

11 *Figure 4.37: A UPGM and two failed attempts to represent it as a DPGM. From Figure 3.10 of [KF09a].*  
 12 Used with kind permission of Daphne Koller.

13  
 14  
 15  
 16 the arrows, we get  $A - C - B$ , which asserts  $A \perp B|C$  and  $A \not\perp B$ , which is not consistent with  
 17 the independence statements encoded by the DPGM. In fact, there is no UPGM that can precisely  
 18 represent all and only the two CI statements encoded by a v-structure. In general, CI properties in  
 19 UPGMs are monotonic, in the following sense: if  $A \perp B|C$ , then  $A \perp B|(C \cup D)$ . But in DPGMs,  
 20 CI properties can be non-monotonic, since conditioning on extra variables can eliminate conditional  
 21 independencies due to explaining away.

22 As an example of some CI relationships that can be perfectly modeled by a UPGM but not a  
 23 DPGM, consider the 4-cycle shown in Figure 4.37(a). One attempt to model this with a DPGM is  
 24 shown in Figure 4.37(b). This correctly asserts that  $A \perp C|B,D$ . However, it incorrectly asserts  
 25 that  $B \perp D|A$ . Figure 4.37(c) is another incorrect DPGM: it correctly encodes  $A \perp C|B,D$ , but  
 26 incorrectly encodes  $B \perp D$ . In fact there is no DPGM that can precisely represent all and only the  
 27 CI statements encoded by this UPGM.

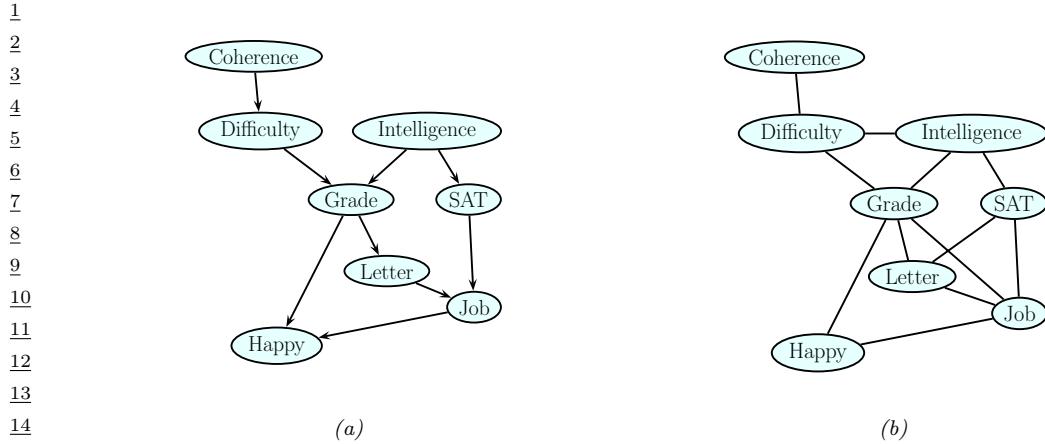
28 Some distributions can be perfectly modeled by either a DPGM or a UPGM; the resulting graphs  
 29 are called **decomposable** or **chordal**. Roughly speaking, this means the following: if we collapse  
 30 together all the variables in each maximal clique, to make “mega-variables”, the resulting graph will  
 31 be a tree. Of course, if the graph is already a tree (which includes chains as a special case), it will  
 32 already be chordal.

### 33 4.5.2 Converting between a directed and undirected model

34 Although DPGMs and UPGMs are not in general equivalent, if we are willing to allow the graph to  
 35 encode fewer CI properties than may strictly hold, then we can safely convert one to the other, as we  
 36 explain below.

#### 40 4.5.2.1 Converting a DPGM to a UPGM

41 We can easily convert a DPGM to a UPGM as follows. First, any “unmarried” parents that share a  
 42 child must get “married”, by adding an edge between them; this process is known as **moralization**.  
 43 Then we can drop the arrows, resulting in an undirected graph. The reason we need to do this is to  
 44 ensure that the CI properties of the UGM match those of the DGM, as explained in Section 4.3.6.2.  
 45 It also ensures there is a clique that can “store” the CPDs of each family.



