an SSM and apply  
25 particle filtering. At the end, we extract the final joint target distribution,  $\tilde{\gamma}_T(\mathbf{z}_{1:T}) = p(\mathbf{z}_{1:T})p(\mathcal{D}|\mathbf{z}_T)$ ,  
26 from which we can compute the marginal target distribution  $\tilde{\gamma}_T(\mathbf{z}_T) = p(\mathbf{z}_T, \mathcal{D})$ , from which we can  
27 get the posterior  $p(\mathbf{z}|\mathcal{D})$  by normalizing. We give the details in Section 13.6.  
28

## 29 13.2 Particle filtering

30 In this section, we cover the basics of SMC for state space models, culminating in a method known  
31 as the **particle filter**.  
32

### 33 13.2.1 Importance sampling

34 We start by reviewing the self-normalized importance sampling method (SNIS, Section 11.5), which  
35 is the foundation of the particle filter.  
36

Suppose we are interested in estimating the expectation of some function  $\varphi_t$  with respect to a target distribution  $\pi_t$ , which we denote by

$$\pi_t(\varphi) \triangleq \mathbb{E}_{\pi_t} [\varphi_t(\mathbf{z}_{1:t})] = \int \frac{\tilde{\gamma}_t(\mathbf{z}_{1:t})}{Z_t} \varphi_t(\mathbf{z}_{1:t}) d\mathbf{z}_{1:t} \quad (13.7)$$

where  $Z_t = \int \tilde{\gamma}_t(\mathbf{z}_{1:t}) d\mathbf{z}_{1:t}$ . Suppose we use SNIS with proposal  $q_t(\mathbf{z}_{1:t})$ . We then get the following approximation:

$$\pi_t(\varphi) \approx \frac{1}{\hat{Z}_t} \frac{1}{N_s} \sum_{i=1}^{N_s} \tilde{w}_t(\mathbf{z}_{1:t}^i) \varphi_t(\mathbf{z}_{1:t}^i) \quad (13.8)$$

where  $\mathbf{z}_{1:t}^i \stackrel{\text{iid}}{\sim} q_t$  are independent samples from the proposal,  $\tilde{w}_t^i$  are the **unnormalized weights** defined by

$$\tilde{w}_t^i = \frac{\tilde{\gamma}_t(\mathbf{z}_{1:t}^i)}{q_t(\mathbf{z}_{1:t}^i)} \quad (13.9)$$

and  $\hat{Z}_t$  is the approximate normalization constant defined by

$$\hat{Z}_t \triangleq \frac{1}{N_s} \sum_{i=1}^{N_s} \tilde{w}_t^i \quad (13.10)$$

To simplify notation, let us define the **normalized weights** by

$$W_t^i = \frac{\tilde{w}_t^i}{\sum_j \tilde{w}_t^j} \quad (13.11)$$

Then we can write

$$\mathbb{E}_{\pi_t} [\varphi_t(\mathbf{z}_{1:t})] \approx \sum_{i=1}^{N_s} W_t^i \varphi_t(\mathbf{z}_{1:t}^i) \quad (13.12)$$

Alternatively, instead of computing the expectation of a specific target function, we can just approximate the target distribution itself, using a sum of weighted samples:

$$\pi_t(\mathbf{z}_{1:t}) \approx \sum_{i=1}^{N_s} W_t^i \delta(\mathbf{z}_{1:t} - \mathbf{z}_{1:t}^i) \triangleq \hat{\pi}_t(\mathbf{z}_{1:t}) \quad (13.13)$$

The problem with importance sampling when applied in the context of sequential models is that the dimensionality of the state space is very large, and increases with  $t$ . This makes it very hard to define a good proposal that covers the high probability regions, resulting in most samples getting negligible weight. In the sections below, we discuss solutions to this problem.

---

### 13.2.2 Sequential importance sampling

In this section, we discuss **sequential importance sampling** or **SIS**, in which the proposal has the following autoregressive structure:

$$q_t(\mathbf{z}_{1:t}) = q_{t-1}(\mathbf{z}_{1:t-1})q_t(\mathbf{z}_t|\mathbf{z}_{1:t-1}) \quad (13.14)$$

We can obtain samples from  $q_{t-1}(\mathbf{z}_{1:t-1})$  by reusing the  $\mathbf{z}_{1:t-1}^i$  samples, which we then extend by one step by sampling from the conditional  $q_t(\mathbf{z}_t|\mathbf{z}_{1:t-1}^i)$ . We can think of this as “growing” the chain (sequence of states). The unnormalized weights can be computed recursively as follows:

$$\tilde{w}_t(\mathbf{z}_{1:t}) = \frac{\tilde{\gamma}_t(\mathbf{z}_{1:t})}{q_t(\mathbf{z}_{1:t})} = \frac{\tilde{\gamma}_{t-1}(\mathbf{z}_{1:t-1})}{\tilde{\gamma}_{t-1}(\mathbf{z}_{1:t-1})} \frac{\tilde{\gamma}_t(\mathbf{z}_{1:t})}{q_t(\mathbf{z}_t|\mathbf{z}_{1:t-1})q_{t-1}(\mathbf{z}_{1:t-1})} \quad (13.15)$$

$$= \frac{\tilde{\gamma}_{t-1}(\mathbf{z}_{1:t-1})}{q_{t-1}(\mathbf{z}_{1:t-1})} \frac{\tilde{\gamma}_t(\mathbf{z}_{1:t})}{\tilde{\gamma}_{t-1}(\mathbf{z}_{1:t-1})q_t(\mathbf{z}_t|\mathbf{z}_{1:t-1})} \quad (13.16)$$

$$= \tilde{w}_{t-1}(\mathbf{z}_{1:t-1}) \frac{\tilde{\gamma}_t(\mathbf{z}_{1:t})}{\tilde{\gamma}_{t-1}(\mathbf{z}_{1:t-1})q_t(\mathbf{z}_t|\mathbf{z}_{1:t-1})} \quad (13.17)$$

The ratio factors are sometimes called the **incremental importance weights**:

$$\alpha_t(\mathbf{z}_{1:t}) = \frac{\tilde{\gamma}_t(\mathbf{z}_{1:t})}{\tilde{\gamma}_{t-1}(\mathbf{z}_{1:t-1})q_t(\mathbf{z}_t|\mathbf{z}_{1:t-1})} \quad (13.18)$$

See Algorithm 31 for pseudocode for the resulting SIS algorithm. (In practice we compute the weights in log-space, and convert back using the log-sum-exp trick.)

Note that, in the special case of state space models, the weight computation can be further simplified. In particular, suppose we have

$$\tilde{\gamma}_t(\mathbf{z}_{1:t}) = p(\mathbf{z}_{1:t}, \mathbf{y}_{1:t}) = p(\mathbf{y}_t|\mathbf{z}_{1:t})p(\mathbf{z}_t|\mathbf{z}_{1:t-1})p(\mathbf{z}_{1:t-1}, \mathbf{y}_{1:t-1}) \quad (13.19)$$

$$= p(\mathbf{y}_t|\mathbf{z}_{1:t})p(\mathbf{z}_t|\mathbf{z}_{1:t-1})\tilde{\gamma}_{t-1}(\mathbf{z}_{1:t-1}) \quad (13.20)$$

Then the incremental weight is given by

$$\alpha_t(\mathbf{z}_{1:t}) = \frac{p(\mathbf{y}_t|\mathbf{z}_{1:t})p(\mathbf{z}_t|\mathbf{z}_{1:t-1})\tilde{\gamma}_{t-1}(\mathbf{z}_{1:t-1})}{\tilde{\gamma}_{t-1}(\mathbf{z}_{1:t-1})q_t(\mathbf{z}_t|\mathbf{z}_{1:t-1})} = \frac{p(\mathbf{y}_t|\mathbf{z}_{1:t})p(\mathbf{z}_t|\mathbf{z}_{1:t-1})}{q_t(\mathbf{z}_t|\mathbf{z}_{1:t-1})} \quad (13.21)$$

Unfortunately SIS suffers from a problem known as **weight degeneracy** or **particle impoverishment**, in which most of the weights become very small (near zero), so the posterior ends up being approximated by a single particle. This is illustrated in Figure 13.2a, where we apply SIS to the non-Markovian example in Equation (13.6) using  $N_s = 5$  particles. The reason for degeneracy is that each particle has to “explain” (generate) the entire sequence of observations. Each sequence of guessed states becomes increasingly improbable over time, due to the product of likelihood terms, and the differences between the weights of each hypothesis will grow exponentially. Of course, there has to be a best sequence amongst the set of candidates, so when we normalize the weights, the best one will get weight 1 and the rest will get weight 0. But this is a waste of most of the particles. We discuss a solution to this in Section 13.2.3.

---

**Algorithm 31:** Sequential importance sampling (SIS)

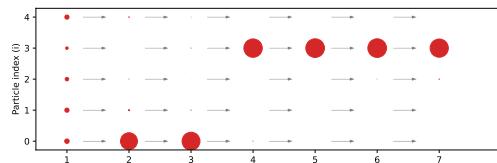
---

```

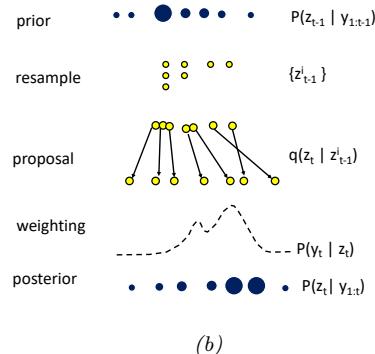
1
2   Initialization:  $\mathbf{z}_1^i \sim q_1(\mathbf{z}_1)$ ,  $\tilde{w}_1^i = \frac{\tilde{\gamma}_1(\mathbf{z}_1^i)}{q_1(\mathbf{z}_1^i)}$ ,  $W_1^i = \frac{\tilde{w}_1^i}{\sum_j \tilde{w}_1^j}$ ,  $\hat{\pi}_1(\mathbf{z}_1) = \sum_{i=1}^{N_s} W_1^i \delta(\mathbf{z}_1 - \mathbf{z}_1^i)$ 
3
4   for  $t = 2 : T$  do
5     for  $i = 1 : N_s$  do
6       Sample  $\mathbf{z}_t^i \sim q_t(\mathbf{z}_t | \mathbf{z}_{1:t-1}^i)$ 
7       Compute incremental weight  $\alpha_t^i = \frac{\tilde{\gamma}_t(\mathbf{z}_{1:t}^i)}{\tilde{\gamma}_{t-1}(\mathbf{z}_{1:t-1}^i) q_t(\mathbf{z}_t | \mathbf{z}_{1:t-1}^i)}$ 
8       Compute unnormalized weight  $\tilde{w}_t^i = \tilde{w}_{t-1}^i \alpha_t^i$ 
9
10    Compute normalized weights  $W_t^i = \frac{\tilde{w}_t^i}{\sum_j \tilde{w}_t^j}$  for  $i = 1 : N_s$ 
11
12    Compute MC posterior  $\hat{\pi}_t(\mathbf{z}_{1:t}) = \sum_{i=1}^{N_s} W_t^i \delta(\mathbf{z}_{1:t} - \mathbf{z}_{1:t}^i)$ 
13
14

```

---



(a)



(b)

Figure 13.2: (a) Illustration of weight degeneracy for SIS applied to the model in Equation (13.6). with parameters  $(\phi, q, \beta, r) = (0.9, 10.0, 0.5, 1.0)$ . We use  $T = 6$  steps and  $N_s = 5$  samples. We see that as  $t$  increases, almost all the probability mass concentrates on particle 3. Generated by `sis_vs_smcl.ipynb`. Adapted from Figure 2 of [NLS19]. (b) Illustration of the bootstrap particle filtering algorithm.

### 13.2.3 Sequential importance sampling with resampling

In this section, we describe **sequential importance sampling with resampling (SISR)**. The basic idea is this: instead of “growing” all of the old particle sequences by one step, we first select the  $N_s$  “fittest” particles, by sampling from the old posterior, and then we let these survivors grow by one step.

In more detail, at step  $t$ , we sample from

$$q_t^{\text{SISR}}(\mathbf{z}_{1:t}) = \hat{\pi}_{t-1}(\mathbf{z}_{1:t-1}) q_t(\mathbf{z}_t | \mathbf{z}_{1:t-1}) \quad (13.22)$$

where  $\hat{\pi}_{t-1}(\mathbf{z}_{1:t-1})$  is the previous weighted posterior approximation. By contrast, in SIS, we sample from

$$q_t^{\text{SIS}}(\mathbf{z}_{1:t}) = q_{t-1}^{\text{SIS}}(\mathbf{z}_{1:t-1}) q_t(\mathbf{z}_t | \mathbf{z}_{1:t-1}) \quad (13.23)$$

---

**Algorithm 32:** Sequential importance sampling with resampling (SISR)

---

```

1 Initialization:  $\mathbf{z}_1^i \sim q_1(\mathbf{z}_1)$ ,  $\tilde{w}_1^i = \frac{\tilde{\gamma}_1(\mathbf{z}_1^i)}{q_1(\mathbf{z}_1^i)}$ ,  $W_1^i = \frac{\tilde{w}_1^i}{\sum_j \tilde{w}_1^j}$ ,  $\hat{\pi}_1(\mathbf{z}_1) = \sum_{i=1}^{N_s} W_1^i \delta(\mathbf{z}_1 - \mathbf{z}_1^i)$ 
2 for  $t = 2 : T$  do
3   Compute ancestors  $\mathbf{a}_{t-1}^{1:N_s} = \text{resample}(\tilde{w}_{t-1}^{1:N_s})$ 
4   Select  $\mathbf{z}_{t-1}^{1:N_s} = \text{permute}(\mathbf{a}_{t-1}^{1:N_s}, \mathbf{z}_{t-1}^{1:N_s})$ 
5   Reset unnormalized weights  $\tilde{w}_{t-1}^{1:N_s} = 1/N_s$ 
6   for  $i = 1 : N_s$  do
7     Sample  $\mathbf{z}_t^i \sim q_t(\mathbf{z}_t | \mathbf{z}_{1:t-1}^i)$ 
8     Compute unnormalized weight  $\tilde{w}_t^i = \alpha_t^i = \frac{\tilde{\gamma}_t(\mathbf{z}_{1:t}^i)}{\tilde{\gamma}_{t-1}(\mathbf{z}_{1:t-1}^i) q_t(\mathbf{z}_t^i | \mathbf{z}_{1:t-1}^i)}$ 
9   Compute normalized weights  $W_t^i = \frac{\tilde{w}_t^i}{\sum_j \tilde{w}_t^j}$  for  $i = 1 : N_s$ 
10  Compute MC posterior  $\hat{\pi}_t(\mathbf{z}_{1:t}) = \sum_{i=1}^{N_s} W_t^i \delta(\mathbf{z}_{1:t} - \mathbf{z}_{1:t}^i)$ 

```

---

We can sample from Equation (13.22) in two steps. First we **resample**  $N_s$  samples from  $\hat{\pi}_{t-1}(\mathbf{z}_{1:t-1})$  to get a *uniformly weighted* set of new samples  $\mathbf{z}_{1:t-1}^i$ . (See Section 13.2.4 for details on how to do this.) Then we extend each sample using  $\mathbf{z}_t^i \sim q_t(\mathbf{z}_t | \mathbf{z}_{1:t-1}^i)$ , and concatenate  $\mathbf{z}_t^i$  to  $\mathbf{z}_{1:t-1}^i$ ,

After making a proposal, we compute the unnormalized weights. We use the standard SNIS method, except we “pretend” that the proposal is given by  $\tilde{\gamma}_{t-1}(\mathbf{z}_{1:t-1}^i) q_t(\mathbf{z}_t^i | \mathbf{z}_{1:t-1}^i)$  even though we used  $\hat{\pi}_{t-1}(\mathbf{z}_{1:t-1}^i) q_t(\mathbf{z}_t^i | \mathbf{z}_{1:t-1}^i)$ . The intuitive reason why this is valid is because the previous weighted approximation,  $\hat{\pi}_{t-1}(\mathbf{z}_{1:t-1}^i)$ , was an unbiased estimate of the previous target distribution,  $\tilde{\gamma}_{t-1}(\mathbf{z}_{1:t-1})$ . (See e.g., [CP20b] for more theoretical details.) We then compute the unnormalized weights, which are the same as the incremental weights, since the resampling step sets  $\tilde{w}_{t-1}^i = 1$ . We then normalize these weights and compute the new approximationg to the target posterior  $\hat{\pi}_t(\mathbf{z}_{1:t})$ . See Algorithm 32 for the pseudocode.

### 13.2.3.1 Bootstrap filter

We now consider a special case of SISR, in which the model is an SSM, and the proposal distribution is equal to the dynamical prior:

$$q_t(\mathbf{z}_t | \mathbf{z}_{1:t-1}) = p(\mathbf{z}_t | \mathbf{z}_{1:t-1}) \quad (13.24)$$

In this case, the corresponding incremental weight in Equation (13.21) simplifies to

$$\alpha_t(\mathbf{z}_{1:t}) = \frac{p(\mathbf{y}_t | \mathbf{z}_{1:t}) p(\mathbf{z}_t | \mathbf{z}_{1:t-1})}{q(\mathbf{z}_t | \mathbf{z}_{1:t-1})} = \frac{p(\mathbf{y}_t | \mathbf{z}_t) p(\mathbf{z}_t | \mathbf{z}_{t-1})}{p(\mathbf{z}_t | \mathbf{z}_{t-1})} = p(\mathbf{y}_t | \mathbf{z}_{1:t}) \quad (13.25)$$

This special case is called the **bootstrap filter** [Gor93] or the **survival of the fittest** algorithm [KKR95]. (In the computer vision literature, this is called the **condensation** algorithm, which stands for “conditional density propagation” [IB98].) See Figure 13.2b for an illustration of how this algorithm works, and Figure 13.1b for some sample results on real data.

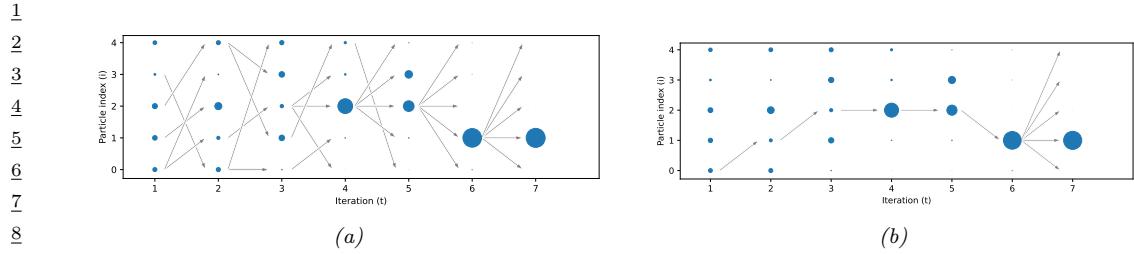


Figure 13.3: (a) Illustration of diversity of samples in SMC applied to the model in Equation (13.6). (b) Illustration of the path degeneracy problem. Generated by `sis_vs_smc.ipynb`. Adapted from Figure 3 of [NLS19].

The bootstrap filter is useful for models where we can sample from the dynamics, but cannot evaluate the transition model pointwise. This occurs in certain implicit dynamical models, such as those defined using differential equations (see e.g., [IBK06]); such models are often used in epidemiology. However, in general it is much more efficient to use proposals that take the current evidence  $\mathbf{y}_t$  into account. We discuss ways to approximate such “locally optimal” proposals in Section 13.3.

20

### 13.2.3.2 Path degeneracy problem

In Figure 13.3a we show how particle filtering can result in a much more diverse set of active particles, with more balanced weights when applied to the non-Markovian example in Equation (13.6).

While particle filtering does not suffer from weight degeneracy, it does suffer from another problem known as **path degeneracy**. This refers to the fact that the number of particles that “survive” (have non-negligible weight) over many steps may drop rapidly over time, resulting in a loss of diversity when we try to represent the distribution over the past. We illustrate this in Figure 13.3b, where we only include arrows for samples that have been resampled at each step up until the final step. We see that we have  $N_s = 5$  identical copies of  $\mathbf{z}_1^1$  in the final set of surviving sequences. (The time at which all the paths meet at a common ancestor, when tracing backwards in time, is known as the **coalescence** time.) We discuss some ways to ameliorate this issue in Section 13.2.4 and Section 13.2.5.

34

### 13.2.3.3 Estimating the normalizing constant

We can use particle filtering to approximate the normalization constant  $Z_T = p(\mathbf{y}_{1:T}) = \prod_{t=1}^T p(\mathbf{y}_t | \mathbf{y}_{1:t-1})$  as follows:

$$\hat{Z}_T = \prod_{t=1}^T \hat{Z}_t \quad (13.26)$$

where, from Equation (13.10), we have

$$\hat{Z}_t = \frac{1}{N_s} \sum_{i=1}^{N_s} \tilde{w}_t^i = \hat{Z}_{t-1} \left( \widehat{Z_t / Z_{t-1}} \right) \quad (13.27)$$

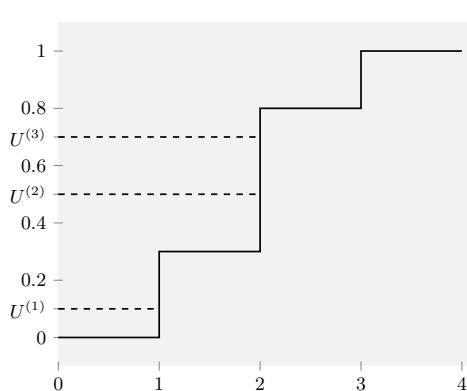


Figure 13.4: Illustration of how to sample from the empirical CDF  $P(x) = \sum_{n=1}^N W_n \mathbb{I}(x \geq n)$  shown in black. The height of step  $n$  is  $W_n$ . If  $U^m$  picks step  $n$ , then we set the ancestor of  $m$  to be  $n$ , i.e.,  $A^m = n$ . In this example,  $A^{1:3} = (1, 2, 2)$ . Adapted from Figure 9.3 of [CP20b].

where

$$\widehat{Z_t / Z_{t-1}} = \frac{\sum_{i=1}^{N_s} \tilde{w}_t^i}{\sum_{i=1}^{N_s} \tilde{w}_{t-1}^i} \quad (13.28)$$

This estimate of the marginal likelihood is very useful for tasks such as parameter estimation.

#### 13.2.4 Resampling methods

In this section, we discuss various resampling methods, which can help reduce weight degeneracy and path degeneracy.

##### 13.2.4.1 Multinomial resampling

The simplest approach to resampling is known as **multinomial resampling**. This works as follows. First we form the cumulative distribution from the weights  $W_{t-1}^{1:N_s}$ , as illustrated by the staircase in Figure 13.4. Then we sample  $N_s$  uniform random variables,  $u^i \sim \{0, 1\}$ . Finally, we see which bin (interval)  $u^i$  lands in; if it falls in bin  $a$ , we assign the new sample  $\mathbf{z}_{t-1}^i$  to be the same as the old  $\mathbf{z}_{t-1}^a$ . We say that  $a$  is the **ancestor** of  $i$ . For precisely, we say  $a$  is the ancestor of sample  $i$  if

$$\sum_{j=1}^{a-1} W_{t-1}^j \leq u^i < \sum_{j=1}^a W_{t-1}^j \quad (13.29)$$

See Listing 13.1 for some Python code.

Listing 13.1: Multinomial resampling

```
def multinomial_resampling(w):
    N = w.shape[0]
    u = np.random.rand(N)
```

```

1
2     bins = np.cumsum(w)
3     ancestors = np.digitize(u, bins)
4     return ancestors

```

Although this is a simple method, it can introduce a lot of variance into the representation of the distribution. For example, suppose all the weights are equal,  $W^n = 1/N$ . Let  $\mathcal{W}^n = \sum_{m=1}^N \mathbb{I}(A^m = n)$  be the number of “offspring” for particle  $n$  (i.e., the number of times this particle is chosen in the resampling step). We have  $\mathcal{W}^n \sim \text{Bin}(N, 1/N)$ , so  $P(\mathcal{W}^n = 0) = (1 - 1/N)^N \approx e^{-1} \approx 0.37$ . So there is a 37% chance that any given particle will disappear even though they all had the same initial weight. In the sections below, we discuss some **low variance resampling** methods.

#### 13.2.4.2 Stratified resampling

A simple approach to improve on multinomial resampling is to use **stratified resampling**, in which we divide the unit interval into  $N_s$  strata,  $(0, 1/N_s]$ ,  $(1/N_s, 2/N_s]$ , up to  $(1 - 1/N_s, 1]$ . We then generate  $u^i \sim \text{Unif}((i - 1)/N_s, i/N_s)$  and derive the corresponding ancestor indexes using Equation (13.29). See Listing 13.2 for some Python code.

*Listing 13.2: Stratified resampling*

```

18 def stratified_resampling(w):
19     N = w.shape[0]
20     u = (np.arange(N) + np.random.rand(N))/N
21     bins = np.cumsum(w)
22     ancestors = np.digitize(u, bins)
23     return ancestors

```

#### 13.2.4.3 Systematic resampling

We can further reduce the variance by forcing all the samples from  $\hat{\pi}_{t-1}$  to be deterministically generated from a shared random source,  $u \sim \text{Unif}(0, 1)$ , by computing

$$u^i = \frac{i - 1}{N_s} + \frac{u}{N_s} \quad (13.30)$$

We derive the corresponding ancestor indexes using Equation (13.29). See Listing 13.3 for some Python code. (The only difference from Listing 13.2 is the use of `np.random.rand()` instead of `np.random.rand(N)`.)

*Listing 13.3: Systematic resampling*

```

35 def systematic_resampling(w):
36     N = w.shape[0]
37     u = (np.arange(N) + np.random.rand())/N
38     bins = np.cumsum(w)
39     ancestors = np.digitize(u, bins)
40     return ancestors

```

#### 13.2.4.4 Comparison

It can be proved that all of the above methods are unbiased. It can also be proved that stratified resampling is lower variance than multinomial resampling. Empirically it seems that systematic resampling is lower variance than other methods [HSG06]. A more complex resampling scheme, that is guaranteed to converge and which is also low variance, is described in [GCW19].

---

### 13.2.5 Adaptive resampling

The resampling step can result in loss of diversity, since each ancestor may generate multiple children, and some may generate no children, since the ancestor indices  $A_t^n$  are sampled independently; this is the path degeneracy problem mentioned above. On the other hand, if we never resample, we end up with SIS, which suffers from weight degeneracy (particles with negligible weight). A compromise is to use **adaptive resampling**, in which we resample whenever the **effective sample size** or ESS drops below some minimum, such as  $N/2$ . A common way to define the ESS is as follows:<sup>1</sup>

$$\text{ESS}(W^{1:N}) = \frac{1}{\sum_{n=1}^N (W^n)^2} \quad (13.31)$$

Alternatively we can compute the ESS using the unnormalized weights:

$$\text{ESS}(\tilde{w}^{1:N}) = \frac{\left(\sum_{n=1}^N \tilde{w}^n\right)^2}{\sum_{n=1}^N (\tilde{w}^n)^2} \quad (13.32)$$

Note that if we have  $k$  weights with  $\tilde{w}^n = 1$  and  $N - k$  weights with  $\tilde{w}^n = 0$ , then the ESS is  $k$ ; thus ESS is between 1 and  $N$ .

The pseudocode for SISR with adaptive resampling is given in Algorithm 33. (We use the notation of [Law+22, App. B], in which we first sample new extensions of the sequences, and then optionally resample the sequences at the end of each step.)

---

**Algorithm 33:** SISR with adaptive resampling (generic SMC)

---

```

1 Initialization:  $\tilde{w}_0^{1:N_s} = 1$ ,  $\hat{Z}_0 = 1$ 
2 for  $t = 1 : T$  do
3   for  $i = 1 : N_s$  do
4     Sample particle  $\mathbf{z}_t^i \sim q_t(\mathbf{z}_t | \mathbf{z}_{1:t-1}^i)$ 
5     Compute incremental weight  $\alpha_t^i = \frac{\tilde{\gamma}_t(\mathbf{z}_{1:t}^i)}{\tilde{\gamma}_{t-1}(\mathbf{z}_{1:t-1}^i)q_t(\mathbf{z}_t^i | \mathbf{z}_{1:t-1}^i)}$ 
6     Compute unnormalized weight  $\tilde{w}_t^i = \tilde{w}_{t-1}^i \alpha_t^i$ 
7   Estimate normalization constant:  $\widehat{Z_t / Z_{t-1}} = \frac{\sum_{i=1}^{N_s} \tilde{w}_t^i}{\sum_{i=1}^{N_s} \tilde{w}_{t-1}^i}$ ,  $\hat{Z}_t = \hat{Z}_{t-1} (\widehat{Z_t / Z_{t-1}})$ 
8   if  $\text{ESS}(\tilde{w}_{t-1}^{1:N}) < \text{ESS}_{\min}$  then
9     Compute ancestors  $\mathbf{a}_t^{1:N_s} = \text{resample}(\tilde{w}_t^{1:N_s})$ 
10    Select  $\mathbf{z}_t^{1:N_s} = \text{permute}(\mathbf{a}_t^{1:N_s}, \mathbf{z}_t^{1:N_s})$ 
11    Reset unnormalized weights  $\tilde{w}_t^{1:N_s} = 1/N_s$ 
12    Compute normalized weights  $W_t^i = \frac{\tilde{w}_t^i}{\sum_j \tilde{w}_t^j}$  for  $i = 1 : N_s$ 
13    Compute MC posterior  $\hat{\pi}_t(\mathbf{z}_{1:t}) = \sum_{i=1}^{N_s} W_t^i \delta(\mathbf{z}_{1:t} - \mathbf{z}_{1:t}^i)$ 

```

---

<sup>1</sup> Note that the ESS used in SMC is different than the ESS used in MCMC (Section 12.6.3); the latter takes into account auto-correlation of the MCMC samples.

1 **13.3 Proposal distributions**

3 The efficiency of PF is crucially dependent on the quality of the proposal distribution. We discuss  
4 some options below.

6 **13.3.1 Locally optimal proposal**

8 We define the (one-step) **locally optimal proposal distribution**  $q_t^*(\mathbf{z}_t | \mathbf{z}_{1:t-1})$  to be the one that  
9 minimizes

$$\underline{11} \quad D_{\text{KL}}(\pi_{t-1}(\mathbf{z}_{1:t-1})q_t(\mathbf{z}_t | \mathbf{z}_{1:t-1}) \| \pi_t(\mathbf{z}_{1:t})) \quad (13.33)$$

$$\underline{12} \quad = \mathbb{E}_{\pi_{t-1}q_t} [\log \{\pi_{t-1}(\mathbf{z}_{1:t-1})q_t(\mathbf{z}_t | \mathbf{z}_{1:t-1})\} - \log \pi_t(\mathbf{z}_{1:t})] \quad (13.34)$$

$$\underline{14} \quad = \mathbb{E}_{\pi_{t-1}q_t} [\log q_t(\mathbf{z}_t | \mathbf{z}_{1:t-1}) - \log \pi_t(\mathbf{z}_t | \mathbf{z}_{1:t-1})] + \text{const} \quad (13.35)$$

$$\underline{15} \quad = \mathbb{E}_{\pi_{t-1}q_t} [D_{\text{KL}}(q_t(\mathbf{z}_t | \mathbf{z}_{1:t-1}) \| \pi_t(\mathbf{z}_t | \mathbf{z}_{1:t-1}))] + \text{const} \quad (13.36)$$

16 The KL is minimized by choosing

$$\underline{19} \quad q_t^*(\mathbf{z}_t | \mathbf{z}_{1:t-1}) = \pi_t(\mathbf{z}_t | \mathbf{z}_{1:t-1}) = \frac{\tilde{\gamma}_t(\mathbf{z}_{1:t})}{\tilde{\gamma}_t(\mathbf{z}_{1:t-1})} \quad (13.37)$$

21 where  $\tilde{\gamma}_t(\mathbf{z}_{1:t-1}) = \int \tilde{\gamma}_t(\mathbf{z}_{1:t}) d\mathbf{z}_t$  is the probability of the past sequence under the current target  
22 distribution.

23 Note that the subscript  $t$  specifies the  $t$ 'th distribution, so in the context of SSMs, we have  
24  $\pi_t(\mathbf{z}_t | \mathbf{z}_{1:t-1}) = p(\mathbf{z}_t | \mathbf{z}_{1:t-1}, \mathbf{y}_{1:t})$ . Thus we see that when proposing  $\mathbf{z}_t$ , we should condition on all  
25 the data, including the most recent observation,  $\mathbf{y}_t$ ; this is called a **guided particle filter**, and will  
26 will be better than the bootstrap filter, which proposes from the prior.

27 In general, it is intractable to compute the locally optimal proposal, so we consider various  
28 approximations below.

30 **13.3.2 Proposals based on the extended and unscented Kalman filter**

32 One way to approximate the locally optimal proposal distribution is based on the extended Kalman  
33 filter (Section 8.5.2) or the unscented Kalman filter (Section 13.3.2, which gives rise to the **extended**  
34 **particle filter** [DGA00] and **unscented particle filter** [Mer+00] respectively. To explain these  
35 methods, we follow the presentation of [NLS19, p36]. As usual, we assume the dynamical system  
36 can be written as  $\mathbf{z}_t = \mathbf{f}(\mathbf{z}_{t-1}) + \mathbf{q}_t$  and  $\mathbf{y}_t = \mathbf{h}(\mathbf{z}_t) + \mathbf{r}_t$ , where  $\mathbf{q}_t$  is the system noise and  $\mathbf{r}_t$  is  
37 the observation noise. The EKF and UKF approximations assume that the joint distribution over  
38 neighboring time steps, given the  $i$ 'th history, is Gaussian:

$$\underline{40} \quad p(\mathbf{z}_t, \mathbf{y}_t | \mathbf{z}_{1:t-1}^i) \approx \mathcal{N} \left( \begin{pmatrix} \mathbf{z}_t \\ \mathbf{y}_t \end{pmatrix} | \hat{\mu}^i, \hat{\Sigma}^i \right) \quad (13.38)$$

43 where

$$\underline{44} \quad \hat{\mu}^i = \begin{pmatrix} \hat{\mu}_z^i \\ \hat{\mu}_y^i \end{pmatrix}, \hat{\Sigma}^i = \begin{pmatrix} \hat{\Sigma}_{zz}^i & \hat{\Sigma}_{zy}^i \\ \hat{\Sigma}_{yz}^i & \hat{\Sigma}_{yy}^i \end{pmatrix} \quad (13.39)$$

1  
2 (See Section 8.8 for details.)

3 The EKF and UKF compute  $\hat{\mu}^i$  and  $\hat{\Sigma}^i$  differently. In the EKF, we linearize  $f$  and  $h$ , and assume  
4 the noise terms are Gaussian. We then compute  $p(z_t, y_t | z_{1:t-1}^i)$  exactly for this linearized model  
5 (see Section 8.5.1). In the UKF, we propagate sigma points through  $f$  and  $h$ , and approximate the  
6 resulting means and covariances using the unscented transform, which can be more accurate (see  
7 Section 8.7). Once we have computed  $\hat{\mu}^i$  and  $\hat{\Sigma}^i$ , we can use standard rules for Gaussian conditioning  
8 to compute the approximate proposal as follows:

$$\underline{9} \quad q(z_t | z_{1:t-1}^i, y_t) \approx \mathcal{N}(z_t | \mu_t^i, \Sigma_t^i) \quad (13.40)$$

$$\underline{10} \quad \mu_t^i = \hat{\mu}_z^i + \hat{\Sigma}_{zy}^i (\hat{\Sigma}_{yy}^i)^{-1} (y_t - \hat{\mu}_y^i) \quad (13.41)$$

$$\underline{11} \quad \Sigma_t^i = \hat{\Sigma}_{zz}^i - \hat{\Sigma}_{zy}^i (\hat{\Sigma}_{yy}^i)^{-1} \hat{\Sigma}_{yz}^i \quad (13.42)$$

12 Note that the linearization (or sigma point) approximation needs to be performed for each particle  
13 separately.

### 14 13.3.3 Proposals based on the Laplace approximation

15 To handle non-Gaussian likelihoods in an SSM, we can use the Laplace approximation (Section 7.4.3),  
16 as suggested in [DGA00]. In particular, consider an SSM with linear-Gaussian latent dynamics and a  
17 GLM likelihood. At each step, we compute the maximum  $z_t^* = \text{argmax} \log p(y_t | z_t)$  as step  $t$  (e.g.,  
18 using Newton-Raphson), and then approximate the likelihood using

$$\underline{19} \quad p(y_t | z_t) \approx \mathcal{N}(z_t | z_t^*, -\mathbf{H}_t^*) \quad (13.43)$$

20 where  $\mathbf{H}_t^*$  is the Hessian of the log-likelihood at the mode. We now compute  $p(z_t | z_{t-1}^i, y_t)$  using the  
21 update step of the Kalman filter, using the same equations as in Section 13.3.2. This combination is  
22 called the the **Laplace Gaussian filter** [Koy+10]. We give an example in Section 13.3.3.1.

#### 23 13.3.3.1 Example: neural decoding

24 In this section, we give an example where we apply the Laplace approximation to an SSM with  
25 linear-Gaussian dynamics and a Poisson likelihood. The application arises from neuroscience. In  
26 particular, assume we record the **neural spike trains** as a monkey moves its hand around in space.  
27 Let  $z_t \in \mathbb{R}^6$  represent the 3d location and velocity of the hand. We model the dynamics of the hand  
28 using a simple Brownian random walk model [CP20b, p157]:

$$\underline{29} \quad \begin{pmatrix} z_t(i) \\ z_t(i+3) \end{pmatrix} | z_{t-1} \sim \mathcal{N}_2 \left( \begin{pmatrix} 1 & \Delta \\ 0 & 1 \end{pmatrix} \begin{pmatrix} z_{t-1}(i) \\ z_{t-1}(i+3) \end{pmatrix}, \sigma^2 \mathbf{Q} \right), i = 1 : 3 \quad (13.44)$$

30 where the covariance of the noise is given by the following, assuming a discretization step of  $\Delta$ :

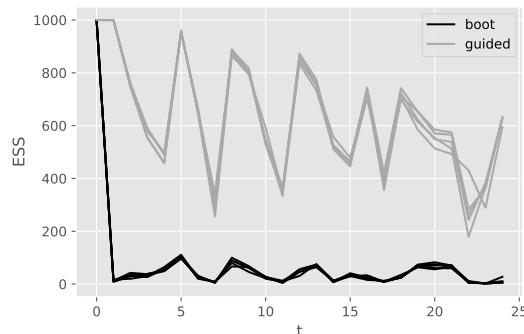
$$\underline{31} \quad \mathbf{Q} = \begin{pmatrix} \Delta^3/3 & \Delta^2/2 \\ \Delta^2/2 & \Delta \end{pmatrix} \quad (13.45)$$

32 We assume the  $k$ 'th observation at time  $t$  is the number of spikes for neuron  $k$  in this sensing  
33 interval:

$$\underline{34} \quad p(y_t(k) | z_t) = \text{Poi}(\lambda_k(z_t)) \quad (13.46)$$

$$\underline{35} \quad \log \lambda_k(z_t) = \alpha_k + \beta_k^\top z_t \quad (13.47)$$

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13



14 *Figure 13.5: Effective sample size at each step for the bootstrap particle filter and a guided particle fil-  
15 ter for a Gaussian SSM with Poisson likelihood. Adapted from Figure 10.4 of [CP20b]. Generated by  
16 [pf\\_guided\\_neural\\_decoding.ipynb](#).*

17

18

19

20 Our goal is to compute  $p(\mathbf{z}_t | \mathbf{y}_{1:t})$ , which lets us infer the position of the hand from the neural code.  
21 (Apart from its value for furthering basic science, this can be useful for applications such as helping  
22 disabled people control their arms using “mind control”.)

23 To illustrate this, we sample a synthetic dataset from the model, to simulate a “monkey” moving  
24 its arm for  $T = 25$  time steps; this generates  $K = 50$  neuronal counts per time step. We then apply  
25 particle filtering to this dataset (using the true model), using either the bootstrap filter (i.e., proposal  
26 is the random walk prior) or the guided filter (i.e., proposal is the Laplace approximation mentioned  
27 above). In Figure 13.5, we see that the effective sample size of the guided filter is much higher than  
28 for the bootstrap filter.

29

### 30 13.3.4 Proposals based on SMC (nested SMC)

31 It is possible to use SMC as a subroutine to compute a proposal distribution for SMC: at each step  $t$ ,  
32 for each particle  $i$ , we run an SMC algorithm where the target distribution is the optimal proposal,  
33  $p(\mathbf{z}_t | \mathbf{z}_{1:t-1}^i, \mathbf{y}_{1:t})$ . This is called **nested SMC** [NLS15; NLS19].

34 This method can approximate the locally optimal proposal arbitrarily well, since it does not make  
35 any limiting parametric assumptions. However, the method can be slow, although the inner SMC  
36 algorithm can be run in parallel for each outer sample [NLS15; NLS19].

37

## 38 13.4 Rao-Blackwellised particle filtering (RBPF)

39

40 In some models, we can partition the hidden variables into two kinds,  $\mathbf{m}_t$  and  $\mathbf{z}_t$ , such that we can  
41 analytically integrate out  $\mathbf{z}_t$  provided we know the values of  $\mathbf{m}_{1:t}$ . This means we only have to sample  
42  $\mathbf{m}_{1:t}$ , and can represent  $p(\mathbf{z}_t | \mathbf{m}_{1:t}, \mathbf{y}_{1:t})$  parametrically. These hybrid particles are sometimes called  
43 **distributional particles** or **collapsed particles** [KF09a, Sec 12.4]. This combines techniques from  
44 particle filtering (Section 13.2) with deterministic methods such as Kalman filtering (Section 8.3.2).

45 The advantage of this approach is that we reduce the dimensionality of the space in which we are  
46 47

sampling, which reduces the variance of our estimate. This technique is known as **Rao-Blackwellised particle filtering** or **RBPF** for short. (See Section 11.6.2 for more details on Rao-Blackwellisation.) In Section 13.4.1 we give an example of RBPF for inference in a switching linear dynamical systems. In Section 13.4.3 we illustrate RBPF for inference in the SLAM model for a mobile robot.

### 13.4.1 Mixture of Kalman filters

In this section, we consider the application of RBPF to a switching linear dynamical system (Section 29.9). This model has both continuous and discrete latent variables. This can be used to track a system that switches between discrete modes or operating regimes, represented by the discrete variable  $m_t$ .

For notational simplicity, we ignore the control inputs  $\mathbf{u}_t$ . Thus the model is given by

$$p(\mathbf{z}_t | \mathbf{z}_{t-1}, m_t = k) = \mathcal{N}(\mathbf{z}_t | \mathbf{F}_k \mathbf{z}_{t-1}, \mathbf{Q}_k) \quad (13.48)$$

$$p(\mathbf{y}_t | \mathbf{z}_t, m_t = k) = \mathcal{N}(\mathbf{y}_t | \mathbf{H}_k \mathbf{z}_t, \mathbf{R}_k) \quad (13.49)$$

$$p(m_t = k | m_{t-1} = j) = A_{jk} \quad (13.50)$$

We let  $\boldsymbol{\theta}_k = (\mathbf{F}_k, \mathbf{H}_k, \mathbf{Q}_k, \mathbf{R}_k, \mathbf{A}_{:,k})$  represent all the parameters for state  $k$ .

Exact inference is intractable, but if we sample the discrete variables, we can infer the continuous variables conditioned on the discretes exactly, making this a good candidate for RBPF. In particular, if we sample trajectories  $\mathbf{m}_{1:t}^n$ , we can apply a Kalman filter to each particle. This can be thought of as a **mixture of Kalman filters** [CL00]. The resulting belief state is represented by

$$p(\mathbf{z}_t, \mathbf{m}_t | \mathbf{y}_{1:t}) \approx \sum_{n=1}^N W_t^n \delta(\mathbf{m}_t - \mathbf{m}_t^n) \mathcal{N}(\mathbf{z}_t | \boldsymbol{\mu}_t^n, \boldsymbol{\Sigma}_t^n) \quad (13.51)$$

To derive the filtering algorithm, note that the full posterior at time  $t$  can be written as follows:

$$p(\mathbf{m}_{1:t}, \mathbf{z}_{1:t} | \mathbf{y}_{1:t}) = p(\mathbf{z}_{1:t} | \mathbf{m}_{1:t}, \mathbf{y}_{1:t}) p(\mathbf{m}_{1:t} | \mathbf{y}_{1:t}) \quad (13.52)$$

The second term is given by the following:

$$p(\mathbf{m}_{1:t} | \mathbf{y}_{1:t}) \propto p(\mathbf{y}_t | \mathbf{m}_{1:t}, \mathbf{y}_{1:t-1}) p(\mathbf{m}_{1:t} | \mathbf{y}_{1:t-1}) \quad (13.53)$$

$$= p(\mathbf{y}_t | \mathbf{m}_{1:t}, \mathbf{y}_{1:t-1}) p(\mathbf{m}_t | \mathbf{m}_{1:t-1}, \mathbf{y}_{1:t-1}) p(\mathbf{m}_{1:t-1} | \mathbf{y}_{1:t-1}) \quad (13.54)$$

$$= p(\mathbf{y}_t | \mathbf{m}_{1:t}, \mathbf{y}_{1:t-1}) p(\mathbf{m}_t | \mathbf{m}_{t-1}) p(\mathbf{m}_{1:t-1} | \mathbf{y}_{1:t-1}) \quad (13.55)$$

Note that, unlike the case of standard particle filtering, we cannot write  $p(\mathbf{y}_t | \mathbf{m}_{1:t}, \mathbf{y}_{1:t-1}) = p(\mathbf{y}_t | \mathbf{m}_t)$ , since  $\mathbf{m}_t$  does not d-separate the past observations from  $\mathbf{y}_t$ , as is evident from Figure 29.25a.

Suppose we use the following recursive proposal distribution:

$$q(\mathbf{m}_{1:t} | \mathbf{y}_{1:t}) = q(\mathbf{m}_t | \mathbf{m}_{1:t-1}, \mathbf{y}_{1:t}) q(\mathbf{m}_{1:t-1} | \mathbf{y}_{1:t}) \quad (13.56)$$

Then we get the unnormalized importance weights

$$\tilde{w}_t^n \propto \frac{p(\mathbf{y}_t | \mathbf{m}_t^n, \mathbf{m}_{1:t-1}^n, \mathbf{y}_{1:t-1}) p(\mathbf{m}_t^n | \mathbf{m}_{t-1}^n)}{q(\mathbf{m}_t^n | \mathbf{m}_{1:t-1}^n, \mathbf{y}_{1:t})} \tilde{w}_{t-1}^n \quad (13.57)$$

As a special case, suppose we propose from the prior,  $q(m_t | \mathbf{m}_{t-1}^n, \mathbf{y}_{1:t}) = p(m_t | m_{t-1}^n)$ . If we sample discrete state  $k$ , the weight update becomes

$$\tilde{w}_t^n \propto \tilde{w}_{t-1}^n p(\mathbf{y}_t | m_t^n = k, \mathbf{m}_{1:t-1}^n, \mathbf{y}_{1:t-1}) = \tilde{w}_{t-1}^n L_{tk}^n \quad (13.58)$$

where

$$L_{tk}^n = p(\mathbf{y}_t | m_t = k, \mathbf{m}_{1:t-1}^n, \mathbf{y}_{1:t-1}) = \int p(\mathbf{y}_t | m_t = k, \mathbf{z}_t) p(\mathbf{z}_t | m_t = k, \mathbf{y}_{1:t-1}, \mathbf{m}_{1:t-1}^n) d\mathbf{z}_t \quad (13.59)$$

The quantity  $L_{tk}^n$  is the predictive density for the new observation  $\mathbf{y}_t$  conditioned on  $m_t = k$  and the history of previous latents,  $\mathbf{m}_{1:t-1}^n$ . In the case of SLDS models, this can be computed using the normalization constant of the Kalman filter, Equation (8.80). The resulting algorithm is shown in Algorithm 34. The step marked ‘‘KFupdate’’ refers to the Kalman filter update equations in Section 8.3.2, and is applied to each particle separately.

---

**Algorithm 34:** One step of RBPF for SLDS using prior as proposal

---

```

17 1 for  $n = 1 : N$  do
18 2    $k \sim p(m_t | m_{t-1}^n)$ 
19 3    $m_t^n := k$ 
20 4    $(\boldsymbol{\mu}_t^n, \boldsymbol{\Sigma}_t^n, L_{tk}^n) = \text{KFupdate}(\boldsymbol{\mu}_{t-1}^n, \boldsymbol{\Sigma}_{t-1}^n, \mathbf{y}_t, \boldsymbol{\theta}_k)$ 
21 5    $\tilde{w}_t^n = \tilde{w}_{t-1}^n L_{tk}^n$ 
22 6   Compute ESS = ESS( $\tilde{w}_t^{1:N_s}$ )
23 7   if  $ESS < ESS_{\min}$  then
24 8      $\mathbf{a}_t^{1:N} = \text{Resample}(\tilde{w}_t^{1:N})$ 
25 9      $(\mathbf{m}_t^{1:N_s}, \boldsymbol{\mu}_t^{1:N_s}, \boldsymbol{\Sigma}_t^{1:N_s}) = \text{permute}(\mathbf{a}_t, \mathbf{m}_t^{1:N_s}, \boldsymbol{\mu}_t^{1:N_s}, \boldsymbol{\Sigma}_t^{1:N_s})$ 
26 10     $\tilde{w}_t^n = 1/N_s$ 

```

---

### 13.4.1.1 Improvements

An improved version of the algorithm can be developed based on the fact that we are sampling a discrete state space. At each step, we propagate each of the  $N$  old particles through all  $K$  possible transition models. We then compute the weight for all  $NK$  new particles, and sample from this to get the final set of  $N$  particles. This latter step can be done using the **optimal resampling** method of [FC03], which will stochastically select the particles with the largest weight, while also ensuring the result is an unbiased approximation. In addition, this approach ensures that we do not have duplicate particles, which is wasteful and unnecessary when the state space is discrete.