20 Let us consider an example from [KF09a]. We will use the (full version of the student network  
21 shown in Figure 4.38(a). The corresponding joint has the following form:

$$P(C, D, I, G, S, L, J, H) \quad (4.149)$$

$$= P(C)P(D|C)P(I)P(G|I, D)P(S|I)P(L|G)P(J|L, S)P(H|G, J) \quad (4.150)$$

26 Next, we define a potential or factor for every CPD, yielding

$$p(C, D, I, G, S, L, J, H) = \psi_C(C)\psi_D(D, C)\psi_I(I)\psi_G(G, I, D) \quad (4.151)$$

$$\psi_S(S, I)\psi_L(L, G)\psi_J(J, L, S)\psi_H(H, G, J) \quad (4.152)$$

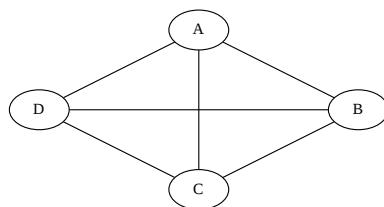
30 All the potentials are **locally normalized**, since they are CPDs, there is no need for a global  
31 normalization constant, so  $Z = 1$ . The corresponding undirected graph is shown in Figure 4.38(b).  
32 We see the that we have added D-I, G-J, and L-S moralization edges.<sup>11</sup>

### 35 4.5.2.2 Converting a UPGM to a DPGM

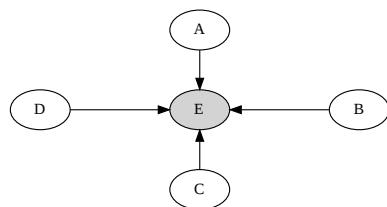
36 To convert a UPGM to a DPGM, we proceed as follows. For each potential function  $\psi_c(\mathbf{x}_c; \boldsymbol{\theta}_c)$ , we  
37 create a “dummy node”, call it  $Y_c$ , which is “clamped” to a special observed state, call it  $y_c^*$ . We  
38 then define  $p(Y_c = y_c^* | \mathbf{x}_c) = \psi_c(\mathbf{x}_c; \boldsymbol{\theta}_c)$ . This “local evidence” CPD encodes the same factor as in the  
39 DGM. The overall joint has the form  $p_{\text{undir}}(\mathbf{x}) \propto p_{\text{dir}}(\mathbf{x}, \mathbf{y}^*)$ .

40 As an example, consider the UPGM in Figure 4.39(a), which defines the joint  $p(A, B, C, D) = \psi(A, B, C, D)/Z$ . We can represent this as a DPGM by adding a dummy  $E$  node, which is a child  
41 of all the other nodes. We set  $E = 1$  and define the CPD  $p(E = 1 | A, B, C, D) \propto \psi(A, B, C, D)$ . By  
42 conditioning on this observed child, all the parents become dependent, as in the UGM.  
43

45 11. We will see this example again in Section 9.4, where we use it to illustrate the variable elimination inference  
46 algorithm.

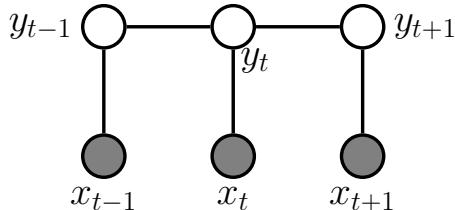


(a)

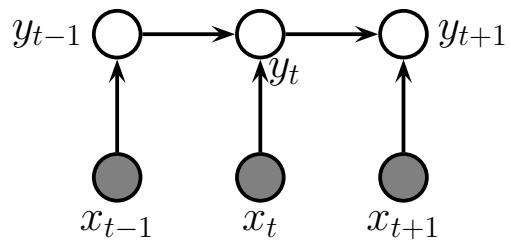


(b)

Figure 4.39: (a) An undirected graphical model. (b) A directed equivalent, obtained by adding a dummy observed child node.



(a)



(b)

Figure 4.40: Two discriminative models for sequential data. (a) An undirected model (CRF). (b) A directed model (MEMM).

### 4.5.3 Conditional directed vs undirected PGMs and the label bias problem

Directed and undirected models behave somewhat differently in the conditional (discriminative) setting. As an example of this, let us compare the 1d undirected CRF in Figure 4.40a with the directed Markov chain in Section 4.5.3. (This latter model is called a maximum entropy Markov model (MEMM), which is a reference to the connection with maxent models discussed in Section 4.3.4.) The MEMM suffers from a subtle problem compared to the CRF known (rather obscurely) as the **label bias** problem [LMP01]. The problem is that local features at time  $t$  do not influence states prior to time  $t$ . That is,  $y_{t-1} \perp \mathbf{x}_t | y_t$ , thus blocking information flow backwards in time.

To understand what this means in practice, consider the part of speech tagging task which we discussed in Section 4.4.1.1. Suppose we see the word “banks”; this could be a verb (as in “he banks at Chase”), or a noun (as in “the river banks were overflowing”). Locally the part of speech tag for the word is ambiguous. However, suppose that later in the sentence, we see the word “fishing”; this gives us enough context to infer that the sense of “banks” is “river banks” and not “financial banks”. However, in an MEMM the “fishing” evidence will not flow backwards, so we will not be able to infer the correct label for “banks”. The CRF does not have this problem.