### 13.4.2 Example: tracking a maneuvering object

In this section we give an example of RBPF for an SLDS from [DGK01]. Our goal is to track an object that has the following motion model:

$$p(\mathbf{z}_t | \mathbf{z}_{t-1}, m_t = k) = \mathcal{N}(\mathbf{z}_t | \mathbf{F}\mathbf{z}_{t-1} + \mathbf{b}_k, \mathbf{Q}) \quad (13.60)$$

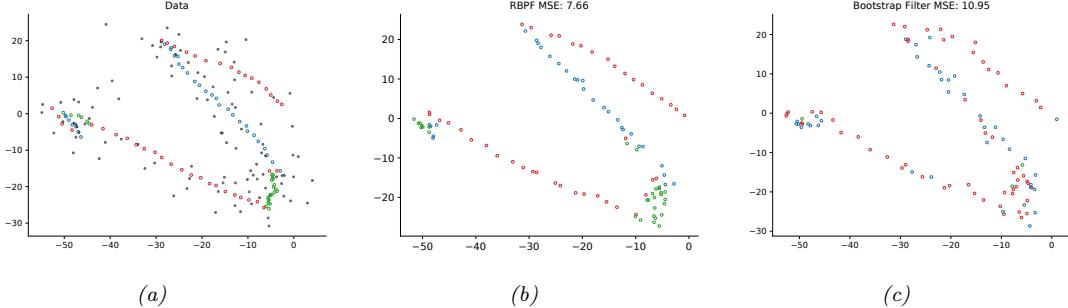


Figure 13.6: Illustration of state estimation for a switching linear model. (a) Black dots are observations, hollow circles are the true location, colors represent the discrete state. (b) Estimate from RBPF. Generated by [rbpf\\_maneuver.ipynb](#). (c) Estimate from bootstrap filter. Generated by [bootstrap\\_filter\\_maneuver.ipynb](#).

where  $\mathbf{z}_t = (x_{1t}, \dot{x}_{1t}, x_{2t}, \dot{x}_{2t})$  contains the 2d position and velocity. We define the observation matrix by  $\mathbf{H} = \mathbf{I}$  and the observation covariance by  $\mathbf{R} = 10\text{diag}(2, 1, 2, 1)$ . We define the dynamics matrix by

$$\mathbf{F} = \begin{pmatrix} 1 & \Delta & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & \Delta \\ 0 & 0 & 0 & 1 \end{pmatrix} \quad (13.61)$$

where  $\Delta = 0.1$ . We set the noise covariance to  $\mathbf{Q} = 0.2\mathbf{I}$  and the input bias vectors for each state to  $\mathbf{b}_1 = (0, 0, 0, 0)$ ,  $\mathbf{b}_2 = (-1.225, -0.35, 1.225, 0.35)$  and  $\mathbf{b}_3 = (1.225, 0.35, -1.225, -0.35)$ . Thus the system will turn in different directions depending on the discrete state. The discrete state transition matrix is given by

$$\mathbf{A} = \begin{pmatrix} 0.8 & 0.1 & 0.1 \\ 0.1 & 0.8 & 0.1 \\ 0.1 & 0.1 & 0.8 \end{pmatrix} \quad (13.62)$$

Figure 13.6a shows some observations, and the true state of the system, from a sample run, for 100 steps. The colors denote the discrete state, and the location of the symbol denotes the  $(x, y)$  location. The small dots represent noisy observations. Figure 13.6b shows the estimate of the state computed using RBPF with the optimal proposal with 1000 particles. In Figure 13.6c, we show the analogous estimate using the bootstrap filter, which does much worse.

In Figure 13.7a and Figure 13.7b, we show the posterior marginals of the  $(x, y)$  locations over time. In Figure 13.7c we show the true discrete state, and in Figure 13.7d we show the posterior marginal over discrete states. The overall state classification error rate is 29%, but it seems that occasionally misclassifying isolated time steps does not significantly hurt estimation of the continuous states, as we can see from Figure 13.6b.

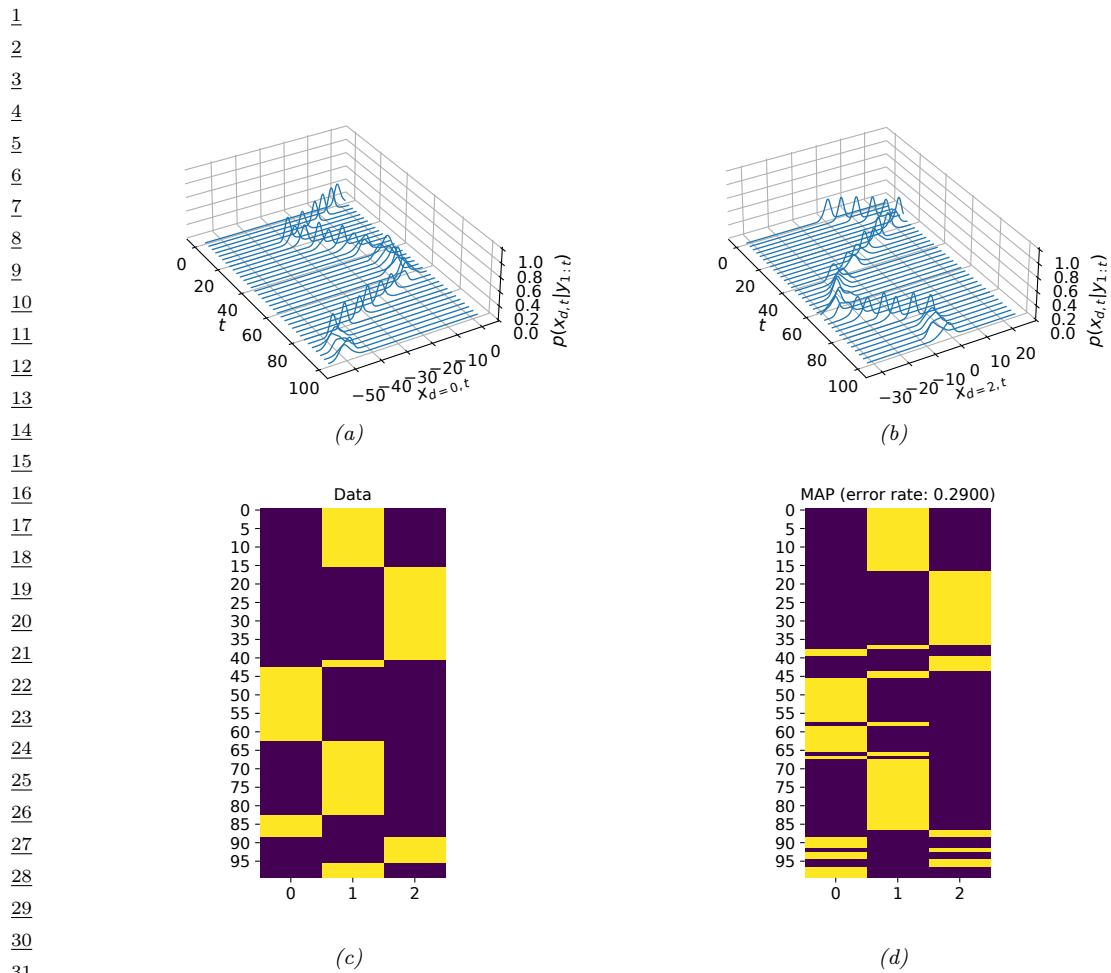


Figure 13.7: Visualizing the posterior from the RBPF algorithm. Top row: Posterior marginals of the location of the object over time, derived from the mixture of Gaussian representation for (a)  $x$  location (dimension 0), (b)  $y$  location (dimension 2). Bottom row: visualization of the true (c) and predicted (d) discrete states. Generated by [rbpf\\_maneuver.ipynb](#).

### 13.4.3 Example: FastSLAM

Consider a robot moving around an environment, such as a maze or indoor office environment. It needs to learn a map of the environment, and keep track of its location (pose) within that map. This problem is known as **simultaneous localization and mapping**, or **SLAM** for short. SLAM is widely used in mobile robotics (see e.g., [SC86; CN01; TBF06] for details). It is also useful in augmented reality, where the task is to recursively estimate the 3d pose of a handheld camera with respect to a set of 2d visual landmarks (this is known as **visual SLAM**, [TUI17; SMT18; Cza+20; DH22]).

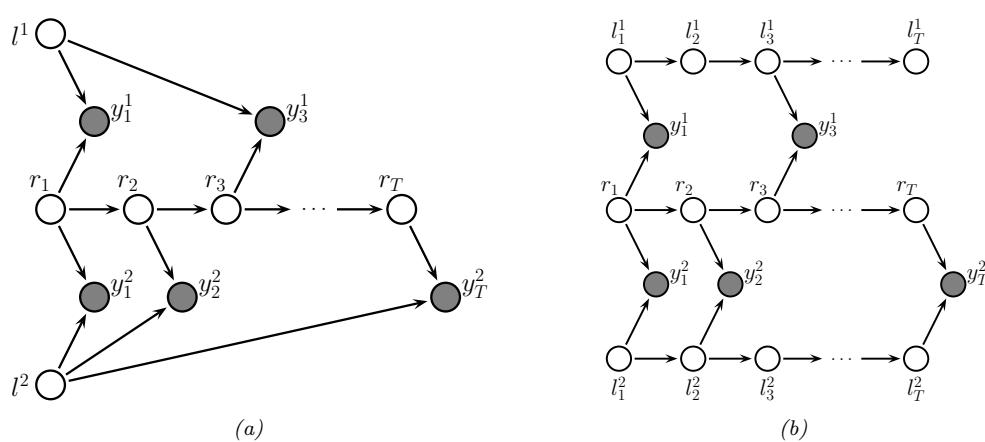


Figure 13.8: Graphical model representing the SLAM problem.  $\mathbf{l}_t^k$  is the location of landmark  $k$  at time  $t$ ,  $\mathbf{r}_t$  is the location of the robot at time  $t$ , and  $\mathbf{y}_t$  is the observation vector. In the model on the left, the landmarks are static (so they act like global shared parameters), on the right, their location can change over time. The robot's observations are based on the distance to the nearest landmarks from the current state, denoted  $f(\mathbf{r}_t, \mathbf{l}_t^k)$ . The number of observations per time step is variable, depending on how many landmarks are within the range of the sensor. Adapted from Figure 15.A.3 of [KF09a].

Let us assume we can represent the map as the 2d locations of a set of  $K$  landmarks, denote them by  $\mathbf{l}^1, \dots, \mathbf{l}^K$  (each is a vector in  $\mathbb{R}^2$ ). (We can use data association to figure out which landmark generated each observation, as discussed in Section 29.9.3.2.) Let  $\mathbf{r}_t$  represent the unknown location of the robot at time  $t$ . Let  $\mathbf{z}_t = (\mathbf{r}_t, \mathbf{l}_t^{1:K})$  be the combined state space. We can then perform online inference so that the robot can update its estimate of its own location, and the landmark locations.

The state transition model is defined as

$$p(\mathbf{z}_t | \mathbf{z}_{t-1}, \mathbf{u}_t) = p(\mathbf{r}_t | \mathbf{r}_{t-1}, \mathbf{l}_{t-1}^{1:K}, \mathbf{u}_t) \prod_{k=1}^K p(\mathbf{l}_t^k | \mathbf{l}_{t-1}^k) \quad (13.63)$$

where  $p(\mathbf{r}_t | \mathbf{r}_{t-1}, \mathbf{l}_{t-1}^{1:K}, \mathbf{u}_t)$  specifies how the robot moves given the control signal  $\mathbf{u}_t$  and the location of the obstacles  $\mathbf{l}_{t-1}^{1:K}$ . (Note that in this section, we assume that a human is joysticking the robot through the environment, so  $\mathbf{u}_{1:t}$  is given as input, i.e., we do not address the decision-theoretic issue of choosing where to move.)

If the obstacles (landmarks) are static, we can define  $p(\mathbf{l}_t^k | \mathbf{l}_{t-1}^k) = \delta(\mathbf{l}_t^k - \mathbf{l}_{t-1}^k)$ , which is equivalent to treating the map as an unknown parameter that is shared globally across all time steps. More generally, we can let the landmark locations evolve over time [Mur00].

The observations  $\mathbf{y}_t$  measure the distance from  $\mathbf{r}_t$  to the set of closest landmarks. Figure 13.8 shows the corresponding graphical model for the case where  $K = 2$ , and where on the first step it sees landmarks 1 and 2, then just landmark 2, then just landmark 1, etc.

If all the CPDs are linear-Gaussian, then we can use a Kalman filter to maintain our belief state about the location of the robot and the location of the landmarks,  $p(\mathbf{z}_t | \mathbf{y}_{1:t}, \mathbf{u}_{1:t})$ . In the more general case of a nonlinear model, we can use the EKF (Section 8.5.2) or UKF (Section 8.7.2).

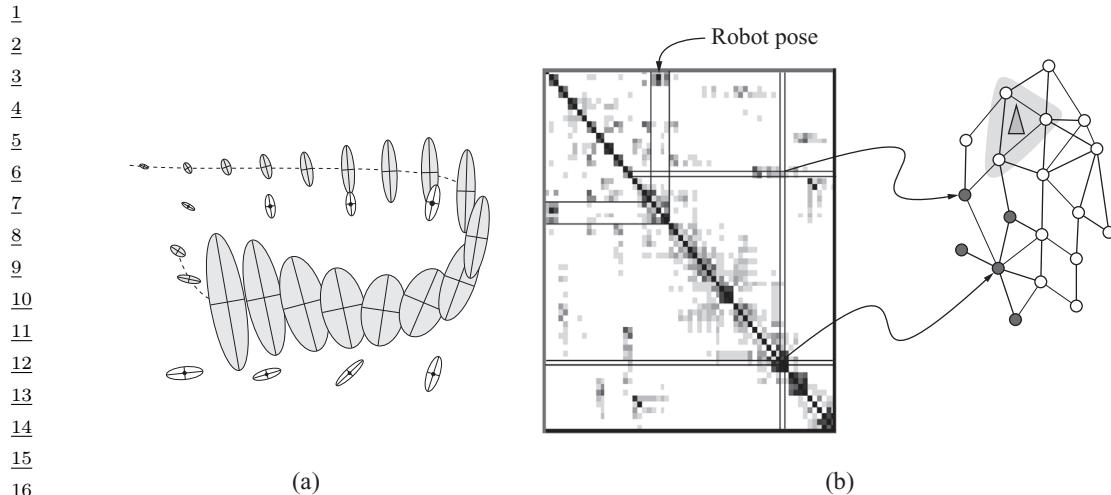


Figure 13.9: Illustration of the SLAM problem. (a) A robot starts at the top left and moves clockwise in a circle back to where it started. We see how the posterior uncertainty about the robot's location increases and then decreases as it returns to a familiar location, closing the loop. If we performed smoothing, this new information would propagate backwards in time to disambiguate the entire trajectory. (b) We show the precision matrix, representing sparse correlations between the landmarks, and between the landmarks and the robot's position (pose). The conditional independencies encoded by the sparse precision matrix can be visualized as a Gaussian graphical model, as shown on the right. From Figure 15.A.3 of [KF09a]. Used with kind permission of Daphne Koller.

Over time, the uncertainty in the robot's location will increase, due to wheel slippage etc., but when the robot returns to a familiar location, its uncertainty will decrease again. This is called **closing the loop**, and is illustrated in Figure 13.9(a), where we see the uncertainty ellipses, representing  $\text{Cov} [\mathbf{z}_t | \mathbf{y}_{1:t}, \mathbf{u}_{1:t}]$ , grow and then shrink.

In addition to visualizing the uncertainty of the robot's location, we can visualize the uncertainty about the map. To do this, consider the posterior precision matrix,  $\mathbf{\Lambda}_t = \mathbf{\Sigma}_t^{-1}$ . Zeros in the precision matrix correspond to absent edges in the corresponding undirected Gaussian graphical model (GGM, see Section 4.3.5). Initially all the beliefs about landmark locations are uncorrelated (by assumption), so the GGM is a disconnected graph, and  $\mathbf{\Lambda}_t$  is diagonal. However, as the robot moves about, it will induce correlation between nearby landmarks. Intuitively this is because the robot is estimating its position based on distance to the landmarks, but the landmarks' locations are being estimated based on the robot's position, so they all become interdependent. This can be seen more clearly from the graphical model in Figure 13.8: it is clear that  $\mathbf{l}^1$  and  $\mathbf{l}^2$  are not d-separated by  $\mathbf{y}_{1:t}$ , because there is a path between them via the unknown sequence of  $\mathbf{r}_{1:t}$  nodes. Consequently, the precision matrix becomes denser over time. As a consequence of the precision matrix becoming denser, each inference step takes  $O(K^3)$  time. This prevents the method from being applied to large maps.

One way to speed this up is based on the following observation: conditional on knowing the robot's path,  $\mathbf{r}_{1:t}$ , the landmark locations are independent, i.e.,  $p(\mathbf{l}_t | \mathbf{r}_{1:t}, \mathbf{y}_{1:t}) = \prod_{k=1}^K p(\mathbf{l}_t^k | \mathbf{r}_{1:t}, \mathbf{y}_{1:t})$ . This can be seen by looking at the DGM in Figure 13.8. We can therefore sample the trajectory using

1 some proposal, and apply (2d) Kalman filtering to each landmark independently. This is an example  
2 of RBPF, and reduces the inference cost to  $O(NK)$ , where  $N$  is the number of particles and  $K$  is  
3 the number of landmarks.  
4

5 The overall cost of this technique is  $O(NK)$  per step. Fortunately, the number of particles  $N$   
6 needed for good performance is quite small, so the algorithm is essentially linear in the number of  
7 landmarks, making it quite scalable. This idea was first suggested in [Mur00], who applied it to  
8 grid-structured occupancy grids (and used the HMM filter for each particle). It was subsequently  
9 extended to landmark-based maps in [Thr+04], using the Kalman filter for each particle; they called  
10 the technique **FastSLAM**.

## 11 13.5 Extensions of the particle filter

12 There are many extensions to the basic particle filtering algorithm, such as the following:  
13

- 14 • We can increase particle diversity by applying one or more steps of MCMC sampling (Section 12.2)  
15 at each PF step using  $\pi_t(\mathbf{z}_t)$  as the target distribution. This is called the **resample-move**  
16 algorithm [DJ11]. It is also possible to use SMC instead of MCMC to diversify the samples  
17 [GM17].  
18
- 19 • We can extend PF to the case of offline inference; this is called **particle smoothing** (see e.g.,  
20 [Kla+06]).  
21
- 22 • We can extend PF to inference in general graphical models (not just chains) by combining PF  
23 with loopy belief propagation (Section 9.3); this is called **non-parametric BP** or **particle BP**  
24 (see e.g., [Sud+03; Isa03; Sud+10; Pac+14]).  
25
- 26 • We can extend PF to perform inference in static models (e.g., for parameter inference), as we  
27 discuss in Section 13.6.  
28

## 30 13.6 SMC samplers

31 In this section, we discuss **SMC samplers** (sequential Monte Carlo samplers), which are a way  
32 to apply particle filters to sample from a generic target distribution,  $\pi(\mathbf{z}) = \tilde{\gamma}(\mathbf{z})/Z$ , rather than  
33 requiring the model to be an SSM. Thus SMC is an alternative to MCMC.  
34

35 The advantages of SMC samplers over MCMC are as follows: we can estimate the normalizing  
36 constant  $Z$ ; we can more easily develop adaptive versions that tune the transition kernel using the  
37 current set of samples; and the method is easier to parallelize (see e.g., [CCS22; Gre+22]).  
38

39 The method works by defining a sequence of intermediate distributions,  $\pi_t(\mathbf{z}_t)$ , which we expand  
40 to a sequence of distributions over all the past variables,  $\bar{\pi}_t(\mathbf{z}_{1:t})$ . We then use the particle filtering  
41 algorithm to sample from each of these intermediate distributions. By marginalizing all but the final  
42 state, we recover samples from the target distribution,  $\pi(\mathbf{z}) = \sum_{\mathbf{z}_{1:T-1}} \bar{\pi}_T(\mathbf{z}_{1:T})$ , as we explain below.  
43 (For more details, see e.g., [Dai+20a; CP20b].)

### 44 13.6.1 Ingredients of an SMC sampler

45 To define an SMC sampler, we need to specify several ingredients:  
46

- A sequence of distributions defined on the same state space,  $\pi_t(\mathbf{z}_t) = \tilde{\gamma}_t(\mathbf{z}_t)/Z_t$ , for  $t = 0 : T$ ;
- A **forwards kernel**  $M_t(\mathbf{z}_t|\mathbf{z}_{t-1})$  (often written as  $M_t(\mathbf{z}_{t-1}, \mathbf{z}_t)$ ), which satisfies  $\sum_{\mathbf{z}_t} M_t(\mathbf{z}_t|\mathbf{z}_{t-1}) = 1$ . This can be used to propose new samples from our current estimate when we apply particle filtering.
- A **backwards kernel**  $L_t(\mathbf{z}_t|\mathbf{z}_{t+1})$  (often written as  $L(\mathbf{z}_t, \mathbf{z}_{t+1})$ ), which satisfies  $\sum_{\mathbf{z}_t} L_t(\mathbf{z}_t|\mathbf{z}_{t+1}) = 1$ . This allows us to create a sequence of variables by working backwards in time from the final target value to the first time step. In particular, we create the following joint distribution:

$$\bar{\pi}_t(\mathbf{z}_{1:t}) = \pi_t(\mathbf{z}_t) \prod_{s=1}^{t-1} L_s(\mathbf{z}_s|\mathbf{z}_{s+1}) \quad (13.64)$$

This satisfies  $\sum_{\mathbf{z}_{1:t-1}} \bar{\pi}_t(\mathbf{z}_{1:t}) = \pi_t(\mathbf{z}_t)$ , so if we apply particle filtering to this for  $t = 1 : T$ , then samples from the “end” of such sequences will be from the target distribution  $\pi_t$ .

With the above ingredients, we can compute the incremental weight at step  $t$  using

$$\alpha_t = \frac{\bar{\pi}_t(\mathbf{z}_{1:t})}{\bar{\pi}_{t-1}(\mathbf{z}_{1:t-1})M_t(\mathbf{z}_t|\mathbf{z}_{t-1})} \propto \frac{\tilde{\gamma}_t(\mathbf{z}_t)}{\tilde{\gamma}_{t-1}(\mathbf{z}_{t-1})} \frac{L_{t-1}(\mathbf{z}_{t-1}|\mathbf{z}_t)}{M_t(\mathbf{z}_t|\mathbf{z}_{t-1})} \quad (13.65)$$

This can be plugged into the generic SMC algorithm, Algorithm 33.

We still have to specify the forwards and backwards kernels. We will assume the forwards kernel  $M_t$  is an MCMC kernel that leaves  $\pi_t$  invariant. We can then define the backwards kernel to be the **time reversal** of the forwards kernel. More precisely, suppose we define  $L_{t-1}$  so it satisfies

$$\pi_t(\mathbf{z}_t)L_{t-1}(\mathbf{z}_{t-1}|\mathbf{z}_t) = \pi_t(\mathbf{z}_{t-1})M_t(\mathbf{z}_t|\mathbf{z}_{t-1}) \quad (13.66)$$

In this case, the incremental weight simplifies as follows:

$$\alpha_t = \frac{Z_t \pi_t(\mathbf{z}_t) L_{t-1}(\mathbf{z}_{t-1}|\mathbf{z}_t)}{Z_{t-1} \pi_{t-1}(\mathbf{z}_{t-1}) M_t(\mathbf{z}_t|\mathbf{z}_{t-1})} \quad (13.67)$$

$$= \frac{Z_t \pi_t(\mathbf{z}_{t-1}) M_t(\mathbf{z}_t|\mathbf{z}_{t-1})}{Z_{t-1} \pi_{t-1}(\mathbf{z}_{t-1}) M_t(\mathbf{z}_t|\mathbf{z}_{t-1})} \quad (13.68)$$

$$= \frac{\tilde{\gamma}_t(\mathbf{z}_{t-1})}{\tilde{\gamma}_{t-1}(\mathbf{z}_{t-1})} \quad (13.69)$$

We can use any kind of MCMC kernel for  $M_t$ . For example, if the parameters are real valued and unconstrained, we can use a Markov kernel that corresponds to  $K$  steps of a random walk Metropolis-Hastings sampler. We can set the covariance of the proposal to  $\delta^2 \hat{\Sigma}_{t-1}$ , where  $\hat{\Sigma}_{t-1}$  is the empirical covariance of the weighted samples from the previous step,  $(W_{t-1}^{1:N}, \mathbf{z}_{t-1}^{1:N})$ , and  $\delta = 2.38D^{-3/2}$  (which is the optimal scaling parameter for RWMH). In high dimensional problems, we can use gradient based Markov kernels, such as HMC [BCJ20] and NUTS [Dev+21]. For binary state spaces, we can use the method of [SC13].

1 **13.6.2 Likelihood tempering (geometric path)**

2 There are many ways to specify the intermediate target distributions. In the **geometric path**  
3 method, we specify the intermediate distributions to be  
4

5  $\tilde{\gamma}_t(\mathbf{z}) = \tilde{\gamma}_0(\mathbf{z})^{1-\lambda_t} \tilde{\gamma}(\mathbf{z})^{\lambda_t}$  (13.70)

6 where  $0 = \lambda_0 < \lambda_1 < \dots < \lambda_T = 1$  are **inverse temperature** parameters, and  $\tilde{\gamma}_0$  is the initial  
7 proposal. If we apply particle filtering to this model, but “turn off” the resampling step, the method  
8 becomes equivalent to **annealed importance sampling** (Section 11.5.4).  
9

10 In the context of Bayesian parameter inference, we often denote the latent variable  $\mathbf{z}$  by  $\boldsymbol{\theta}$ , we  
11 define  $\tilde{\gamma}_0(\boldsymbol{\theta}) \propto \pi_0(\boldsymbol{\theta})$  as the prior, and  $\tilde{\gamma}(\boldsymbol{\theta}) = \pi_0(\boldsymbol{\theta})p(\mathcal{D}|\boldsymbol{\theta})$  as the posterior. We can then define the  
12 intermediate distributions to be  
13

14  $\tilde{\gamma}_t(\boldsymbol{\theta}) = \pi_0(\boldsymbol{\theta})^{1-\lambda_t} \pi_0(\boldsymbol{\theta})^{\lambda_t} p(\mathcal{D}|\boldsymbol{\theta})^{\lambda_t} = \pi_0(\boldsymbol{\theta})^{1-\lambda_t} \exp[-\lambda_t \mathcal{E}(\boldsymbol{\theta})]$  (13.71)

15 where  $\mathcal{E}(\boldsymbol{\theta}) = -\log p(\mathcal{D}, \boldsymbol{\theta})$  is the energy (potential) function. The incremental weights are given by  
16

17  $\alpha_t(\boldsymbol{\theta}) = \frac{\pi_0(\boldsymbol{\theta})^{1-\lambda_t} \exp[-\lambda_t \mathcal{E}(\boldsymbol{\theta})]}{\pi_0(\boldsymbol{\theta})^{1-\lambda_t} \exp[-\lambda_{t-1} \mathcal{E}(\boldsymbol{\theta})]} = \exp[-\delta_t \mathcal{E}(\boldsymbol{\theta})]$  (13.72)

18 where  $\lambda_t = \lambda_{t-1} + \delta_t$ .  
19

20 For this method to work well, it is important to choose the  $\lambda_t$  so that the successive distributions  
21 are “equidistant”; this is called **adaptive tempering**. In the case of a Gaussian prior and Gaussian  
22 energy, one can show [CP20b] that this can be achieved by picking  $\lambda_t = (1 + \gamma)^{t+1} - 1$ , where  $\gamma > 0$   
23 is some constant. Thus we should increase  $\lambda$  slowly at first, and then make bigger and bigger steps.  
24

25 In practice we can estimate  $\lambda_t$  by setting  $\lambda_t = \lambda_{t-1} + \delta_t$ , where  
26

27  $\delta_t = \underset{\delta \in [0, 1 - \lambda_{t-1}]}{\operatorname{argmin}} (\operatorname{ESSLW}(\{-\delta \mathcal{E}(\boldsymbol{\theta}_t^n)\}) - \operatorname{ESS}_{\min})$  (13.73)

28 where  $\operatorname{ESSLW}(\{l_n\}) = \operatorname{ESS}(\{e^{l_n}\})$  computes the ESS (Equation (13.32)) from the log weights,  
29  $l_n = \log \tilde{w}^n$ . This ensures the change in the ESS across steps is close to the desired minimum ESS,  
30 typically  $0.5N$ . (If there is no solution for  $\delta$  in the interval, we set  $\delta_t = 1 - \lambda_{t-1}$ .) See Algorithm 35  
31 for the overall algorithm.  
32

33 **13.6.2.1 Example: sampling from a 1d bimodal distribution**

34 Consider the simple distribution  
35

36  $p(\boldsymbol{\theta}) \propto \mathcal{N}(\boldsymbol{\theta} | \mathbf{0}, \mathbf{I}) \exp(-\mathcal{E}(\boldsymbol{\theta}))$  (13.74)

37 where  $\mathcal{E}(\boldsymbol{\theta}) = c(\|\boldsymbol{\theta}\|^2 - 1)^2$ . We plot this in 1d in Figure 13.10a for  $c = 5$ ; we see that it has a  
38 bimodal shape, since the low energy states correspond to parameter vectors whose norm is close to 1.  
39

40 SMC is particularly useful for sampling from multimodal distributions, which can be provably  
41 hard to efficiently sample from using other methods, including HMC [MPS18], since gradients only  
42 provide local information about the curvature. As an example, in Figure 13.11a and Figure 13.11b  
43 we show the result of applying HMC (Section 12.5) and NUTS (Section 12.5.4.1) to this problem.  
44 We see that both algorithms get stuck near the initial state of  $\boldsymbol{\theta}_0 = \mathbf{1}$ .  
45

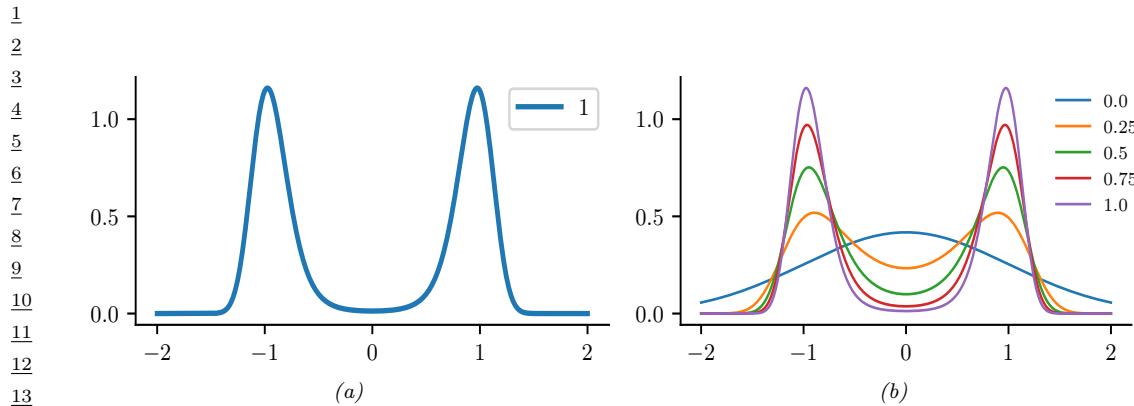


Figure 13.10: (a) Illustration of a bimodal target distribution. (b) Tempered versions of the target at different inverse temperatures, from  $\lambda_T = 1$  down to  $\lambda_1 = 0$ . Generated by [smc\\_tempered\\_1d\\_bimodal.ipynb](#).

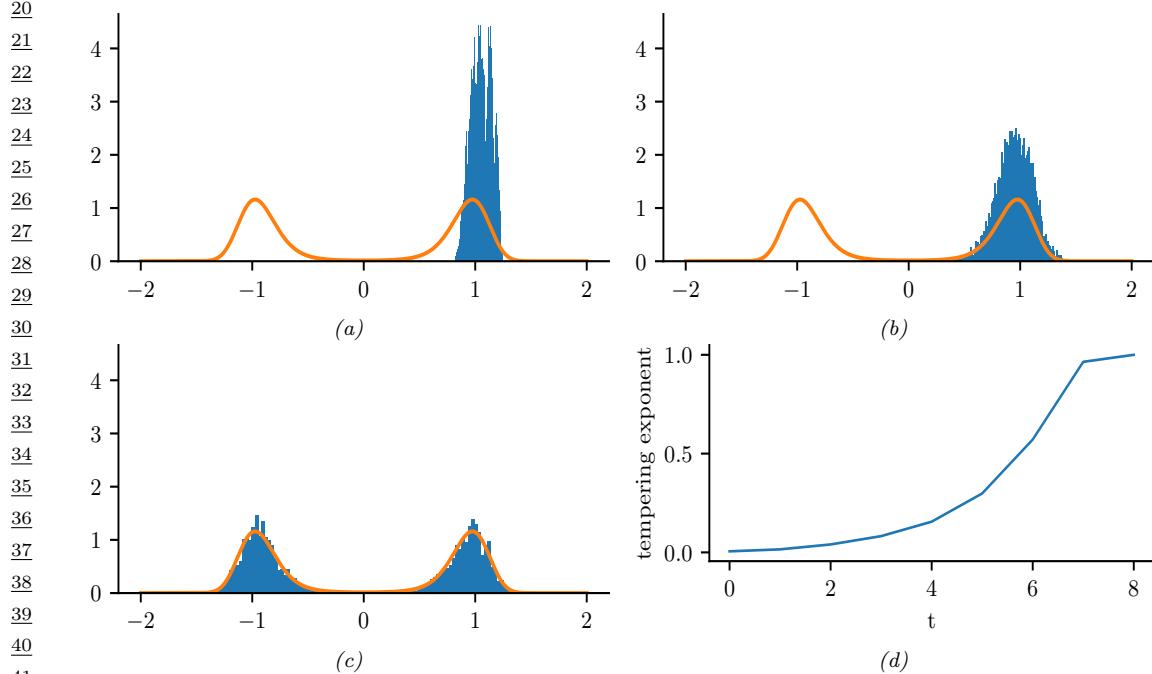


Figure 13.11: Sampling from the bimodal distribution in Figure 13.10a. (a) HMC. (b) NUTS. (c) Tempered SMC with HMC kernel (single step). (d) Adaptive inverse temperature schedule. Generated by [smc\\_tempered\\_1d\\_bimodal.ipynb](#).

---

**Algorithm 35:** SMC with adaptive tempering
 

---

```

1    $\lambda_{-1} = 0, t = -1, W_{-1}^n = 1$ 
2   while  $\lambda_t < 1$  do
3      $t = t + 1$ 
4     if  $t = 0$  then
5        $\boldsymbol{\theta}_0^n \sim \pi_0(\boldsymbol{\theta})$ 
6     else
7        $A_t^{1:N} = \text{Resample}(W_{t-1}^{1:N})$ 
8        $\boldsymbol{\theta}_t^n \sim M_{\lambda_{t-1}}(\boldsymbol{\theta}_{t-1}^{A_t^n}, \cdot)$ 
9     Compute  $\delta_t$  using Equation (13.73)
10     $\lambda_t = \lambda_{t-1} + \delta_t$ 
11     $\tilde{w}_t^n = \exp[-\delta \mathcal{E}(\boldsymbol{\theta}_t^n)]$ 
12     $W_t^n = \tilde{w}_t^n / (\sum_{m=1}^N \tilde{w}_t^m)$ 
13
14
15
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

```

---

In Figure 13.10b, we show tempered versions of the target distribution at 5 different temperatures, chosen uniformly in the interval  $[0, 1]$ . We see that at  $\lambda_1 = 0$ , the tempered target is equal to the Gaussian prior (blue line), which is easy to sample from. Each subsequent distribution is close to the previous one, so SMC can track the change until it ends up at the target distribution with  $\lambda_T = 1$ , as shown in Figure 13.11c.

These SMC results were obtained using the adaptive tempering scheme described above. In Figure 13.11d we see that initially the temperature is small, and then it increases exponentially. The algorithm takes 8 steps until  $\lambda_T \geq 1$ .

### 13.6.3 Data tempering

If we have a set of iid observations, we can define the  $t$ 'th target to be

$$\tilde{\gamma}_t(\boldsymbol{\theta}) = p(\boldsymbol{\theta})p(\mathbf{y}_{1:t}|\boldsymbol{\theta}) \quad (13.75)$$

We can now apply SMC to this model. From Equation (13.69), the incremental weight becomes

$$\alpha_t(\boldsymbol{\theta}) = \frac{\tilde{\gamma}_t(\mathbf{z}_{t-1})}{\tilde{\gamma}_{t-1}(\mathbf{z}_{t-1})} = \frac{p(\boldsymbol{\theta})p(\mathbf{y}_{1:t}|\boldsymbol{\theta})}{p(\boldsymbol{\theta})p(\mathbf{y}_{1:t-1}|\boldsymbol{\theta})} = p(\mathbf{y}_t|\mathbf{y}_{1:t-1}, \boldsymbol{\theta}) \quad (13.76)$$

This can be plugged into the generic SMC algorithm in Algorithm 33.

Unfortunately, to sample from the MCMC kernel will typically take  $O(t)$  time, since the MH accept/reject step requires computing  $p(\boldsymbol{\theta}') \prod_{i=1}^t p(\mathbf{y}_{1:i}|\boldsymbol{\theta}')$  for any proposed  $\boldsymbol{\theta}'$ . Hence the total cost is  $O(T^2)$  if there are  $T$  observations. To reduce this, we can only sample parameters at times  $t$  when the ESS drops below a certain level; in the remaining steps, we just grow the sequence deterministically by repeating the previously sampled value. This technique was proposed in [Cho02], who called it the **iterated batch importance sampling** or **IBIS** algorithm.

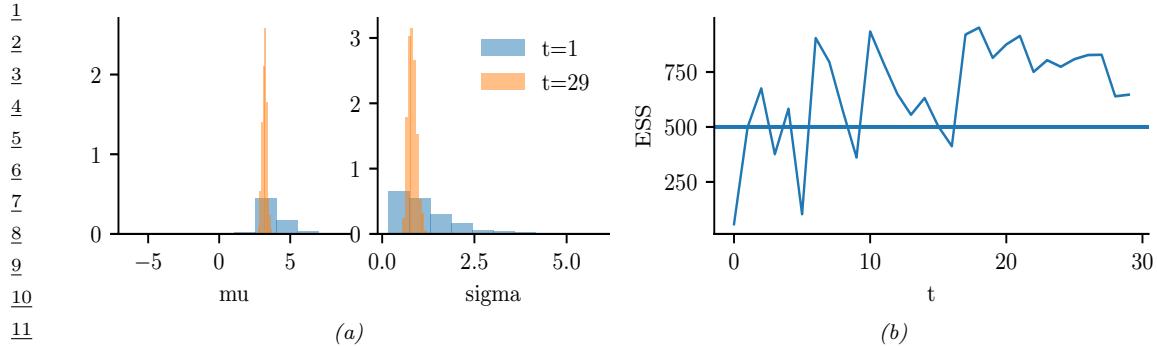


Figure 13.12: Illustration of IBIS applied to 30 samples from  $\mathcal{N}(\mu = 3.14, \sigma = 1)$ . (a) Posterior approximation after  $t = 1$  and  $t = 29$  observations. (b) Effective sample size over time. The sudden jumps up occur whenever resampling is triggered, which happens when the ESS drops below 500. Generated by `smc_ibis_1d.ipynb`.

### 13.6.3.1 Example: IBIS for a 1d Gaussian

In this section, we give a simple example of IBIS applied to data from a 1d Gaussian,  $y_t \sim \mathcal{N}(\mu = 3.14, \sigma = 1)$  for  $t = 1 : 30$ . The unknowns are  $\theta = (\mu, \sigma)$ . The prior is  $p(\theta) = \mathcal{N}(\mu|0, 1)\text{Ga}(\sigma|a = 1, b = 1)$ . We use IBIS with an adaptive RWMH kernel. We use  $N = 20$  particles, each updated for  $K = 50$  MCMC steps, so we collect 1000 samples per time step.

Figure 13.12a shows the approximate posterior after  $t = 1$  and  $t = 29$  time steps. We see that the posterior concentrates on the true values of  $\mu = 3.14$  and  $\sigma = 1$ .

Figure 13.12b plots the ESS vs time. The number of particles is 1000, and resampling (and MCMC moves) is triggered whenever this drops below 500. We see that we only need to invoke MCMC updates 3 times.

### 13.6.4 Sampling rare events and extrema

Suppose we want to sample values from  $\pi_0(\theta)$  conditioned on the event that  $S(\theta) > \lambda^*$ , where  $S$  is some score or “fitness” function. If  $\lambda^*$  is in the tail of the score distribution, this corresponds to sampling a **rare event**, which can be hard.

One approach is to use SMC to sample from a sequence of distributions with gradually increasing thresholds:

$$\pi_t(\theta) = \frac{1}{Z_t} \mathbb{I}(S(\theta) \geq \lambda_t) \pi_0(\theta) \quad (13.77)$$

with  $\lambda_0 < \dots < \lambda_T = \lambda^*$ . We can then use likelihood tempering, where the “likelihood” is the function

$$G_t(\theta_t) = \mathbb{I}(S(\theta_t) \geq \lambda_t) \quad (13.78)$$

We can use SMC to generate samples from the final distribution  $\pi_T$ . We may also be interested in estimating

$$Z_T = p(S(\theta) \geq \lambda_T) \quad (13.79)$$

1 where the probability is taken wrt  $\pi_0(\boldsymbol{\theta})$ .  
 2

3 We can adaptively set the thresholds  $\lambda_t$  as follows: at each step, sort the samples by their score,  
 4 and set  $\lambda_t$  to the  $\alpha$ 'th highest quantile. For example, if we set  $\alpha = 0.5$ , we keep the top 50% fittest  
 5 particles. This ensures the ESS equals the minimum threshold at each step. For details, see [Cér+12].

6 Note that this method is very similar to the **cross-entropy method** (Section 6.9.5). The difference  
 7 is that CEM fits a parametric distribution (e.g., a Gaussian) to the particles at each step and samples  
 8 from that, rather than using a Markov kernel.  
 9

### 10 13.6.5 SMC-ABC and likelihood-free inference

11 The term **likelihood-free inference** refers to estimating the parameters  $\boldsymbol{\theta}$  of a black box from  
 12 which we can sample data,  $\mathbf{y} \sim p(\cdot|\boldsymbol{\theta})$ , but where we cannot evaluate  $p(\mathbf{y}|\boldsymbol{\theta})$  pointwise. Such models  
 13 are called simulators, so this approach to inference is also called **simulation-based inference** (see  
 14 e.g., [Nea+08; CBL20; Gou+96]). These models are also called **implicit models** (see Section 26.1).

15 If we want to approximate the posterior of a model with no known likelihood, we can use  
 16 **Approximate Bayesian Computation** or **ABC** (see e.g., [Bea19; SFB18; Gut+14; Pes+21]).  
 17 In this setting, we sample both parameters  $\boldsymbol{\theta}$  and synthetic data  $\mathbf{y}$  such that the synthetic data  
 18 (generated from  $\boldsymbol{\theta}$ ) is sufficiently close to the observed data  $\mathbf{y}^*$ , as judged by some distance score,  
 19  $d(\mathbf{y}, \mathbf{y}^*) < \epsilon$ . (For high dimensional problems, we typically require  $d(\mathbf{s}(\mathbf{y}), \mathbf{s}(\mathbf{y}^*)) < \epsilon$ , where  $\mathbf{s}(\mathbf{y})$  is  
 20 a low-dimensional summary statistic of the data.)  
 21

22 In **SMC-ABC**, we gradually decrease the discrepancy  $\epsilon$  to get a series of distributions as follows:

$$23 \quad \pi_t(\boldsymbol{\theta}, \mathbf{y}) = \frac{1}{Z_t} \pi_0(\boldsymbol{\theta}) p(\mathbf{y}|\boldsymbol{\theta}) \mathbb{I}(d(\mathbf{y}, \mathbf{y}^*) < \epsilon_t) \quad (13.80)$$

26 where  $\epsilon_0 > \epsilon_1 > \dots$ . This is similar to the rare event SMC samplers in Section 13.6.4, except that  
 27 we can't directly evaluate the quality of a candidate  $\boldsymbol{\theta}$ , instead we must first convert it to data space  
 28 and make the comparison there. For details, see [DMDJ12].  
 29

30 Although SMC-ABC is popular in some fields, such as genetics and epidemiology, this method  
 31 is quite slow and does not scale to high dimensional problems. In such settings, a more efficient  
 32 approach is to train a generative model to **emulate** the simulator; if this model is parametric with a  
 33 tractable likelihood (e.g., a flow model), we can use the usual methods for posterior inference of its  
 34 parameters (including gradient based methods like HMC). See e.g., [Bre+20a] for details.  
 35

### 36 13.6.6 SMC<sup>2</sup>

37 We have seen how SMC can be a useful alternative to MCMC. However it requires that we can  
 38 efficiently evaluate the likelihood ratio terms  $\frac{\gamma_t(\boldsymbol{\theta}_t)}{\gamma_{t-1}(\boldsymbol{\theta}_t)}$ . In cases where this is not possible (e.g., for  
 39 latent variable models), we can use SMC (specifically the estimate  $\hat{Z}_t$  in Equation (13.10)) as a  
 40 subroutine to approximate these likelihoods. This is called **SMC<sup>2</sup>**. For details, see [CP20b, Ch. 18].  
 41

### 42 13.6.7 Variational filtering SMC

43 One way to improve SMC is to learn a proposal distribution (e.g., using a neural network) such that  
 44 the approximate posterior,  $\hat{\pi}_T(\mathbf{z}_{1:T}; \boldsymbol{\phi}, \boldsymbol{\theta})$ , is close to the target posterior,  $\pi_T(\mathbf{z}_{1:T}; \boldsymbol{\theta})$ , where  $\boldsymbol{\theta}$  are  
 45

the model parameters, and  $\phi$  are the proposal parameters (which may depend on  $\theta$ ). One can show [Nae+18] that the KL divergence between these distributions can be bounded as follows:

$$0 \leq D_{\text{KL}}(\mathbb{E}[\hat{\pi}_T(\mathbf{z}_{1:T})] \parallel \pi_T(\mathbf{z}_{1:T})) \leq -\mathbb{E}\left[\log \frac{\hat{Z}_T}{Z_T}\right] \quad (13.81)$$

where

$$Z_T(\theta) = p_\theta(\mathbf{y}_{1:T}) = \int p_\theta(\mathbf{z}_{1:T}, \mathbf{y}_{1:T}) d\mathbf{z}_{1:T} \quad (13.82)$$

Hence

$$\mathbb{E}\left[\log \hat{Z}_T(\theta, \phi)\right] \leq \mathbb{E}[\log Z_T(\theta)] = \log Z_T(\theta) \quad (13.83)$$

Thus we can use SMC sampling to compute an unbiased approximation to  $\mathbb{E}[\log \hat{Z}_T(\theta, \phi)]$ , which is a lower bound on the evidence (log marginal likelihood).

We can now maximize this lower bound wrt  $\theta$  and  $\phi$  using SGD, as a way to learn both proposals and the model. Unfortunately, computing the gradient of the bound is tricky, since the resampling step is non-differentiable. However, in practice one can ignore the dependence of the resampling operator on the parameters, or one can use differentiable approximations (see e.g., [Ros+22]). This overall approach was independently proposed in several papers: the **FIVO** (filtering variational objective) paper [Mad+17], the **variational SMC** paper [Nae+18] and the **auto-encoding SMC** paper [Le+18].

### 13.6.8 Variational smoothing SMC

The methods in Section 13.6.7 use SMC in which the target distributions are defined to be the filtered distributions,  $\pi_t(\mathbf{z}_{1:t}) = p_\theta(\mathbf{z}_{1:t} | \mathbf{y}_{1:t})$ ; this is called **filtering SMC**. Unfortunately, this can work poorly when fitting models to offline sequence data, since at time  $t$ , all future observations are ignored in the objective, no matter how good the proposal. This can create situations where future observations are unlikely given the current set of sampled trajectories, which can result in particle impoverishment and high variance in the estimate of the lower bound.

Recently, a new method called **SIXO** (smoothing inference with twisted objectives) was proposed in [Law+22] that uses the smoothing distributions as targets,  $\pi_t(\mathbf{z}_{1:t}) = p_\theta(\mathbf{z}_{1:t} | \mathbf{y}_{1:T})$ , to create a much lower variance variational lower bound. Of course it is impossible to directly compute this posterior, but we can approximate it using **twisted particle filters** [WL14a; AL+16]. In this approach, we approximate the (unnormalized) posterior using

$$p_\theta(\mathbf{z}_{1:t}, \mathbf{y}_{1:T}) = p_\theta(\mathbf{z}_{1:t}, \mathbf{y}_{1:t}) p_\theta(\mathbf{y}_{t+1:T} | \mathbf{z}_{1:t}, \mathbf{y}_{1:t}) \quad (13.84)$$

$$= p_\theta(\mathbf{z}_{1:t}, \mathbf{y}_{1:t}) p_\theta(\mathbf{y}_{t+1:T} | \mathbf{z}_t) \quad (13.85)$$

$$\approx p_\theta(\mathbf{z}_{1:t}, \mathbf{y}_{1:t}) r_\psi(\mathbf{y}_{t+1:T}, \mathbf{z}_t) \quad (13.86)$$

where  $r_\psi(\mathbf{y}_{t+1:T}, \mathbf{z}_t) \approx p_\theta(\mathbf{y}_{t+1:T} | \mathbf{z}_t)$  is the **twisting function**, which acts as a “**lookahead function**”.

1 One way to approximate the twisting function is to note that  
 2

$$3 \quad p_{\theta}(\mathbf{y}_{t+1:T} | \mathbf{z}_t) = \frac{p_{\theta}(\mathbf{z}_t | \mathbf{y}_{t+1:T}) p_{\theta}(\mathbf{y}_{t+1:T})}{p_{\theta}(\mathbf{z}_t)} \propto \frac{p_{\theta}(\mathbf{z}_t | \mathbf{y}_{t+1:T})}{p_{\theta}(\mathbf{z}_t)} \quad (13.87)$$

4 where we drop terms that are independent of  $\mathbf{z}_t$  since such terms will cancel out when we normalize  
 5 the sampling weights. We can approximate the density ratio using the binary classifier method of  
 6 Section 2.7.5. To do this, we define one distribution to be  $p_1 = p_{\theta}(\mathbf{z}_t, \mathbf{y}_{t+1:T})$  and the other to be  
 7  $p_2 = p_{\theta}(\mathbf{z}_t) p_{\theta}(\mathbf{y}_{t+1:T})$ , so that  $p_1/p_2 = \frac{p_{\theta}(\mathbf{z}_t | \mathbf{y}_{t+1:T})}{p_{\theta}(\mathbf{z}_t)}$ . We can easily draw a sample  
 8  $(\mathbf{z}_{1:T}, \mathbf{y}_{1:T}) \sim p_{\theta}$  using ancestral sampling, from which we can compute  $(\mathbf{z}_t, \mathbf{y}_{t+1:T}) \sim p_1$  by marginalization. We can  
 9 also sample a fresh sequence from  $(\tilde{\mathbf{z}}_{1:T}, \tilde{\mathbf{y}}_{1:T}) \sim p_{\theta}$  from which we can compute  $(\tilde{\mathbf{z}}_t, \tilde{\mathbf{y}}_{t+1:T}) \sim p_2$   
 10 by marginalization. We then use  $(\mathbf{z}_t, \mathbf{y}_{t+1:T})$  as a positive example and  $(\tilde{\mathbf{z}}_t, \tilde{\mathbf{y}}_{t+1:T})$  as a negative example when training the binary classifier,  $r_{\psi}(\mathbf{y}_{t+1:T}, \mathbf{z}_t)$ .

11 Once we have updated the twisting parameters  $\psi$ , we can rerun SMC to get a tighter lower bound  
 12 on the log marginal likelihood, which we can then optimize wrt the model parameters  $\theta$  and proposal  
 13 parameters  $\phi$ . Thus the overall method is a stochastic variational EM-like method for optimziing  
 14 the bound

$$15 \quad \mathcal{L}_{\text{SIXO}}(\theta, \phi, \psi, \mathbf{y}_{1:T}) \triangleq \mathbb{E} \left[ \log \hat{Z}_{\text{SIXO}}(\theta, \phi, \psi, \mathbf{y}_{1:T}) \right] \quad (13.88)$$

$$16 \quad \leq \log \mathbb{E} \left[ \hat{Z}_{\text{SIXO}}(\theta, \phi, \psi, \mathbf{y}_{1:T}) \right] = \log p_{\theta}(\mathbf{y}_{1:T}) \quad (13.89)$$

17 In [Law+22] they prove the following: suppose the true model  $p^*$  is an SSM in which the optimal  
 18 proposal function for the model satisfies  $p^*(\mathbf{z}_t | \mathbf{z}_{1:t-1}, \mathbf{y}_{1:T}) \in \mathcal{Q}$ , and the optimal lookahead function  
 19 for the model satisfies  $p^*(\mathbf{y}_{t+1:T} | \mathbf{z}_t) \in \mathcal{R}$ . Furthermore, assume the SIXO objective has a unique  
 20 maximizer. Then, at the optimum, we have that the learned proposal  $q_{\phi^*}(\mathbf{z}_t | \mathbf{z}_{1:t-1}, \mathbf{y}_{1:T}) \in \mathcal{Q}$  is  
 21 equal to the optimal proposal, the learned twisting function  $r_{\psi^*}(\mathbf{y}_{t+1:T}, \mathbf{z}_t) \in \mathcal{R}$  is equal to the  
 22 optimal lookahead, and the lower bound is tight (i.e.,  $\mathcal{L}_{\text{SIXO}}(\theta^*, \phi^*, \psi^*) = p^*(\mathbf{y}_{1:T})$ ) for any number  
 23 of samples  $N_s \geq 1$  and for any kind of SSM  $p^*$ . (This is in contrast to the FIVO bound, whiere the  
 24 bound does not usually become tight.)

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47



PART III

## Prediction



# 14 Predictive models: an overview

## 14.1 Introduction

The vast majority of machine learning is concerned with tackling a single problem, namely learning to predict outputs  $\mathbf{y}$  from inputs  $\mathbf{x}$  using some function  $f$  that is estimated from a labeled training set  $\mathcal{D} = \{(\mathbf{x}_n, \mathbf{y}_n) : n = 1 : N\}$ , for  $\mathbf{x}_n \in \mathcal{X} \subseteq \mathbb{R}^D$  and  $\mathbf{y}_n \in \mathcal{Y} \subseteq \mathbb{R}^C$ . We can model our uncertainty about the correct output for a given input using a conditional probability model of the form  $p(\mathbf{y}|f(\mathbf{x}))$ . When  $\mathcal{Y}$  is a discrete set of labels, this is called (in the ML literature) a **discriminative model**, since it lets us discriminate (distinguish) between the different possible values of  $\mathbf{y}$ . If the output is real-valued,  $\mathcal{Y} = \mathbb{R}$ , this is called a **regression model**. (In the statistics literature, the term “regression model” is used in both cases, even if  $\mathcal{Y}$  is a discrete set.) We will use the more generic term **“predictive model”** to refer to such models.

A predictive model can be considered as a special case of a conditional generative model (discussed in Chapter 20). In a predictive model, the output is usually low dimensional, and there is a single best answer that we want to predict. However, in most generative models, the output is usually high dimensional, such as images or sentences, and there may be many correct outputs for any given input. We will discuss a variety of types of predictive model in Section 14.1.1, but we defer the details to subsequent chapters. The rest of this chapter then discusses issues that are relevant to all types of predictive model, regardless of the specific form, such as evaluation.

### 14.1.1 Types of model

There are many different kinds of predictive model  $p(\mathbf{y}|\mathbf{x})$ . The biggest distinction is between **parametric models**, that have a fixed number of parameters independent of the size of the training set, and **non-parametric models** that have a variable number of parameters that grows with the size of the training set. Non-parametric models are usually more flexible, but can be slower to use for prediction. Parametric models are usually less flexible, but are faster to use for prediction.

Most non-parametric models are based on comparing a test input  $\mathbf{x}$  to some or all of the stored training examples  $\{\mathbf{x}_n, n = 1 : N\}$ , using some form of similarity,  $s_n = \mathcal{K}(\mathbf{x}, \mathbf{x}_n) \geq 0$ , and then predicting the output using some weighted combination of the training labels, such as  $\hat{\mathbf{y}} = \sum_{n=1}^N s_n \mathbf{y}_n$ . A typical example is a Gaussian process, which we discuss in Chapter 18. Other examples, such as  $K$ -nearest neighbor models, are discussed in the prequel to this book, [Mur22].

Most parametric models have the form  $p(\mathbf{y}|\mathbf{x}) = p(\mathbf{y}|f(\mathbf{x}; \boldsymbol{\theta}))$ , where  $f$  is some kind of function that predicts the parameters (e.g., the mean, or logits) of the output distribution (e.g., Gaussian or categorical). There are many kinds of function we can use. If  $f$  is a linear function of  $\boldsymbol{\theta}$  (i.e.,

$f(\mathbf{x}; \boldsymbol{\theta}) = \boldsymbol{\theta}^\top \phi(\mathbf{x})$  for some *fixed* feature transformation  $\phi$ ), then the model is called a generalized linear model or GLM, which we discuss in Chapter 15. If  $f$  is a non-linear, but differentiable, function of  $\boldsymbol{\theta}$  (e.g.,  $f(\mathbf{x}; \boldsymbol{\theta}) = \boldsymbol{\theta}_2^\top \phi(\mathbf{x}; \boldsymbol{\theta}_1)$  for some learnable function  $\phi(\mathbf{x}; \boldsymbol{\theta}_1)$ ), then it is common to represent  $f$  using a neural network (Chapter 16). Other types of predictive model, such as decision trees and random forests, are discussed in the prequel to this book, [Mur22].

### 14.1.2 Model fitting using ERM, MLE and MAP

In this section, we briefly discuss some methods used for fitting (parametric) models. The most common approach is to use **maximum likelihood estimation** or **MLE**, which amounts to solving the following optimization problem:

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmax}} p(\mathcal{D}|\boldsymbol{\theta}) = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmax}} \log p(\mathcal{D}|\boldsymbol{\theta}) \quad (14.1)$$

If the dataset is  $N$  iid data samples, the likelihood decomposes into a product of terms,  $p(\mathcal{D}|\boldsymbol{\theta}) = \prod_{n=1}^N p(\mathbf{y}_n|\mathbf{x}_n, \boldsymbol{\theta})$ . Thus we can instead minimize the following (scaled) **negative log likelihood**:

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmin}} \frac{1}{N} \sum_{n=1}^N [-\log p(\mathbf{y}_n|\mathbf{x}_n, \boldsymbol{\theta})] \quad (14.2)$$

We can generalize this by replacing the **log loss**  $\ell_n(\boldsymbol{\theta}) = -\log p(\mathbf{y}_n|\mathbf{x}_n, \boldsymbol{\theta})$  with a more general loss function to get

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmin}} r(\boldsymbol{\theta}) \quad (14.3)$$

where  $r(\boldsymbol{\theta})$  is the **empirical risk**

$$r(\boldsymbol{\theta}) = \frac{1}{N} \sum_{n=1}^N \ell_n(\boldsymbol{\theta}) \quad (14.4)$$

This approach is called **empirical risk minimization** or **ERM**.

ERM can easily result in **overfitting**, so it is common to add a penalty or regularizer term to get

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmin}} r(\boldsymbol{\theta}) + \lambda C(\boldsymbol{\theta}) \quad (14.5)$$

where  $\lambda \geq 0$  controls the degree of regularization, and  $C(\boldsymbol{\theta})$  is some complexity measure. If we use log loss, and we define  $C(\boldsymbol{\theta}) = -\log \pi_0(\boldsymbol{\theta})$ , where  $\pi_0(\boldsymbol{\theta})$  is some prior distribution, and we use  $\lambda = 1$ , we recover the **MAP estimate**

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmax}} \log p(\mathcal{D}|\boldsymbol{\theta}) + \log \pi_0(\boldsymbol{\theta}) \quad (14.6)$$

This can be solved using standard optimization methods (see Chapter 6).

---

### 14.1.3 Model fitting using Bayes, VI and generalized Bayes

Another way to prevent overfitting is to estimate a *probability distribution over parameters*,  $q(\boldsymbol{\theta})$ , instead of a point estimate. That is, we can try to estimate the ERM in expectation:

$$\hat{q} = \operatorname{argmin}_{q \in \mathcal{P}(\Theta)} \mathbb{E}_{q(\boldsymbol{\theta})} [r(\boldsymbol{\theta})] \quad (14.7)$$

If  $\mathcal{P}(\Theta)$  is the space of all probability distributions over parameters, then the solution will converge to a delta function that puts all its probability on the MLE. Thus this approach, on its own, will not prevent overfitting. However, we can regularize the problem by preventing the distribution from moving too far from the prior. If we measure the divergence between  $q$  and the prior using KL divergence, we get

$$\hat{q} = \operatorname{argmin}_{q \in \mathcal{P}(\Theta)} \mathbb{E}_{q(\boldsymbol{\theta})} [r(\boldsymbol{\theta})] + \frac{1}{\lambda} D_{\text{KL}}(q \parallel \pi_0) \quad (14.8)$$

The solution to this problem is known as the **Gibbs posterior**, and is given by the following:

$$\hat{q}(\boldsymbol{\theta}) = \frac{e^{-\lambda r(\boldsymbol{\theta})} \pi_0(\boldsymbol{\theta})}{\int e^{-\lambda r(\boldsymbol{\theta}')} \pi_0(\boldsymbol{\theta}') d\boldsymbol{\theta}'} \quad (14.9)$$

This is widely used in the **PAC-Bayes** community (see e.g., [Alq21]).

Now suppose we use log loss, and set  $\lambda = N$ , to get

$$\hat{q}(\boldsymbol{\theta}) = \frac{e^{-\sum_{n=1}^N \log p(\mathbf{y}_n | \mathbf{x}_n, \boldsymbol{\theta})} \pi_0(\boldsymbol{\theta})}{\int e^{-\sum_{n=1}^N \log p(\mathbf{y}_n | \mathbf{x}_n, \boldsymbol{\theta}')} \pi_0(\boldsymbol{\theta}') d\boldsymbol{\theta}'} \quad (14.10)$$

Then the resulting distribution is equivalent to the Bayes posterior:

$$\hat{q}(\boldsymbol{\theta}) = \frac{p(\mathcal{D} | \boldsymbol{\theta}) \pi_0(\boldsymbol{\theta})}{\int p(\mathcal{D} | \boldsymbol{\theta}') \pi_0(\boldsymbol{\theta}') d\boldsymbol{\theta}'} \quad (14.11)$$

Often computing the Bayes posterior is intractable. We can simplify the problem by restricting attention to a limited family of distributions,  $\mathcal{Q}(\Theta) \subset \mathcal{P}(\Theta)$ . This gives rise to the following objective:

$$\hat{q} = \operatorname{argmin}_{q \in \mathcal{Q}(\Theta)} \mathbb{E}_{q(\boldsymbol{\theta})} [-\log p(\mathcal{D} | \boldsymbol{\theta})] + D_{\text{KL}}(q \parallel \pi_0) \quad (14.12)$$

This is known as **variational inference**; see Chapter 10 for details. (See also Section 6.7, where we discuss the Bayesian learning rule.)

We can generalize this by replacing the negative log likelihood with a general risk,  $r(\boldsymbol{\theta})$ . Furthermore, we can replace the KL with a general divergence,  $D(q || \pi_0)$ , which we can weight using a general  $\lambda$ . This gives rise to the following objective:

$$\hat{q} = \operatorname{argmin}_{q \in \mathcal{Q}(\Theta)} \mathbb{E}_{q(\boldsymbol{\theta})} [r(\boldsymbol{\theta})] + \lambda D(q || \pi_0) \quad (14.13)$$

This is called **generalized Bayesian inference** [BHW16; KJD19; KJD21].

1 **14.2 Evaluating predictive models**

3 In this section we discuss how to evaluate the quality of a trained discriminative model.

5 **14.2.1 Proper scoring rules**

7 It is common to measure performance of a predictive model using a **proper scoring rule** [GR07a],  
8 which is defined as follows. Let  $S(p_{\theta}, (y, \mathbf{x}))$  be the score for predictive distribution  $p_{\theta}(y|\mathbf{x})$  when  
9 given an event  $y|\mathbf{x} \sim p^*(y|\mathbf{x})$ , where  $p^*$  is the true conditional distribution. (If we want to evaluate  
10 a Bayesian model, where we marginalize out  $\theta$  rather than condition on it, we just replace  $p_{\theta}(y|\mathbf{x})$   
11 with  $p(y|\mathbf{x}) = \int p_{\theta}(y|\mathbf{x})p(\theta|\mathcal{D})d\theta$ .) The expected score is defined by

13 
$$S(p_{\theta}, p^*) = \int p^*(\mathbf{x})p^*(y|\mathbf{x})S(p_{\theta}, (y, \mathbf{x}))dyd\mathbf{x} \quad (14.14)$$

16 A proper scoring rule is one where  $S(p_{\theta}, p^*) \leq S(p^*, p^*)$ , with equality iff  $p_{\theta}(y|\mathbf{x}) = p^*(y|\mathbf{x})$ . Thus  
17 maximizing such a proper scoring rule will force the model to match the true probabilities.

18 The log-likelihood,  $S(p_{\theta}, (y, \mathbf{x})) = \log p_{\theta}(y|\mathbf{x})$ , is a proper scoring rule. This follows from Gibbs  
19 inequality:

21 
$$S(p_{\theta}, p^*) = \mathbb{E}_{p^*(\mathbf{x})p^*(y|\mathbf{x})} [\log p_{\theta}(y|\mathbf{x})] \leq \mathbb{E}_{p^*(\mathbf{x})p^*(y|\mathbf{x})} [\log p^*(y|\mathbf{x})] \quad (14.15)$$

23 Therefore minimizing the NLL (aka log loss) should result in well-calibrated probabilities. However,  
24 in practice, log-loss can over-emphasize tail probabilities [QC+06].

25 A common alternative is to use the **Brier score** [Bri50], which is defined as follows:

27 
$$S(p_{\theta}, (y, \mathbf{x})) \triangleq \frac{1}{C} \sum_{c=1}^C (p_{\theta}(y=c|\mathbf{x}) - \mathbb{I}(y=c))^2 \quad (14.16)$$

30 This is just the squared error of the predictive distribution  $\mathbf{p} = p(1:C|\mathbf{x})$  compared to the one-hot  
31 label distribution  $\mathbf{y}$ . Since it based on squared error, the Brier score is less sensitive to extremely  
32 rare or extremely common classes. The Brier score is also a proper scoring rule.

34 **14.2.2 Calibration**

36 A model whose predicted probabilities match the empirical frequencies is said to be **calibrated**  
37 [Daw82; NMC05; Guo+17]. For example, if a classifier predicts  $p(y=c|\mathbf{x}) = 0.9$ , then we expect this  
38 to be the true label about 90% of the time. A well-calibrated model is useful to avoid making the  
39 wrong decision when the outcome is too uncertain. In the sections below, we discuss some ways to  
40 measure and improve calibration.

41

42 **14.2.2.1 Expected calibration error**

44 To assess calibration, we divide the predicted probabilities into a finite set of bins or buckets, and then  
45 assess the discrepancy between the empirical probability and the predicted probability by counting.  
46 More precisely, suppose we have  $B$  bins. Let  $\mathcal{B}_b$  be the set of indices of samples whose prediction

47

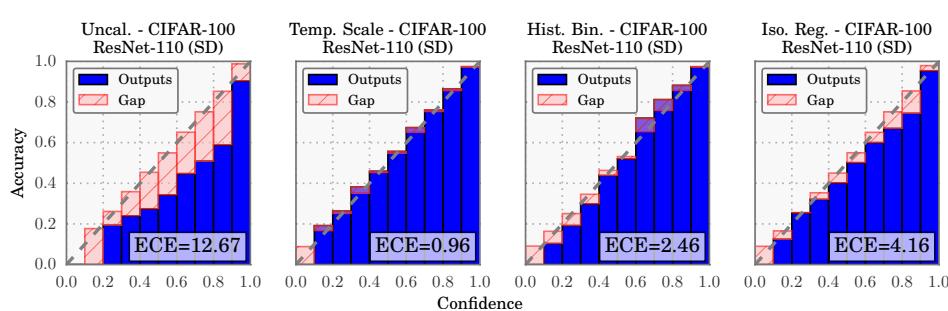


Figure 14.1: Reliability diagrams for the ResNet CNN image classifier [He+16b] applied to CIFAR-100 dataset. ECE is the expected calibration error, and measures the size of the red gap. Methods from left to right: original probabilities; after temperature scaling; after histogram binning; after isotonic regression. From Figure 4 of [Guo+17]. Used with kind permission of Chuan Guo.

confidence falls into the interval  $I_b = (\frac{b-1}{B}, \frac{b}{B}]$ . Here we use uniform bin widths, but we could also define the bins so that we can get an equal number of samples in each one.

Let  $f(\mathbf{x})_c = p(y = c|\mathbf{x})$ ,  $\hat{y}_n = \text{argmax}_{c \in \{1, \dots, C\}} f(\mathbf{x}_n)_c$ , and  $\hat{p}_n = \max_{c \in \{1, \dots, C\}} f(\mathbf{x}_n)_c$ . The accuracy within bin  $b$  is defined as

$$\text{acc}(\mathcal{B}_b) = \frac{1}{|\mathcal{B}_b|} \sum_{n \in \mathcal{B}_b} \mathbb{I}(\hat{y}_n = y_n) \quad (14.17)$$

The average confidence within this bin is defined as

$$\text{conf}(\mathcal{B}_b) = \frac{1}{|\mathcal{B}_b|} \sum_{n \in \mathcal{B}_b} \hat{p}_n \quad (14.18)$$

If we plot accuracy vs confidence, we get a **reliability diagram**, as shown in Figure 14.1. The gap between the accuracy and confidence is shown in the red bars. We can measure this using the **expected calibration error (ECE)** [NCH15]:

$$\text{ECE}(f) = \sum_{b=1}^B \frac{|\mathcal{B}_b|}{B} |\text{acc}(\mathcal{B}_b) - \text{conf}(\mathcal{B}_b)| \quad (14.19)$$

In the multiclass case, the ECE only looks at the error of the MAP (top label) prediction. We can extend the metric to look at all the classes using the **marginal calibration error**, proposed in [KLM19]:

$$\text{MCE} = \sum_{c=1}^C w_c \mathbb{E} [(p(Y = c|f(\mathbf{x})_c) - f(\mathbf{x})_c)^2] \quad (14.20)$$

$$= \sum_{c=1}^C w_c \sum_{b=1}^B \frac{|\mathcal{B}_{b,c}|}{B} (\text{acc}(\mathcal{B}_{b,c}) - \text{conf}(\mathcal{B}_{b,c}))^2 \quad (14.21)$$

where  $\mathcal{B}_{b,c}$  is the  $b$ 'th bin for class  $c$ , and  $w_c \in [0, 1]$  denotes the importance of class  $c$ . (We can set  $w_c = 1/C$  if all classes are equally important.) In [Nix+19], they call this metric **static calibration error**; they show that certain methods that have good ECE may have poor MCE. Other multi-class calibration metrics are discussed in [WLZ19].

### 14.2.2.2 Improving calibration

In principle, training a classifier so it optimizes a proper scoring rule (such as NLL) should automatically result in a well-calibrated classifier. In practice, however, unbalanced datasets can result in poorly calibrated predictions. Below we discuss various ways for improving the calibration of probabilistic classifiers, following [Guo+17].

### 14.2.2.3 Platt scaling

Let  $z$  be the log-odds, or logit, and  $p = \sigma(z)$ , produced by a probabilistic binary classifier. We wish to convert this to a more calibrated value  $q$ . The simplest way to do this is known as **Platt scaling**, and was proposed in [Pla00]. The idea is to compute  $q = \sigma(az + b)$ , where  $a$  and  $b$  are estimated via maximum likelihood on a validation set.

In the multiclass case, we can extend Platt scaling by using matrix scaling:  $\mathbf{q} = \text{softmax}(\mathbf{W}\mathbf{z} + \mathbf{b})$ , where we estimate  $\mathbf{W}$  and  $\mathbf{b}$  via maximum likelihood on a validation set. Since  $\mathbf{W}$  has  $K \times K$  parameters, where  $K$  is the number of classes, this method can easily overfit, so in practice we restrict  $\mathbf{W}$  to be diagonal.

### 14.2.2.4 Nonparametric (histogram) methods

Platt scaling makes a strong assumption about how the shape of the calibration curve. A more flexible, nonparametric, method is to partition the predicted probabilities into bins,  $p_m$ , and to estimate an empirical probability  $q_m$  for each such bin; we then replace  $p_m$  with  $q_m$ ; this is known as **histogram binning** [ZE01a]. We can regularize this method by requiring that  $q = f(p)$  be a piecewise constant, monotonically non-decreasing function; this is known as **isotonic regression** [ZE01a]. An alternative approach, known as the **scaling-binning calibrator**, is to apply a scaling method (such as Platt scaling), and then to apply histogram binning to that. This has the advantage of using the average of the scaled probabilities in each bin instead of the average of the observed binary labels (see Figure 14.2). In [KLM19], they prove that this results in better calibration, due to the lower variance of the estimator.

In the multiclass case,  $\mathbf{z}$  is the vector of logits, and  $\mathbf{p} = \text{softmax}(\mathbf{z})$  is the vector of probabilities. We wish to convert this to a better calibrated version,  $\mathbf{q}$ . [ZE01b] propose to extend histogram binning and isotonic regression to this case by applying the above binary method to each of the  $K$  one-vs-rest problems, where  $K$  is the number of classes. However, this requires  $K$  separate calibration models, and results in an unnormalized probability distribution.

### 14.2.2.5 Temperature scaling

In [Guo+17], they noticed empirically that the diagonal version of Platt scaling, when applied to a variety of DNNs, often ended learning a vector of the form  $\mathbf{w} = (c, c, \dots, c)$ , for some constant  $c$ . This suggests a simpler form of scaling, which they call **temperature scaling**:  $\mathbf{q} = \text{softmax}(\mathbf{z}/T)$ ,

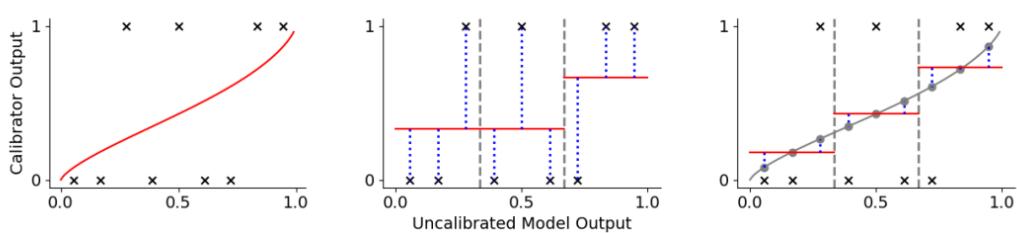


Figure 14.2: Visualization of 3 different approaches to calibrating a binary probabilistic classifier. Black crosses are the observed binary labels, red lines are the calibrated outputs. (a) Platt scaling. (b) Histogram binning with 3 bins. The output in each bin is the average of the binary labels in each bin. (c) The scaling-binning calibrator. This first applies Platt scaling, and then computes the average of the scaled points (gray circles) in each bin. From Figure 1 of [KLM19]. Used with kind permission of Ananya Kumar.

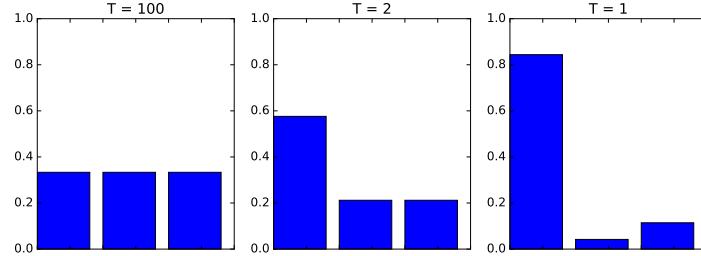


Figure 14.3: Softmax distribution  $\text{softmax}(\mathbf{a}/T)$ , where  $\mathbf{a} = (3, 0, 1)$ , at temperatures of  $T = 100$ ,  $T = 2$  and  $T = 1$ . When the temperature is high (left), the distribution is uniform, whereas when the temperature is low (right), the distribution is “spiky”, with most of its mass on the largest element. Generated by softmax\_plot.ipynb.

where  $T > 0$  is a temperature parameter, which can be estimated by maximum likelihood on the validation set. The effect of this temperature parameter is to make the distribution less peaky, as shown in Figure 14.3. [Guo+17] show empirically that this method produces the lowest ECE on a variety of DNN classification problems (see Figure 14.1 for a visualization). Furthermore, it is much simpler and faster than the other methods.

Note that Platt scaling and temperature scaling do not affect the identity of the most probable class label, so these methods have no impact on classification accuracy. However, they do improve calibration performance. A more recent multi-class calibration method is discussed in [Kul+19].

#### 14.2.2.6 Label smoothing

When training classifiers, we usually represent the true target label as a one-hot vector, say  $\mathbf{y} = (0, 1, 0)$  to represent class 2 out of 3. We can improve results if we “spread” some of the probability mass across all the bins. For example we may use  $\mathbf{y} = (0.1, 0.8, 0.1)$ . This is called **label smoothing** and

1 often results in better-calibrated models [MKH19].  
2

3  
4 **14.2.2.7 Bayesian methods**

5 Bayesian approaches to fitting classifiers often result in more calibrated predictions, since they  
6 represent uncertainty in the parameters. See Section 17.3.8 for an example. However, [Ova+19]  
7 shows that well-calibrated models (even Bayesian ones) often become mis-calibrated when applied to  
8 inputs that come from a different distribution (see Section 19.2 for details).  
9

10  
11 **14.2.3 Beyond evaluating marginal probabilities**

12 Calibration (Section 14.2.2) focuses on assessing properties of the marginal predictive distribution  
13  $p(y|\mathbf{x})$ . But this can sometimes be insufficient to distinguish between a good and bad model, especially  
14 in the context of online learning and sequential decision making, as pointed out in [Lu+22; Osb+21;  
15 WSG21; KKG22]. For example, consider two learning agents who observe a sequence of coin tosses.  
16 Let the outcome at time  $t$  be  $Y_t \sim \text{Ber}(\theta)$ , where  $\theta$  is the unknown parameter. Agent 1 believes  
17  $\theta = 2/3$ , whereas agent 2 believes either  $\theta = 0$  or  $\theta = 1$ , but is not sure which, and puts probabilities  
18  $1/3$  and  $2/3$  on these events. Thus both agents, despite having different models, make identical  
19 predictions for the next outcome:  $p(Y_1^i = 0) = 1/3$  for agents  $i = 1, 2$ . However, the predictions of  
20 the two agents about a *sequence* of  $\tau$  future outcomes is very different: In particular, agent 1 predicts  
21 each individual coin toss is a random Bernoulli event, where the probability is due to irreducible  
22 noise or **aleatoric uncertainty**:  
23

$$\begin{aligned} \text{p}(Y_1^1 = 0, \dots, Y_\tau^1 = 0) &= \frac{1}{3^\tau} \end{aligned} \tag{14.22}$$

26 By contrast, agent 2 predicts that the sequence will either be all heads or all tails, where the  
27 probability is induced by **epistemic uncertainty** about the true parameters:  
28

$$\begin{aligned} \text{p}(Y_1^2 = y_1, \dots, Y_\tau^2 = y_\tau) &= \begin{cases} 1/3 & \text{if } y_1 = \dots = y_\tau = 0 \\ 2/3 & \text{if } y_1 = \dots = y_\tau = 1 \\ 0 & \text{otherwise} \end{cases} \end{aligned} \tag{14.23}$$

33 The difference in beliefs between these agents will impact their behavior. For example, in a casino,  
34 agent 1 incurs little risk on repeatedly betting on heads in the long run, but for agent 2, this would  
35 be a very unwise strategy, and some initial information gathering (exploration) would be worthwhile.  
36

37 Based on the above, we see that it is useful to evaluate *joint* predictive distributions when assessing  
38 predictive models. In [Lu+22; Osb+21] they propose to evaluate the posterior predictive distributions  
39 over  $\tau$  outcomes  $\mathbf{y} = Y_{T+1:T+\tau}$ , given a set of  $\tau$  inputs  $\mathbf{x} = X_{T:T+\tau-1}$ , and the past  $T$  data samples,  
40  $\mathcal{D}_T = \{(X_t, Y_{t+1}) : t = 0, 1, \dots, T - 1\}$ . The Bayes optimal predictive distribution is

$$\begin{aligned} \text{P}_T^B &= p(\mathbf{y}|\mathbf{x}, \mathcal{D}_T) \end{aligned} \tag{14.24}$$

43 This is usually intractable to compute. Instead the agent will use an approximate distribution, known  
44 as a **belief state**, which we denote by  
45

$$\begin{aligned} Q_T &= p(\mathbf{y}|\mathbf{x}, \mathcal{D}_T) \end{aligned} \tag{14.25}$$

The natural performance metric is the KL between these distributions. Since this depend on the inputs  $\mathbf{x}$  and  $\mathcal{D}_T = (X_{0:T-1}, Y_{1:T})$ , we will averaged the KL over these values, which are drawn iid from the true data generating distribution, which we denote by

$$P(X, Y, \mathcal{E}) = P(X|\mathcal{E})P(Y|X, \mathcal{E})P(\mathcal{E}) \quad (14.26)$$

where  $\mathcal{E}$  is the true but unknown environment. Thus we define our metric as

$$d_{B,Q}^{KL} = \mathbb{E}_{P(\mathbf{x}, \mathcal{D}_T)} [D_{\text{KL}}(P^B(\mathbf{y}|\mathbf{x}, \mathcal{D}_T) \| Q(\mathbf{y}|\mathbf{x}, \mathcal{D}_T))] \quad (14.27)$$

where

$$P(\mathbf{x}, \mathcal{D}_T, \mathcal{E}) = P(\mathcal{E}) \underbrace{\left[ \prod_{t=0}^{T-1} P(X_t|\mathcal{E})P(Y_{t+1}|X_t, \mathcal{E}) \right]}_{P(\mathcal{D}_T|\mathcal{E})} \underbrace{\left[ \prod_{t=T}^{T+\tau-1} P(x_t|\mathcal{E}) \right]}_{P(\mathbf{x}|\mathcal{E})} \quad (14.28)$$

and  $P(\mathbf{x}, \mathcal{D}_T)$  marginalizes this over environments.

Unfortunately, it is usually intractable to compute the exact Bayes posterior,  $P_T^B$ , so we cannot evaluate  $d_{B,Q}^{KL}$ . However, in Section 14.2.3.1, we show that

$$d_{B,Q}^{KL} = d_{\mathcal{E},Q}^{KL} - \mathbb{I}(\mathcal{E}; \mathbf{y}|\mathcal{D}_T, \mathbf{x}) \quad (14.29)$$

where the second term is a constant wrt the agent, and the first term is given by

$$d_{\mathcal{E},Q}^{KL} = \mathbb{E}_{P(\mathbf{x}, \mathcal{D}_T, \mathcal{E})} [D_{\text{KL}}(P(\mathbf{y}|\mathbf{x}, \mathcal{E}) \| Q(\mathbf{y}|\mathbf{x}, \mathcal{D}_T))] \quad (14.30)$$

$$= \mathbb{E}_{P(\mathbf{y}|\mathbf{x}, \mathcal{E})P(\mathbf{x}, \mathcal{D}_T, \mathcal{E})} \left[ \log \frac{P(\mathbf{y}|\mathbf{x}, \mathcal{E})}{Q(\mathbf{y}|\mathbf{x}, \mathcal{D}_T)} \right] \quad (14.31)$$

Hence if we rank agents in terms of  $d_{\mathcal{E},Q}^{KL}$ , it will give the same results as ranking them by  $d_{B,Q}^{KL}$ .

To compute  $d_{\mathcal{E},Q}^{KL}$  in practice, we can use a Monte Carlo approximation: we just have to sample  $J$  environments,  $\mathcal{E}^j \sim P(\mathcal{E})$ , sample a training set  $\mathcal{D}_T$  from each environment,  $\mathcal{D}_T^j \sim P(\mathcal{D}_T|\mathcal{E}^j)$ , and then sample  $N$  data vectors of length  $\tau$ ,  $(\mathbf{x}_n^j, \mathbf{y}_n^j) \sim P(X_{T:T+\tau-1}, Y_{T+1:T+\tau}|\mathcal{E}^j)$ . We can then compute

$$\hat{d}_{\mathcal{E},Q}^{KL} = \frac{1}{JN} \sum_{j=1}^J \sum_{n=1}^N \left[ \log P(\mathbf{y}_n^j|\mathbf{x}_n^j, \mathcal{E}^j) - \log Q(\mathbf{y}_n^j|\mathbf{x}_n^j, \mathcal{D}_T^j) \right] \quad (14.32)$$

where

$$p_{jn} = P(\mathbf{y}_n^j|\mathbf{x}_n^j, \mathcal{E}^j) = \prod_{t=T}^{T+\tau-1} P(Y_{n,t+1}^j|X_{n,t}^j, \mathcal{E}^j) \quad (14.33)$$

$$q_{jn} = Q(\mathbf{y}_n^j|\mathbf{x}_n^j, \mathcal{D}_T^j) = \int Q(\mathbf{y}_n^j|\mathbf{x}_n^j, \boldsymbol{\theta})Q(\boldsymbol{\theta}|\mathcal{D}_T^j)d\boldsymbol{\theta} \quad (14.34)$$

$$\approx \frac{1}{M} \sum_{m=1}^M \prod_{t=T}^{T+\tau-1} Q(Y_{n,t+1}^j|X_{n,t}^j, \boldsymbol{\theta}_m^j) \quad (14.35)$$

1 where  $\theta_m^j \sim Q(\theta|\mathcal{D}_T^j)$  is a sample from the agent's posterior over the environment.  
2

3 The above assumes that  $P(Y|X)$  is known; this will be the case if we use a synthetic data generator,  
4 as in the the “neural testbed” in [Osb+21]. If we just have an  $J$  empirical distributions for  $P^j(X, Y)$ ,  
5 we can replace the KL with the cross entropy, which only differs by an additive constant:

$$\underline{d}_{\mathcal{E},Q}^{KL} = \mathbb{E}_{P(\mathbf{x}, \mathcal{D}_T, \mathcal{E})} [D_{KL}(P(\mathbf{y}|\mathbf{x}, \mathcal{E}) \| Q(\mathbf{y}|\mathbf{x}, \mathcal{D}_T))] \quad (14.36)$$

$$= \underbrace{\mathbb{E}_{P(\mathbf{x}, \mathbf{y}, \mathcal{E})} [\log P(\mathbf{y}|\mathbf{x}, \mathcal{E})]}_{\text{const}} - \underbrace{\mathbb{E}_{P(\mathbf{x}, \mathbf{y}, \mathcal{D}_T|\mathcal{E})P(\mathcal{E})} [\log Q(\mathbf{y}|\mathbf{x}, \mathcal{D}_T)]}_{d_{\mathcal{E},Q}^{CE}} \quad (14.37)$$

11 where the latter term is just the empirical negative log likelihood (NLL) of the agent on samples  
12 from the environment. Hence if we rank agents in terms of their NLL or cross entropy  $d_{\mathcal{E},Q}^{CE}$  we will  
13 get the same results as ranking them by  $d_{\mathcal{E},Q}^{KL}$ , which will in turn give the same results as ranking  
14 them by  $d_{B,Q}^{KL}$ .

15 In practice we can approximate the cross entropy as follows:

$$\hat{d}_{\mathcal{E},Q}^{CE} = -\frac{1}{JN} \sum_{j=1}^J \sum_{n=1}^N \log Q(\mathbf{y}_n^j | \mathbf{x}_n^j, \mathcal{D}_T^j) \quad (14.38)$$

20 where  $\mathcal{D}_T^j \sim P^j$ , and  $(\mathbf{x}_n^j, \mathbf{y}_n^j) \sim P^j$ .

21 An alternative to estimating the KL or NLL is to evaluate the joint predictive accuracy by using it  
22 in a downstream task. In [Osb+21], they show that good predictive accuracy (for  $\tau > 1$ ) correlates  
23 with good performance on a bandit problem (see Section 34.4). In [WSG21] they show that good  
24 predictive accuracy (for  $\tau > 1$ ) results in good performance on a transductive active learning task.  
25

#### 26 14.2.3.1 Proof of claim

28 We now prove Equation (14.29), based on [Lu+21a]. First note that

$$\underline{d}_{\mathcal{E},Q}^{KL} = \mathbb{E}_{P(\mathbf{x}, \mathcal{D}_T, \mathcal{E})P(\mathbf{y}|\mathbf{x}, \mathcal{E})} \left[ \log \frac{P(\mathbf{y}|\mathbf{x}, \mathcal{E})}{Q(\mathbf{y}|\mathbf{x}, \mathcal{D}_T)} \right] \quad (14.39)$$

$$= \mathbb{E} \left[ \log \frac{P(\mathbf{y}|\mathbf{x}, \mathcal{D}_T)}{Q(\mathbf{y}|\mathbf{x}, \mathcal{D}_T)} \right] + \mathbb{E} \left[ \log \frac{P(\mathbf{y}|\mathbf{x}, \mathcal{E})}{P(\mathbf{y}|\mathbf{x}, \mathcal{D}_T)} \right] \quad (14.40)$$

34 For the first term in Equation (14.40) we have

$$\mathbb{E} \left[ \log \frac{P(\mathbf{y}|\mathbf{x}, \mathcal{D}_T)}{Q(\mathbf{y}|\mathbf{x}, \mathcal{D}_T)} \right] = \sum P(\mathbf{x}, \mathbf{y}, \mathcal{D}_T) \log \frac{P(\mathbf{y}|\mathbf{x}, \mathcal{D}_T)}{Q(\mathbf{y}|\mathbf{x}, \mathcal{D}_T)} \quad (14.41)$$

$$= \sum P(\mathbf{x}, \mathcal{D}_T) \sum P(\mathbf{y}|\mathbf{x}, \mathcal{D}_T) \log \frac{P(\mathbf{y}|\mathbf{x}, \mathcal{D}_T)}{Q(\mathbf{y}|\mathbf{x}, \mathcal{D}_T)} \quad (14.42)$$

$$= \mathbb{E}_{P(\mathbf{x}, \mathcal{D}_T)} [D_{KL}(P(\mathbf{y}|\mathbf{x}, \mathcal{D}_T) \| Q(\mathbf{y}|\mathbf{x}, \mathcal{D}_T))] = d_{B,Q}^{KL} \quad (14.43)$$

42 We now show that the second term in Equation (14.40) reduces to the mutual information. We  
43 exploit the fact that

$$P(\mathbf{y}|\mathbf{x}, \mathcal{E}) = P(\mathbf{y}|\mathcal{D}_T, \mathbf{x}, \mathcal{E}) = \frac{P(\mathcal{E}, \mathbf{y}|\mathcal{D}_T, \mathbf{x})}{P(\mathcal{E}|\mathcal{D}_T, \mathbf{x})} \quad (14.44)$$



Figure 14.4: Prediction set examples on Imagenet. We show three progressively more difficult examples of the class fox squirrel and the prediction sets generated by conformal prediction. From Figure 1 of [AB21]. Used with kind permission of Anastasios Angelopoulos.

since  $\mathcal{D}_T$  has no new information in beyond  $\mathcal{E}$ . From this we get

$$\mathbb{E} \left[ \log \frac{P(\mathbf{y}|\mathbf{x}, \mathcal{E})}{P(\mathbf{y}|\mathbf{x}, \mathcal{D}_T)} \right] = \mathbb{E} \left[ \log \frac{P(\mathcal{E}, \mathbf{y}|\mathcal{D}_T, \mathbf{x})/P(\mathcal{E}|\mathcal{D}_T, \mathbf{x})}{P(\mathbf{y}|\mathcal{D}, \mathcal{D}_T)} \right] \quad (14.45)$$

$$= \sum P(\mathcal{D}_T, \mathbf{x}) \sum P(\mathcal{E}, \mathbf{y}|\mathcal{D}_T, \mathbf{x}) \log \frac{P(\mathcal{E}, \mathbf{y}|\mathcal{D}_T, \mathbf{x})}{P(\mathbf{y}|\mathcal{D}_T, \mathbf{x})P(\mathcal{E}|\mathcal{D}_T, \mathbf{x})} \quad (14.46)$$

$$= \mathbb{I}(\mathcal{E}; \mathbf{y}|\mathcal{D}_T, \mathbf{x}) \quad (14.47)$$

Hence

$$d_{\mathcal{E}, Q}^{KL} = d_{B, Q}^{KL} + \mathbb{I}(\mathcal{E}; \mathbf{y}|\mathcal{D}_T, \mathbf{x}) \quad (14.48)$$

as claimed.

### 14.3 Conformal prediction

In this section, we briefly discuss **conformal prediction** [VGS05; SV08; ZFV20; AB21; KSB21; Man22]. This is a simple but effective way to create prediction intervals or sets with guaranteed frequentist coverage probability from any predictive method  $p(y|\mathbf{x})$ . This can be seen as a form of **distribution free uncertainty quantification**, since it works without making assumptions (beyond exchangeability of the data) about the true data generating process or the form of the model.<sup>1</sup> Our presentation is based on the excellent tutorial of [AB21].<sup>2</sup>

In conformal prediction, we start with some heuristic notion of uncertainty — such as the softmax score for a classification problem, or the variance for a regression problem — and we use it to define a **conformal score**  $s(\mathbf{x}, y) \in \mathbb{R}$ , which measures how badly the output  $y$  “conforms” to  $\mathbf{x}$ . (Large

1. The exchangeability assumption rules out time series data, which is serially correlated. However, extensions to conformal prediction have been developed for the time series case, see e.g., [Zaf+22]. The exchangeability assumption also rules out distribution shift, although this has also been partially addressed.

2. See also the easy-to-use **MAPIE** Python library at <https://mapie.readthedocs.io/en/latest/index.html>, and the list of papers at <https://github.com/valeman/awesome-conformal-prediction>.

values of the score are less likely, so it is better to think of it as a non-conformity score.) Next we apply this score to a **calibration set** of  $n$  labeled examples, that was not used to train  $f$ , to get  $\mathcal{S} = \{s_i = s(\mathbf{x}_i, y_i) : i = 1 : n\}$ .<sup>3</sup> The user specifies a desired confidence threshold  $\alpha$ , say 0.1, and we then compute the  $(1 - \alpha)$  quantile  $\hat{q}$  of  $\mathcal{S}$ . (In fact, we should replace  $1 - \alpha$  with  $\frac{\lceil(n+1)(1-\alpha)\rceil}{n}$ , to account for the finite size of  $\mathcal{S}$ .) Finally, given a new test input,  $\mathbf{x}_{n+1}$ , we compute the prediction set to be

$$\mathcal{T}(\mathbf{x}_{n+1}) = \{y : s(\mathbf{x}_{n+1}, y) \leq \hat{q}\} \quad (14.49)$$

Intuitively, we include all the outputs  $y$  that are plausible given the input. See Figure 14.4 for an illustration.

Remarkably, one can show the following general result

$$1 - \alpha \leq P^*(y^{n+1} \in \mathcal{T}(\mathbf{x}_{n+1})) \leq 1 - \alpha + \frac{1}{n+1} \quad (14.50)$$

where the probability is wrt the true distribution  $P^*(\mathbf{x}_{n+1}, y_{n+1})$ . We say that the prediction set has a **coverage** level of  $1 - \alpha$ . This holds for any value of  $n \geq 1$  and  $\alpha \in [0, 1]$ . The only assumption is that the values  $(\mathbf{x}_i, y_i)$  are exchangeable, and hence the calibration scores  $s_i$  are also exchangeable.

To see why this is true, let us sort the scores so  $s_1 < \dots < s_n$ , so  $\hat{q} = s_i$ , where  $i = \frac{\lceil(n+1)(1-\alpha)\rceil}{n}$ . (We assume the scores are distinct, for simplicity.) The score  $s_{n+1}$  is equally likely to fall in anywhere between the calibration points  $s_1, \dots, s_n$ , since the points are exchangeable. Hence

$$P^*(s_{n+1} \leq s_k) = \frac{k}{n+1} \quad (14.51)$$

for any  $k \in \{1, \dots, n+1\}$ . The event  $\{y_{n+1} \in \mathcal{T}(\mathbf{x}_{n+1})\}$  is equivalent to  $\{s_{n+1} \leq \hat{q}\}$ . Hence

$$P^*(y_{n+1} \in \mathcal{T}(\mathbf{x}_{n+1})) = P^*(s_{n+1} \leq \hat{q}) = \frac{\lceil(n+1)(1-\alpha)\rceil}{n+1} \geq 1 - \alpha \quad (14.52)$$

For the proof of the upper bound, see [Lei+18].

Although this result may seem like a “free lunch”, it is worth noting that we can always achieve a desired coverage level by defining the prediction set to be all possible labels. In this case, the prediction set will be independent of the input, but it will cover the true label  $1 - \alpha$  of the time. To rule out some degenerate cases, we seek prediction sets that are as small as possible (although we allow for the set to be larger for harder examples), while meeting the coverage requirement. Achieving this goal requires that we define suitable conformal scores. Below we give some examples of how to compute conformal scores  $s(\mathbf{x}, y)$  for different kinds of problem.<sup>4</sup>

37

### 14.3.1 Conformalizing classification

The simplest way to apply conformal prediction to multiclass classification is to derive the conformal score from the softmax score assigned to the label using  $s(\mathbf{x}, y) = 1 - f(\mathbf{x})_y$ , so large values are

43 3. Using a calibration set is called **split conformal prediction**. If we don't have enough data to adopt this splitting approach, we can use **full conformal prediction** [VGS05], which requires fitting the model  $n$  times using a leave-one-out type procedure.

44 4. It is also possible to learn conformal scores in an end-to-end way, jointly with the predictive model, as discussed in [Stu+22].

47

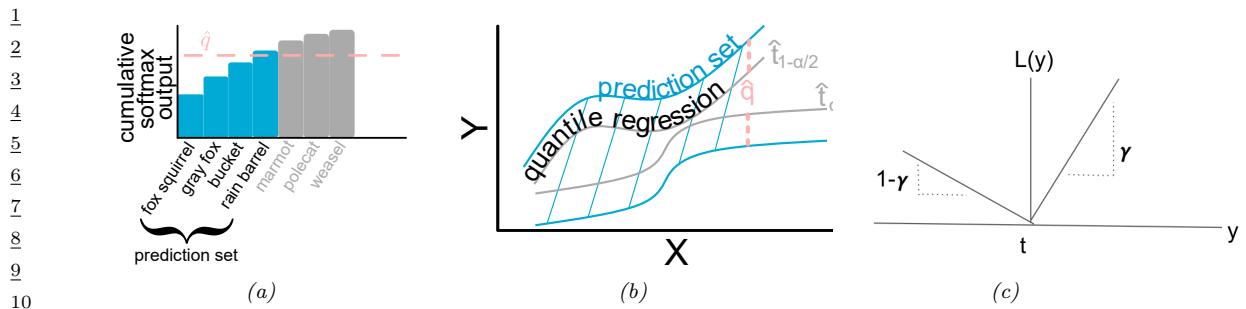


Figure 14.5: (a) Illustration of adaptive prediction set. From Figure 5 of [AB21]. Used with kind permission of Anastasios Angelopoulos. (b) Illustrate of conformalized quantile regression. From Figure 6 of [AB21]. Used with kind permission of Anastasios Angelopoulos. (c) Illustration of pinball loss function.

considered less likely than small values. We compute the threshold  $\hat{q}$  as described above, and then we define the prediction set to be  $\mathcal{T}(\mathbf{x}) = \{y : f(\mathbf{x})_y \geq 1 - \hat{q}\}$ , which matches Equation (14.49). That is, we take the set of all class labels above the specified threshold, as illustrated in Figure 14.4.

Although the above approach produces prediction sets with the smallest average size (as proved in [SLW19]), the size of the set tends to be too large for easy examples and too small for hard examples. We now present an improved method, known as **adaptive prediction sets**, due to [RSC20], which solves this problem. The idea is simple: we sort all the softmax scores,  $f(\mathbf{x})_c$  for  $c = 1 : C$ , to get permutation  $\pi_{1:C}$ , and then we define  $s(\mathbf{x}, y)$  to be the cumulative sum of the scores up until we reach label  $y$ :  $s(\mathbf{x}, y) = \sum_{c=1}^k f(\mathbf{x})_{\pi_c}$ , where  $k = \pi_y$ . We now compute  $\hat{q}$  as before, and define the prediction set  $\mathcal{T}(\mathbf{x})$  to be the set of all labels, sorted in order of decreasing probability, until we cover  $\hat{q}$  of the probability mass. See Figure 14.5a for an illustration. This uses all the softmax scores output by the model, rather than just the top score, which accounts for its improved performance.

### 14.3.2 Conformalizing regression

In this section, we consider conformalized regression problems. Since now  $y \in \mathbb{R}$ , computing the prediction set in Equation (14.49) is expensive, so instead we will compute a prediction interval, specified by a lower and upper bound.

#### 14.3.2.1 Conformalizing quantile regression

In this section, we use **quantile regression** to compute the lower and upper bounds. We first fit a function of the form  $t_\gamma(\mathbf{x})$ , which predicts the  $\gamma$  quantile of the pdf  $P(Y|\mathbf{x})$ . For example, if we set  $\gamma = 0.5$ , we get the median. If we use  $\gamma = 0.05$  and  $\gamma = 0.95$ , we can get an approximate 90% prediction interval using  $[t_{0.05}(\mathbf{x}), t_{0.95}(\mathbf{x})]$ , as illustrated by the gray lines in Figure 14.5b. To fit the quantile regression model, we just replace squared loss with the **quantile loss**, also called the **pinball loss**, which is defined as

$$\ell_\gamma(y, \hat{t}) = (y - \hat{t})\gamma\mathbb{I}(y > \hat{t}) + (\hat{t} - y)(1 - \gamma)\mathbb{I}(y < \hat{t}) \quad (14.53)$$

where  $y$  is the true output and  $\hat{t}$  is the predicted value at quantile  $\gamma$ . See Figure 14.5c for an illustration.

1 The regression quantiles are only approximately a 90% interval because the model may be  
2 mismatched to the true distribution. However we can use conformal prediction to fix this. In  
3 particular, let us define the conformal score to be  
4

5  $s(\mathbf{x}, y) = \max(\hat{t}_{\alpha/2}(\mathbf{x}) - y, y - \hat{t}_{\alpha/2}(\mathbf{x}))$  (14.54)  
6

7 In other words,  $s(\mathbf{x}, y)$  is a positive measure of how far the value  $y$  is outside the prediction interval,  
8 or is a negative measure if  $y$  is inside the prediction interval. We compute  $\hat{q}$  as before, and define the  
9 conformal prediction interval to be  
10

11  $\mathcal{T}(\mathbf{x}) = [\hat{t}_{\alpha/2}(\mathbf{x}) - \hat{q}, \hat{t}_{\alpha/2}(\mathbf{x}) + \hat{q}]$  (14.55)  
12

13 This makes the quantile regression interval wider if  $\hat{q}$  is positive (if the base method was overconfident),  
14 and narrower if  $\hat{q}$  is negative (if the base method was underconfident). See Figure 14.5b for an  
15 illustration. This approach is called **conformalized quantile regression** or **CQR** [RPC19].  
16

#### 17 14.3.2.2 Conformalizing predicted variances

18 There are many ways to define uncertainty scores  $u(\mathbf{x})$ , such as the predicted standard deviation,  
19 from which we can derive a prediction interval using  
20

21  $\mathcal{T}(\mathbf{x}) = [f(\mathbf{x}) - u(\mathbf{x})\hat{q}, f(\mathbf{x}) + u(\mathbf{x})\hat{q}]$  (14.56)  
22

23 Here  $\hat{q}$  is derived from the quantiles of the following conformal scores  
24

25  $s(\mathbf{x}, y) = \frac{|y - f(\mathbf{x})|}{u(\mathbf{x})}$  (14.57)  
26

27 The interval produced by this method tends to be wider than the one computed by CQR, since it  
28 extends an equal amount above and below the predicted value  $f(\mathbf{x})$ . In addition, the uncertainty  
29 measure  $u(\mathbf{x})$  may not scale properly with  $\alpha$ . Nevertheless, this is a simple post-hoc method that  
30 can be applied to many regression methods without needing to retrain them.  
31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

# 15 Generalized linear models

## 15.1 Introduction

A **generalized linear model** or **GLM** [MN89] is a conditional version of an exponential family distribution (Section 2.3). More precisely, the model has the following form:

$$p(y_n | \mathbf{x}_n, \mathbf{w}, \sigma^2) = \exp \left[ \frac{y_n \eta_n - A(\eta_n)}{\sigma^2} + \log h(y_n, \sigma^2) \right] \quad (15.1)$$

where  $\eta_n = \mathbf{w}^\top \mathbf{x}_n$  is the natural parameter for the distribution,  $A(\eta_n)$  is the log normalizer,  $\mathcal{T}(y) = y$  is the sufficient statistic, and  $\sigma^2$  is the dispersion term. Based on the results in Section 2.3.3, we can show that the mean and variance of the response variable are as follows:

$$\mu_n \triangleq \mathbb{E}[y_n | \mathbf{x}_n, \mathbf{w}, \sigma^2] = A'(\eta_n) \triangleq \ell^{-1}(\eta_n) \quad (15.2)$$

$$\mathbb{V}[y_n | \mathbf{x}_n, \mathbf{w}, \sigma^2] = A''(\eta_n) \sigma^2 \quad (15.3)$$

We will denote the mapping from the linear inputs to the mean of the output using  $\mu_n = \ell^{-1}(\eta_n)$ , where the function  $\ell$  is known as the **link function**, and  $\ell^{-1}$  is known as the **mean function**. This relationship is usually written as follows:

$$\ell(\mu_n) = \eta_n = \mathbf{w}^\top \mathbf{x}_n \quad (15.4)$$

### 15.1.1 Examples

In this section, we give some examples of widely used GLMs.

#### 15.1.1.1 Linear regression

Recall that linear regression has the form

$$p(y_n | \mathbf{x}_n, \mathbf{w}, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left(-\frac{1}{2\sigma^2}(y_n - \mathbf{w}^\top \mathbf{x}_n)^2\right) \quad (15.5)$$

Hence

$$\log p(y_n | \mathbf{x}_n, \mathbf{w}, \sigma^2) = -\frac{1}{2\sigma^2}(y_n - \eta_n)^2 - \frac{1}{2} \log(2\pi\sigma^2) \quad (15.6)$$

1 where  $\eta_n = \mathbf{w}^\top \mathbf{x}_n$ . We can write this in GLM form as follows:

3

$$\log p(y_n | \mathbf{x}_n, \mathbf{w}, \sigma^2) = \frac{y_n \eta_n - \frac{\eta_n^2}{2}}{\sigma^2} - \frac{1}{2} \left( \frac{y_n^2}{\sigma^2} + \log(2\pi\sigma^2) \right) \quad (15.7)$$

4

5 We see that  $A(\eta_n) = \eta_n^2/2$  and hence

6

$$\mathbb{E}[y_n] = \eta_n = \mathbf{w}^\top \mathbf{x}_n \quad (15.8)$$

7

8

$$\mathbb{V}[y_n] = \sigma^2 \quad (15.9)$$

9

10 See Section 15.2 for details on linear regression.

### 11 15.1.1.2 Binomial regression

12 If the response variable is the number of successes in  $N_n$  trials,  $y_n \in \{0, \dots, N_n\}$ , we can use  
13 **binomial regression**, which is defined by

14

$$p(y_n | \mathbf{x}_n, N_n, \mathbf{w}) = \text{Bin}(y_n | \sigma(\mathbf{w}^\top \mathbf{x}_n), N_n) \quad (15.10)$$

15

16 We see that binary logistic regression is the special case when  $N_n = 1$ .

17 The log pdf is given by

18

$$\log p(y_n | \mathbf{x}_n, N_n, \mathbf{w}) = y_n \log \mu_n + (N_n - y_n) \log(1 - \mu_n) + \log \binom{N_n}{y_n} \quad (15.11)$$

19

20

$$= y_n \log\left(\frac{\mu_n}{1 - \mu_n}\right) + N_n \log(1 - \mu_n) + \log \binom{N_n}{y_n} \quad (15.12)$$

21

22 where  $\mu_n = \sigma(\eta_n)$ . To rewrite this in GLM form, let us define

23

$$\eta_n \triangleq \log \left[ \frac{\mu_n}{(1 - \mu_n)} \right] = \log \left[ \frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}_n}} \frac{1 + e^{-\mathbf{w}^\top \mathbf{x}_n}}{e^{-\mathbf{w}^\top \mathbf{x}_n}} \right] = \log \frac{1}{e^{-\mathbf{w}^\top \mathbf{x}_n}} = \mathbf{w}^\top \mathbf{x}_n \quad (15.13)$$

24

25 Hence we can write binomial regression in GLM form as follows

26

$$\log p(y_n | \mathbf{x}_n, N_n, \mathbf{w}) = y_n \eta_n - A(\eta_n) + h(y_n) \quad (15.14)$$

27

28 where  $h(y_n) = \log \binom{N_n}{y_n}$  and

29

$$A(\eta_n) = -N_n \log(1 - \mu_n) = N_n \log(1 + e^{\eta_n}) \quad (15.15)$$

30

31 Hence

32

$$\mathbb{E}[y_n] = \frac{dA}{d\eta_n} = \frac{N_n e^{\eta_n}}{1 + e^{\eta_n}} = \frac{N_n}{1 + e^{-\eta_n}} = N_n \mu_n \quad (15.16)$$

33

34 and

35

$$\mathbb{V}[y_n] = \frac{d^2 A}{d\eta_n^2} = N_n \mu_n (1 - \mu_n) \quad (15.17)$$

36

37 See Supplementary Section 15.3 for an example of binomial regression.

38

---

### 15.1.1.3 Poisson regression

3 If the response variable is an integer count,  $y_n \in \{0, 1, \dots\}$ , we can use **Poisson regression**, which  
 4 is defined by

$$6 p(y_n | \mathbf{x}_n, \mathbf{w}) = \text{Poi}(y_n | \exp(\mathbf{w}^\top \mathbf{x}_n)) \quad (15.18)$$

7 where

$$9 \text{Poi}(y | \mu) = e^{-\mu} \frac{\mu^y}{y!} \quad (15.19)$$

11 is the Poisson distribution. Poisson regression is widely used in bio-statistical applications, where  
 12  $y_n$  might represent the number of diseases of a given person or place, or the number of reads at a  
 13 genomic location in a high-throughput sequencing context (see e.g., [Kua+09]).

14 The log pdf is given by

$$16 \log p(y_n | \mathbf{x}_n, \mathbf{w}) = y_n \log \mu_n - \mu_n - \log(y_n!) \quad (15.20)$$

17 where  $\mu_n = \exp(\mathbf{w}^\top \mathbf{x}_n)$ . Hence in GLM form we have

$$19 \log p(y_n | \mathbf{x}_n, \mathbf{w}) = y_n \eta_n - A(\eta_n) + h(y_n) \quad (15.21)$$

20 where  $\eta_n = \log(\mu_n) = \mathbf{w}^\top \mathbf{x}_n$ ,  $A(\eta_n) = \mu_n = e^{\eta_n}$ , and  $h(y_n) = -\log(y_n!)$ . Hence

$$22 \mathbb{E}[y_n] = \frac{dA}{d\eta_n} = e^{\eta_n} = \mu_n \quad (15.22)$$

24 and

$$26 \mathbb{V}[y_n] = \frac{d^2A}{d\eta_n^2} = e^{2\eta_n} = \mu_n \quad (15.23)$$

### 29 15.1.1.4 Zero-inflated Poisson regression

30 In many forms of count data, the number of observed 0s is larger than what a model might expect,  
 31 even after taking into account the predictors. Intuitively, this is because there may be many ways  
 32 to produce no outcome. For example, consider predicting sales data for a product. If the sales are  
 33 0, does it mean the product is unpopular (so the demand is very low), or was it simply sold out  
 34 (implying the demand is high, but exceed supply)? Similar problems arise in genomics, epidemiology,  
 35 etc.

36 To handle such situations, it is common to use a **zero-inflated Poisson** or **ZIP** model. The  
 37 likelihood for this model is a mixture of two distributions: a spike at 0, and a standard Poisson.  
 38 Formally, we define

$$40 \text{ZIP}(y | \rho, \lambda) = \begin{cases} \rho + (1 - \rho) \exp(-\lambda) & \text{if } y = 0 \\ (1 - \rho) \frac{\lambda^y \exp(-\lambda)}{y!} & \text{if } y > 0 \end{cases} \quad (15.24)$$

43 Here  $\rho$  is the prior probability of picking the spike, and  $\lambda$  is the rate of the Poisson. We see that  
 44 there are two “mechanisms” for generating a 0: either (with probability  $\rho$ ) we chose the spike, or  
 45 (with probability  $1 - \rho$ ) we simply generate a zero count just because the rate of the Poisson is so  
 46 low. (This latter event has probability  $\lambda^0 e^{-\lambda} / 0! = e^{-\lambda}$ .)

---

1 **15.1.2 GLMs with non-canonical link functions**

3 We have seen how the mean parameters of the output distribution are given by  $\mu = \ell^{-1}(\eta)$ , where the  
4 function  $\ell$  is the link function. There are several choices for this function, as we now discuss.

5 The **canonical link function**  $\ell$  satisfies the property that  $\theta = \ell(\mu)$ , where  $\theta$  are the canonical  
6 (natural) parameters. Hence

8 
$$\theta = \ell(\mu) = \ell(\ell^{-1}(\eta)) = \eta \tag{15.25}$$

10 This is what we have assumed so far. For example, for the Bernoulli distribution, the canonical  
11 parameter is the log-odds  $\eta = \log(\mu/(1 - \mu))$ , which is given by the logit transform

12 
$$\eta = \ell(\mu) = \text{logit}(\mu) = \log\left(\frac{\mu}{1 - \mu}\right) \tag{15.26}$$

15 The inverse of this is the sigmoid or logistic function

17 
$$\mu = \ell^{-1}(\eta) = \sigma(\eta) = 1/(1 + e^{-\eta}) \tag{15.27}$$

19 However, we are free to use other kinds of link function. For example, in Section 15.4 we use

21 
$$\eta = \ell(\mu) = \Phi^{-1}(\mu) \tag{15.28}$$

22 
$$\mu = \ell^{-1}(\eta) = \Phi(\eta) \tag{15.29}$$

24 This is known as the **probit link function**.

25 Another link function that is sometimes used for binary responses is the **complementary log-log**  
26 function

27 
$$\eta = \ell(\mu) = \log(-\log(1 - \mu)) \tag{15.30}$$

29 This is used in applications where we either observe 0 events (denoted by  $y = 0$ ) or one or more  
30 (denoted by  $y = 1$ ), where events are assumed to be governed by a Poisson distribution with rate  $\lambda$ .  
31 Let  $E$  be the number of events. The Poisson assumption means  $p(E = 0) = \exp(-\lambda)$  and hence

33 
$$p(y = 0) = (1 - \mu) = p(E = 0) = \exp(-\lambda) \tag{15.31}$$

35 Thus  $\lambda = -\log(1 - \mu)$ . When  $\lambda$  is a function of covariates, we need to ensure it is positive, so we use  
36  $\lambda = e^\eta$ , and hence

38 
$$\eta = \log(\lambda) = \log(-\log(1 - \mu)) \tag{15.32}$$

40 **15.1.3 Maximum likelihood estimation**

41 GLMs can be fit using similar methods to those that we used to fit logistic regression. In particular,  
42 the negative log-likelihood has the following form (ignoring constant terms):

44 
$$\text{NLL}(\mathbf{w}) = -\log p(\mathcal{D}|\mathbf{w}) = -\frac{1}{\sigma^2} \sum_{n=1}^N \ell_n \tag{15.33}$$

1 where

2

$$\ell_n \triangleq \eta_n y_n - A(\eta_n) \quad (15.34)$$

3 where  $\eta_n = \mathbf{w}^\top \mathbf{x}_n$ . For notational simplicity, we will assume  $\sigma^2 = 1$ .

4 We can compute the gradient for a single term as follows:

5

$$\mathbf{g}_n \triangleq \frac{\partial \ell_n}{\partial \mathbf{w}} = \frac{\partial \ell_n}{\partial \eta_n} \frac{\partial \eta_n}{\partial \mathbf{w}} = (y_n - A'(\eta_n)) \mathbf{x}_n = (y_n - \mu_n) \mathbf{x}_n \quad (15.35)$$

6 where  $\mu_n = f(\mathbf{w}^\top \mathbf{x}_n)$ , and  $f$  is the inverse link function that maps from canonical parameters to  
7 mean parameters. (For example, in the case of logistic regression, we have  $\mu_n = \sigma(\mathbf{w}^\top \mathbf{x})$ .)

8 The Hessian is given by

9

$$\mathbf{H} = \frac{\partial^2}{\partial \mathbf{w} \partial \mathbf{w}^\top} \text{NLL}(\mathbf{w}) = - \sum_{n=1}^N \frac{\partial \mathbf{g}_n}{\partial \mathbf{w}^\top} \quad (15.36)$$

10 where

11

$$\frac{\partial \mathbf{g}_n}{\partial \mathbf{w}^\top} = \frac{\partial \mathbf{g}_n}{\partial \mu_n} \frac{\partial \mu_n}{\partial \mathbf{w}^\top} = -\mathbf{x}_n f'(\mathbf{w}^\top \mathbf{x}_n) \mathbf{x}_n^\top \quad (15.37)$$

12 Hence

13

$$\mathbf{H} = \sum_{n=1}^N f'(\eta_n) \mathbf{x}_n \mathbf{x}_n^\top \quad (15.38)$$

14 For example, in the case of logistic regression,  $f(\eta_n) = \sigma(\eta_n) = \mu_n$ , and  $f'(\eta_n) = \mu_n(1 - \mu_n)$ . In  
15 general, we see that the Hessian is positive definite, since  $f'(\eta_n) > 0$ ; hence the negative log likelihood  
16 is convex, so the MLE for a GLM is unique (assuming  $f(\eta_n) > 0$  for all  $n$ ).

17 For small datasets, we can use the **iteratively reweighted least squares** or **IRLS** algorithm,  
18 which is a form Newton's method, to compute the MLE (see e.g., [Mur22, Sec 10.2.6]). For  
19 large datasets, we can use SGD. (In practice it is often useful to combine SGD with methods that  
20 automatically tune the step size, such as [Loi+21].)

### 21 15.1.4 Bayesian inference

22 To perform Bayesian inference of the parameters, we first need to specify a prior. Choosing a suitable  
23 prior depends on the form of link function. For example, a “flat” or “uninformative” prior on the  
24 offset term  $\alpha \in \mathbb{R}$  will not translate to an uninformative prior on the probability scale if we pass  $\alpha$   
25 through a sigmoid, as we discuss in Section 15.3.3.

26 Once we have chosen the prior, we can compute the posterior using a variety of approximate  
27 inference methods. For small datasets, HMC (Section 12.5) is the easiest to use, since you just need  
28 to write down the log likelihood and log prior, and use autograd to compute derivatives and pass  
29 them to the HMC engine. We give some examples in the following sections. For large datasets,  
30 assumed density filtering (Section 8.10.3) stochastic variational inference (Section 10.3.1), or more  
31 specialized algorithms (e.g., [HAB17]) are the best choice.

1 **15.2 Linear regression**

3 **Linear regression** is the simplest case of a GLM. We gave a detailed introduction to this model in  
4 the prequel to this book, [Mur22]. In this section, we discuss this model from a Bayesian perspective.  
5

6 **15.2.1 Conjugate priors**

8 In this section, we derive the posterior for the parameters of the model when we use a conjugate prior.  
9 We first consider the case where just  $\mathbf{w}$  is unknown (so the observation noise variance parameter  $\sigma^2$   
10 is fixed), and then we consider the general case, where both  $\sigma^2$  and  $\mathbf{w}$  are unknown.  
11

12 **15.2.1.1 Noise variance is known**

13 The conjugate prior for linear regression has the following form:

$$\underline{16} \quad p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \tilde{\mathbf{w}}, \breve{\Sigma}) \quad (15.39)$$

17 We often use  $\tilde{\mathbf{w}} = \mathbf{0}$  as the prior mean and  $\breve{\Sigma} = \tau^2 \mathbf{I}_D$  as the prior covariance. (We assume the bias  
18 term is included in the weight vector, but often use a much weaker prior for it, since we typically do  
19 not want to regularize the overall mean level of the output.)  
20

21 To derive the posterior, let us first rewrite the likelihood in terms of an MVN as follows:

$$\underline{22} \quad \ell(\mathbf{w}) = p(\mathcal{D} | \mathbf{w}, \sigma^2) = \prod_{n=1}^N p(y_n | \mathbf{w}^\top \mathbf{x}, \sigma^2) = \mathcal{N}(\mathbf{y} | \mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I}_N) \quad (15.40)$$

25 where  $\mathbf{I}_N$  is the  $N \times N$  identity matrix. We can then use Bayes rule for Gaussians (Equation (2.82))  
26 to derive the posterior, which is as follows:  
27

$$\underline{28} \quad p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \sigma^2) \propto \mathcal{N}(\mathbf{w} | \tilde{\mathbf{w}}, \breve{\Sigma}) \mathcal{N}(\mathbf{y} | \mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I}_N) = \mathcal{N}(\mathbf{w} | \hat{\mathbf{w}}, \hat{\Sigma}) \quad (15.41)$$

$$\underline{29} \quad \hat{\mathbf{w}} \triangleq \hat{\Sigma}^{-1} (\breve{\Sigma}^{-1} \tilde{\mathbf{w}} + \frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{y}) \quad (15.42)$$

$$\underline{31} \quad \hat{\Sigma} \triangleq (\breve{\Sigma}^{-1} + \frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X})^{-1} \quad (15.43)$$

34 where  $\hat{\mathbf{w}}$  is the posterior mean, and  $\hat{\Sigma}$  is the posterior covariance.

35 Now suppose  $\tilde{\mathbf{w}} = \mathbf{0}$  and  $\breve{\Sigma} = \tau^2 \mathbf{I}$ . In this case, the posterior mean becomes

$$\underline{37} \quad \hat{\mathbf{w}} = \frac{1}{\sigma^2} \hat{\Sigma} \mathbf{X}^\top \mathbf{y} = \left( \frac{\sigma^2}{\tau^2} \mathbf{I} + \mathbf{X}^\top \mathbf{X} \right)^{-1} \mathbf{X}^\top \mathbf{y} \quad (15.44)$$

39 If we define  $\lambda = \frac{\sigma^2}{\tau^2}$ , we see this is equivalent to **ridge regression**, which optimizes  
40

$$\underline{41} \quad \mathcal{L}(\mathbf{w}) = \text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|^2 \quad (15.45)$$

43 where RSS is the residual sum of squares:

$$\underline{45} \quad \text{RSS}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 = \frac{1}{2} (\mathbf{X}\mathbf{w} - \mathbf{y})^\top (\mathbf{X}\mathbf{w} - \mathbf{y}) \quad (15.46)$$

1 **15.2.1.2 Noise variance is unknown**

3 In this section, we assume  $\mathbf{w}$  and  $\sigma^2$  are both unknown. The likelihood is given by

$$\ell(\mathbf{w}, \sigma^2) = p(\mathcal{D}|\mathbf{w}, \sigma^2) \propto (\sigma^2)^{-N_D/2} \exp\left(-\frac{1}{2\sigma^2} \sum_{n=1}^{N_D} (y_n - \mathbf{w}^\top \mathbf{x}_n)^2\right) \quad (15.47)$$

8 Since the regression weights now depend on  $\sigma^2$  in the likelihood, the conjugate prior for  $\mathbf{w}$  has the  
9 form

$$p(\mathbf{w}|\sigma^2) = \mathcal{N}(\mathbf{w}|\tilde{\mathbf{w}}, \sigma^2 \tilde{\Sigma}) \quad (15.48)$$

12 For the noise variance  $\sigma^2$ , the conjugate prior is based on the inverse Gamma distribution, which has  
13 the form

$$\text{IG}(\sigma^2|\tilde{a}, \tilde{b}) = \frac{\tilde{b}^{\tilde{a}}}{\Gamma(\tilde{a})} (\sigma^2)^{-(\tilde{a}+1)} \exp(-\frac{\tilde{b}}{\sigma^2}) \quad (15.49)$$

17 (See Section 2.2.3.4 for more details.) Putting these two together, we find that the joint conjugate  
18 prior is the **normal inverse Gamma** distribution:

$$\text{NIG}(\mathbf{w}, \sigma^2|\tilde{\mathbf{w}}, \tilde{\Sigma}, \tilde{a}, \tilde{b}) \triangleq \mathcal{N}(\mathbf{w}|\tilde{\mathbf{w}}, \sigma^2 \tilde{\Sigma}) \text{IG}(\sigma^2|\tilde{a}, \tilde{b}) \quad (15.50)$$

$$\begin{aligned} &= \frac{\tilde{b}^{\tilde{a}}}{(2\pi)^{D/2} |\tilde{\Sigma}|^{1/2} \Gamma(\tilde{a})} (\sigma^2)^{-(\tilde{a}+(D/2)+1)} \\ &\times \exp\left[-\frac{(\mathbf{w}-\tilde{\mathbf{w}})^\top \tilde{\Sigma}^{-1} (\mathbf{w}-\tilde{\mathbf{w}}) + 2\tilde{b}}{2\sigma^2}\right] \end{aligned} \quad (15.51)$$

This results in the following posterior:

$$p(\mathbf{w}, \sigma^2|\mathcal{D}) = \text{NIG}(\mathbf{w}, \sigma^2|\hat{\mathbf{w}}, \hat{\Sigma}, \hat{a}, \hat{b}) \quad (15.52)$$

$$\hat{\mathbf{w}} = \hat{\Sigma} (\tilde{\Sigma}^{-1} \tilde{\mathbf{w}} + \mathbf{X}^\top \mathbf{y}) \quad (15.53)$$

$$\hat{\Sigma} = (\tilde{\Sigma}^{-1} + \mathbf{X}^\top \mathbf{X})^{-1} \quad (15.54)$$

$$\hat{a} = \tilde{a} + N/2 \quad (15.55)$$

$$\hat{b} = \tilde{b} + \frac{1}{2} \left( \tilde{\mathbf{w}}^\top \tilde{\Sigma}^{-1} \tilde{\mathbf{w}} + \mathbf{y}^\top \mathbf{y} - \hat{\mathbf{w}}^\top \hat{\Sigma}^{-1} \hat{\mathbf{w}} \right) \quad (15.56)$$

The expressions for  $\hat{\mathbf{w}}$  and  $\hat{\Sigma}$  are similar to the case where  $\sigma^2$  is known. The expression for  $\hat{a}$  is also intuitive, since it just updates the counts. The expression for  $\hat{b}$  can be interpreted as follows: it is the prior sum of squares,  $\tilde{b}$ , plus the empirical sum of squares,  $\mathbf{y}^\top \mathbf{y}$ , plus a term due to the error in the prior on  $\mathbf{w}$ .

The posterior marginals are as follows. For the variance, we have

$$p(\sigma^2|\mathcal{D}) = \int p(\mathbf{w}|\sigma^2, \mathcal{D}) p(\sigma^2|\mathcal{D}) d\mathbf{w} = \text{IG}(\sigma^2|\hat{a}, \hat{b}) \quad (15.57)$$

For the regression weights, it can be shown that

$$p(\mathbf{w}|\mathcal{D}) = \int p(\mathbf{w}|\sigma^2, \mathcal{D}) p(\sigma^2|\mathcal{D}) d\sigma^2 = \mathcal{T}(\mathbf{w}|\hat{\mathbf{w}}, \frac{\hat{b}}{\hat{a}} \hat{\Sigma}, 2\hat{a}) \quad (15.58)$$

1 **15.2.1.3 Posterior predictive distribution**

3 In machine learning we usually care more about uncertainty (and accuracy) of our predictions, not  
4 our parameter estimates. Fortunately, one can derive the posterior predictive distribution in closed  
5 form. In particular, one can show that, given  $N'$  new test inputs  $\tilde{\mathbf{X}}$ , we have

$$\underline{7} \quad p(\tilde{\mathbf{y}}|\tilde{\mathbf{X}}, \mathcal{D}) = \int \int p(\tilde{\mathbf{y}}|\tilde{\mathbf{X}}, \mathbf{w}, \sigma^2) p(\mathbf{w}, \sigma^2 | \mathcal{D}) d\mathbf{w} d\sigma^2 \quad (15.59)$$

$$\underline{9} \quad = \int \int \mathcal{N}(\tilde{\mathbf{y}}|\tilde{\mathbf{X}}\mathbf{w}, \sigma^2 \mathbf{I}_{N'}) \text{NIG}(\mathbf{w}, \sigma^2 | \hat{\mathbf{w}}, \hat{\Sigma}, \hat{a}, \hat{b}) d\mathbf{w} d\sigma^2 \quad (15.60)$$

$$\underline{11} \quad = \mathcal{T}(\tilde{\mathbf{y}}|\tilde{\mathbf{X}}\hat{\mathbf{w}}, \frac{\hat{b}}{\hat{a}}(\mathbf{I}_{N'} + \tilde{\mathbf{X}}\hat{\Sigma}\tilde{\mathbf{X}}^\top), 2\hat{a}) \quad (15.61)$$

14 The posterior predictive mean is equivalent to “normal” linear regression, but where we plug in  
15  $\hat{\mathbf{w}} = \mathbb{E}[\mathbf{w}|\mathcal{D}]$  instead of the MLE. The posterior predictive variance has two components:  $\hat{b}/\hat{a}\mathbf{I}_{N'}$   
16 due to the measurement noise, and  $\hat{b}/\hat{a}\tilde{\mathbf{X}}\hat{\Sigma}\tilde{\mathbf{X}}^\top$  due to the uncertainty in  $\mathbf{w}$ . This latter term varies  
17 depending on how close the test inputs are to the training data. The results are similar to using a  
18 Gaussian prior (with fixed  $\hat{\sigma}^2$ ), except the predictive distribution is even wider, since we are taking  
19 into account uncertainty about  $\sigma^2$ .

20

21 **15.2.2 Uninformative priors**

23 A common criticism of Bayesian inference is the need to use a prior. This is sometimes thought to  
24 “pollute” the inferences one makes from the data. We can minimize the effect of the prior by using an  
25 uninformative prior, as we discussed in Section 3.6. Below we discuss various uninformative priors  
26 for linear regression.

27

28 **15.2.2.1 Jeffreys prior**

30 From Section 3.6.3.1, we know that the Jeffreys prior for the location parameter has the form  
31  $p(\mathbf{w}) \propto 1$ , and from Section 3.6.3.2, we know that the Jeffreys prior for the scale factor has the  
32 form  $p(\sigma) \propto \sigma^{-1}$ . We can emulate these priors using an improper NIG prior with  $\check{\mathbf{w}} = \mathbf{0}$ ,  $\check{\Sigma} = \infty \mathbf{I}$ ,  
33  $\check{a} = -D/2$  and  $\check{b} = 0$ . The corresponding posterior is given by

$$\underline{35} \quad p(\mathbf{w}, \sigma^2 | \mathcal{D}) = \text{NIG}(\mathbf{w}, \sigma^2 | \hat{\mathbf{w}}, \hat{\Sigma}, \hat{a}, \hat{b}) \quad (15.62)$$

$$\underline{36} \quad \hat{\mathbf{w}} = \hat{\mathbf{w}}_{\text{mle}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \quad (15.63)$$

$$\underline{38} \quad \hat{\Sigma} = (\mathbf{X}^\top \mathbf{X})^{-1} \triangleq \mathbf{C} \quad (15.64)$$

$$\underline{40} \quad \hat{a} = \frac{\nu}{2} \quad (15.65)$$

$$\underline{42} \quad \hat{b} = \frac{s^2 \nu}{2} \quad (15.66)$$

$$\underline{44} \quad s^2 \triangleq \frac{\|\mathbf{y} - \hat{\mathbf{y}}\|^2}{\nu} \quad (15.67)$$

$$\underline{46} \quad \nu = N - D \quad (15.68)$$

Hence the posterior distribution of the weights is given by

$$p(\mathbf{w}|\mathcal{D}) = \mathcal{T}(\mathbf{w}|\hat{\mathbf{w}}, s^2 \mathbf{C}, \nu) \quad (15.69)$$

where  $\hat{\mathbf{w}}$  is the MLE. The marginals for each weight therefore have the form

$$p(w_d|\mathcal{D}) = \mathcal{T}(w_d|\hat{w}_d, s^2 C_{dd}, \nu) \quad (15.70)$$

### 15.2.2.2 Connection to frequentist statistics

Interestingly, the posterior when using Jeffrey's prior is formally equivalent to the **frequentist sampling distribution** of the MLE, which has the form

$$p(\hat{w}_d|\mathcal{D}^*) = \mathcal{T}(\hat{w}_d|w_d, s^2 C_{dd}, \nu) \quad (15.71)$$

where  $\mathcal{D}^* = (\mathbf{X}, \mathbf{y}^*)$  is hypothetical data generated from the true model given the fixed inputs  $\mathbf{X}$ . In books on frequentist statistics, this is more commonly written in the following equivalent way (see e.g., [Ric95, p542]):

$$\frac{\hat{w}_d - w_d}{s\sqrt{C_{dd}}} \sim t_{N-D} \quad (15.72)$$

The sampling distribution is numerically the same as the posterior distribution in Equation (15.70) because  $\mathcal{T}(w|\mu, \sigma^2, \nu) = \mathcal{T}(\mu|w, \sigma^2, \nu)$ . However, it is semantically quite different, since the sampling distribution does not condition on the observed data, but instead is based on hypothetical data drawn from the model. See [BT73, p117] for more discussion of the equivalences between Bayesian and frequentist analysis of simple linear models when using uninformative priors.

### 15.2.2.3 Zellner's *g*-prior

It is often reasonable to assume an uninformative prior on  $\sigma^2$ , since that is just a scalar that does not have much influence on the results, but using an uninformative prior for  $\mathbf{w}$  can be dangerous, since the strength of the prior controls how well regularized the model is, as we know from ridge regression.

A common compromise is to use an NIG prior with  $\check{a} = -D/2$ ,  $\check{b} = 0$  (to ensure  $p(\sigma^2) \propto 1$ ) and  $\check{\mathbf{w}} = \mathbf{0}$  and  $\check{\Sigma} = g(\mathbf{X}^\top \mathbf{X})^{-1}$ , where  $g > 0$  plays a role analogous to  $1/\lambda$  in ridge regression. This is called Zellner's **g-prior** [Zel86].<sup>1</sup> We see that the prior covariance is proportional to  $(\mathbf{X}^\top \mathbf{X})^{-1}$  rather than  $\mathbf{I}$ ; this ensures that the posterior is invariant to scaling of the inputs, e.g., due to a change in the units of measurement [Min00a].

<sup>1</sup> Note this prior is conditioned on the inputs  $\mathbf{X}$ , but not the outputs  $\mathbf{y}$ ; this is totally valid in a conditional (discriminative) model, where all calculations are conditioned on  $\mathbf{X}$ , which is treated like a fixed constant input.

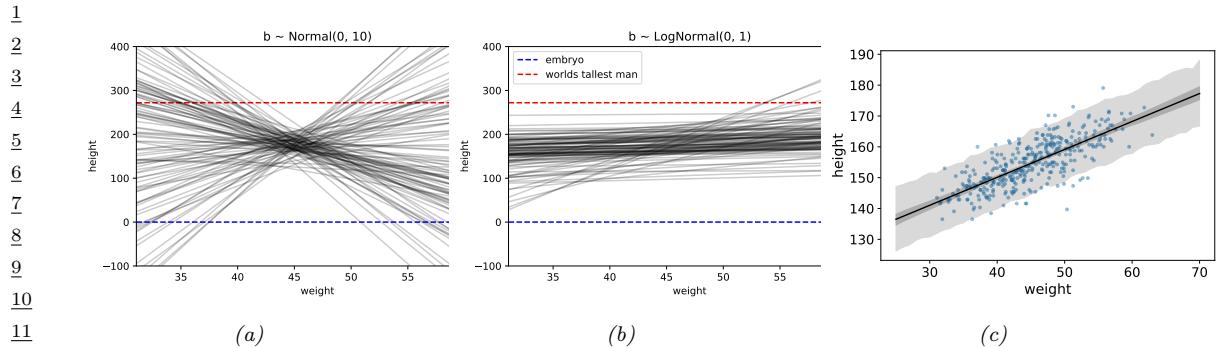


Figure 15.1: Linear regression for predicting height given weight,  $y \sim \mathcal{N}(\alpha + \beta x, \sigma^2)$ . (a) Prior predictive samples using a Gaussian prior for  $\beta$ . (b) Prior predictive samples using a Log-Gaussian prior for  $\beta$ . (c) Posterior predictive samples using the Log-Gaussian prior. The inner shaded band is the 95% credible interval for  $\mu$ , representing epistemic uncertainty. The outer shaded band is the 95% credible interval for the observations  $y$ , which also adds data uncertainty due to  $\sigma$ . Adapted from Figures 4.5 and 4.10 of [McE20]. Generated by [linreg\\_height\\_weight.ipynb](#).

With this prior, the posterior becomes

$$p(\mathbf{w}, \sigma^2 | g, \mathcal{D}) = \text{NIG}(\mathbf{w}, \sigma^2 | \mathbf{w}_N, \mathbf{V}_N, a_N, b_N) \quad (15.73)$$

$$\mathbf{V}_N = \frac{g}{g+1} (\mathbf{X}^T \mathbf{X})^{-1} \quad (15.74)$$

$$\mathbf{w}_N = \frac{g}{g+1} \hat{\mathbf{w}}_{mle} \quad (15.75)$$

$$a_N = N/2 \quad (15.76)$$

$$b_N = \frac{s^2}{2} + \frac{1}{2(g+1)} \hat{\mathbf{w}}_{mle}^T \mathbf{X}^T \mathbf{X} \hat{\mathbf{w}}_{mle} \quad (15.77)$$

Various approaches have been proposed for setting  $g$ , including cross validation, empirical Bayes [Min00a; GF00], hierarchical Bayes [Lia+08], etc.

### 15.2.3 Informative priors

In many problems, it is possible to use domain knowledge to come up with plausible priors. As an example, we consider the problem of predicting the height of a person given their weight. We will use a dataset collected from Kalahari foragers by the anthropologist Nancy Howell (this example is from the book “Statistical Rethinking” [McE20, p93]).

Let  $x_i$  be the weight (in kg) and  $y_i$  be height (in cm) of the  $i$ 'th person, and let  $\bar{x}$  be the mean of the inputs. The observation model is given by

$$y_i \sim \mathcal{N}(\mu_i, \sigma) \quad (15.78)$$

$$\mu_i = \alpha + \beta(x_i - \bar{x}) \quad (15.79)$$

We see that the intercept  $\alpha$  is the predicted output if  $x_i = \bar{x}$ , and the slope  $\beta$  is the predicted change in height per unit change in weight above or below the average weight.

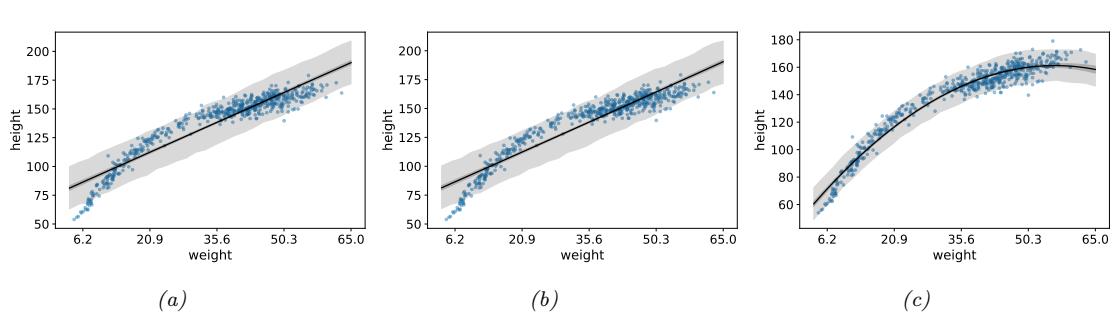


Figure 15.2: Linear regression for predicting height given weight for the full dataset (including children) using polynomial regression. (a) Posterior fit for linear model with log-Gaussian prior for  $\beta_1$ . (b) Posterior fit for quadratic model with log-Gaussian prior for  $\beta_2$ . (c) Posterior fit for quadratic model with Gaussian prior for  $\beta_2$ . Adapted from Figure 4.11 of [McE20]. Generated by [linreg\\_height\\_weight.ipynb](#).

The question is: what priors should we use? To be truly Bayesian, we should set these before looking at the data. A sensible prior for  $\alpha$  is the height of a “typical person”, with some spread. We use  $\alpha \sim \mathcal{N}(178, 20)$ , since the author of the “Rethinking Statistics” book from which this example is taken is 178cm. By using a standard deviation of 20, the prior puts 95% probability on the broad range of  $178 \pm 40$ .

What about the prior for  $\beta$ ? It is tempting to use a **vague prior**, or **weak prior**, such as  $\beta \sim \mathcal{N}(0, 10)$ , which is similar to a flat (uniform) prior, but more concentrated at 0 (a form of mild regularization). To see if this is reasonable, we can compute samples from the **prior predictive distribution**, i.e., we sample  $(\alpha_s, \beta_s) \sim p(\alpha)p(\beta)$ , and then plot  $\alpha_s x + \beta_s$  for a range of  $x$  values, for different samples  $s = 1 : S$ . The results are shown in Figure 15.1a. We see that this is not a very sensible prior. For example, we see that it suggests that it is just as likely for the height to decrease with weight as increase with weight, which is not plausible. In addition, it predicts heights which are larger than the world’s tallest person (272 cm) and smaller than the world’s shortest person (an embryo, of size 0).

We can encode the monotonically increasing relationship between weight and height by restricting  $\beta$  to be positive. An easy way to do this is to use a log-normal or log-Gaussian prior. (If  $\tilde{\beta} = \log(\beta)$  is Gaussian, then  $e^{\tilde{\beta}}$  must be positive.) Specifically, we will assume  $\beta \sim \mathcal{LN}(0, 1)$ . Samples from this prior are shown in Figure 15.1b. This is much more reasonable.

Finally we must choose a prior over  $\sigma$ . In [McE20] they use  $\sigma \sim \text{Unif}(0, 50)$ . This ensures that  $\sigma$  is positive, and that the prior predictive distribution for the output is within 100cm of the average height. However, it is usually easier to specify the expected value for  $\sigma$  than an upper bound. To do this, we can use  $\sigma \sim \text{Expon}(\lambda)$ , where  $\lambda$  is the rate. We then set  $\mathbb{E}[\sigma] = 1/\lambda$  to the value of the standard deviation that we expect. For example, we can use the empirical standard deviation of the data.

Since these priors are no longer conjugate, we cannot compute the posterior in closed form. However, we can use a variety of approximate inference methods. In this simple example, it suffices to use a quadratic (Laplace) approximation (see Section 7.4.3). The results are shown in Figure 15.1c, and look sensible.

So far, we have only considered a subset of the data, corresponding to adults over the age of 18. If

<sup>1</sup> we include children, we find that the mapping from weight to height is nonlinear. This is illustrated  
<sup>2</sup> in Figure 15.2a. We can fix this problem by using **polynomial regression**. For example, consider a  
<sup>3</sup> quadratic expansion of the standardized features  $x_i$ :  
<sup>4</sup>

$$\mu_i = \alpha + \beta_1 x_i + \beta_2 x_i^2 \quad (15.80)$$

<sup>5</sup> If we use a log-Gaussian prior for  $\beta_2$ , we find that the model is too constrained, and it underfits.  
<sup>6</sup> This is illustrated in Figure 15.2b. The reason is that we need to use an inverted quadratic with  
<sup>7</sup> a negative coefficient, but since this is disallowed by the prior, the model ends up not using this  
<sup>8</sup> degree of freedom (we find  $\mathbb{E}[\beta_2|\mathcal{D}] \approx 0.08$ ). If we use a Gaussian prior on  $\beta_2$ , we avoid this problem,  
<sup>9</sup> illustrated in Figure 15.2c.  
<sup>10</sup>

<sup>11</sup> This example shows that it can be useful to think about the functional form of the mapping from  
<sup>12</sup> inputs to outputs in order to specify sensible priors.  
<sup>13</sup>

#### <sup>14</sup> 15.2.4 Spike and slab prior

<sup>15</sup> It is often useful to be able to select a subset of the input features when performing prediction, either  
<sup>16</sup> to reduce overfitting, or to improve interpretability of the model. This can be achieved if we ensure  
<sup>17</sup> that the weight vector  $\mathbf{w}$  is **sparse** (i.e., has many zero elements), since if  $w_d = 0$ , then  $x_d$  plays no  
<sup>18</sup> role in the inner product  $\mathbf{w}^\top \mathbf{x}$ .  
<sup>19</sup>

<sup>20</sup> The canonical way to achieve sparsity when using Bayesian inference is to use a **spike-and-slab**  
<sup>21</sup> (**SS**) prior [MB88], which has the form of a 2 component mixture model, with one component being  
<sup>22</sup> a “spike” at 0, and the other being a uniform “slab” between  $-a$  and  $a$ :  
<sup>23</sup>

$$\begin{aligned} p(\mathbf{w}) &= \prod_{d=1}^D (1 - \pi) \delta(w_d) + \pi \text{Unif}(w_d | -a, a) \end{aligned} \quad (15.81)$$

<sup>24</sup> where  $\pi$  is the prior probability that each coefficient is non-zero. The corresponding log prior on the  
<sup>25</sup> coefficients is thus  
<sup>26</sup>

$$\log p(\mathbf{w}) = \|\mathbf{w}\|_0 \log(1 - \pi) + (D - \|\mathbf{w}\|_0) \log \pi = -\lambda \|\mathbf{w}\|_0 + \text{const} \quad (15.82)$$

<sup>27</sup> where  $\lambda = \log \frac{\pi}{1-\pi}$  controls the sparsity of the model, and  $\|\mathbf{w}\|_0 = \sum_{d=1}^D \mathbb{I}(w_d \neq 0)$  is the  $\ell_0$  **norm**  
<sup>28</sup> of the weights. Thus MAP estimation with a spike and slab prior is equivalent  $\ell_0$  **regularization**;  
<sup>29</sup> this penalizes the number of non-zero coefficients. Interestingly, posterior samples will also be sparse.

<sup>30</sup> By contrast, consider using a Laplace prior. The **lasso** estimator uses MAP estimation, which  
<sup>31</sup> results in a sparse estimate. However, posterior samples are not sparse. Interestingly, [EY09] show  
<sup>32</sup> theoretically (and [SPZ09] confirm experimentally) that using the posterior mean with a spike-and-  
<sup>33</sup> slab prior also results in better prediction accuracy than using the posterior mode with a Laplace  
<sup>34</sup> prior.  
<sup>35</sup>

<sup>36</sup> In practice, we often approximate the uniform slab with a broad Gaussian distribution,  
<sup>37</sup>

$$p(\mathbf{w}) = \prod_d (1 - \pi) \delta(w_d) + \pi \mathcal{N}(w_d | 0, \sigma_w^2) \quad (15.83)$$

<sup>38</sup> As  $\sigma_w^2 \rightarrow \infty$ , the second term approaches a uniform distribution over  $[-\infty, +\infty]$ . We can implement  
<sup>39</sup> the mixture model by associating a binary random variable,  $s_d \sim \text{Ber}(\pi)$ , with each coefficient, to  
<sup>40</sup> indicate if the coefficient is “on” or “off”.  
<sup>41</sup>