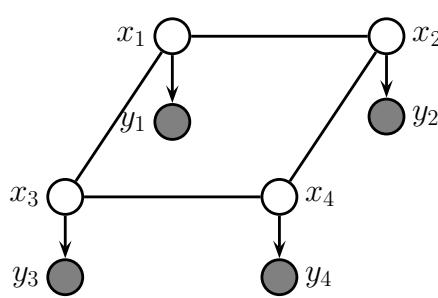


Figure 4.41: A grid-structured MRF with hidden nodes  $x_i$  and local evidence nodes  $y_i$ . The prior  $p(\mathbf{x})$  is an undirected Ising model, and the likelihood  $p(\mathbf{y}|\mathbf{x}) = \prod_i p(y_i|x_i)$  is a directed fully factored model.

The label bias problem in MEMMs occurs because directed models are **locally normalized**, meaning each CPD sums to 1. By contrast, MRFs and CRFs are **globally normalized**, which means that local factors do not need to sum to 1, since the partition function  $Z$ , which sums over all joint configurations, will ensure the model defines a valid distribution.

However, this solution comes at a price: in a CRF, we do not get a valid probability distribution over  $\mathbf{y}_{1:T}$  until we have seen the whole sentence, since only then can we normalize over all configurations. Consequently, CRFs are not as useful as directed probabilistic graphical models (DPGM) for online or real-time inference. Furthermore, the fact that  $Z$  is a function of all the parameters makes CRFs less modular and much slower to train than DPGM's, as we discuss in Section 4.4.3.

#### 4.5.4 Combining directed and undirected graphs

We can also define graphical models that contain directed and undirected edges. We discuss a few examples below.

##### 4.5.4.1 Chain graphs

A **chain graph** is a PGM which may have both directed and undirected edges, but without any directed cycles. A simple example is shown in Figure 4.41, which defines the following joint model:

$$p(\mathbf{x}_{1:D}, \mathbf{y}_{1:D}) = p(\mathbf{x}_{1:D})p(\mathbf{y}_{1:D}|\mathbf{x}_{1:D}) = \left[ \frac{1}{Z} \prod_{i \sim j} \psi_{ij}(x_i, x_j) \right] \left[ \prod_{i=1}^D p(y_i|x_i) \right] \quad (4.153)$$

In this example, the prior  $p(\mathbf{x})$  is specified by a UPGM, and the likelihood  $p(\mathbf{y}|\mathbf{x})$  is specified as a fully factorized DPGM.

More generally, a chain graph can be defined in terms of a **partially directed acyclic graph (PDAG)**. This is a graph which can be decomposed into a directed graph of **chain components**, where the nodes within each chain component are connected with each other only with undirected edges. See Figure 4.42 for an example.

We can use a PDAG to define a joint distribution using  $\prod_i p(C_i|pa(C_i))$ , where each  $C_i$  is a chain component, and each CPD is a conditional random field. For example, referring to Figure 4.42, we

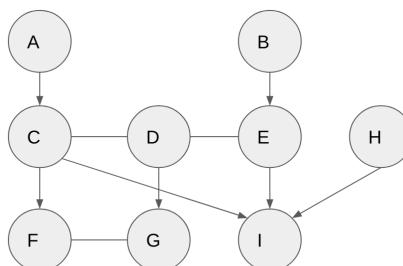


Figure 4.42: A partially directed acyclic graph (PDAG). The chain components are  $\{A\}$ ,  $\{B\}$ ,  $\{C, D, E\}$ ,  $\{F, G\}$ ,  $\{H\}$  and  $\{I\}$ . Adapted from Figure 4.15 of [KF09a].

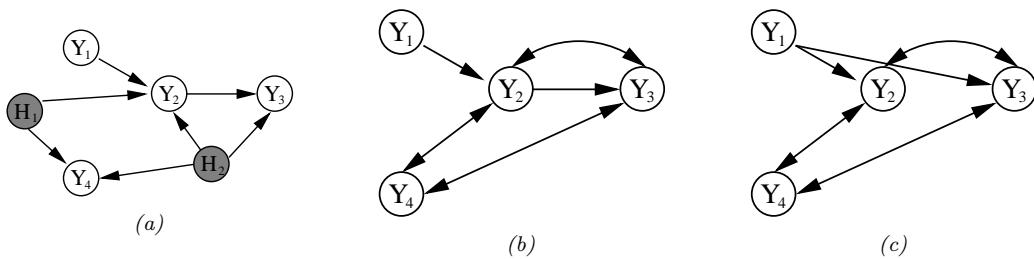


Figure 4.43: (a) A DAG with two hidden variables (shaded). (b) The corresponding ADMG. The bidirected edges reflect correlation due to the hidden variable. (c) A Markov equivalent ADMG. From Figure 3 of [SG09]. Used with kind permission of Ricardo Silva.

have

$$p(A, B, \dots, I) = p(A)p(B)p(C, D, E|A, B)p(F, G|C, D)p(H)p(I|C, E, H) \quad (4.154)$$

$$p(C, D, E|A, B) = \frac{1}{Z(A, B)}\phi(A, C)\phi(B, E)\phi(C, D)\phi(D, E) \quad (4.155)$$

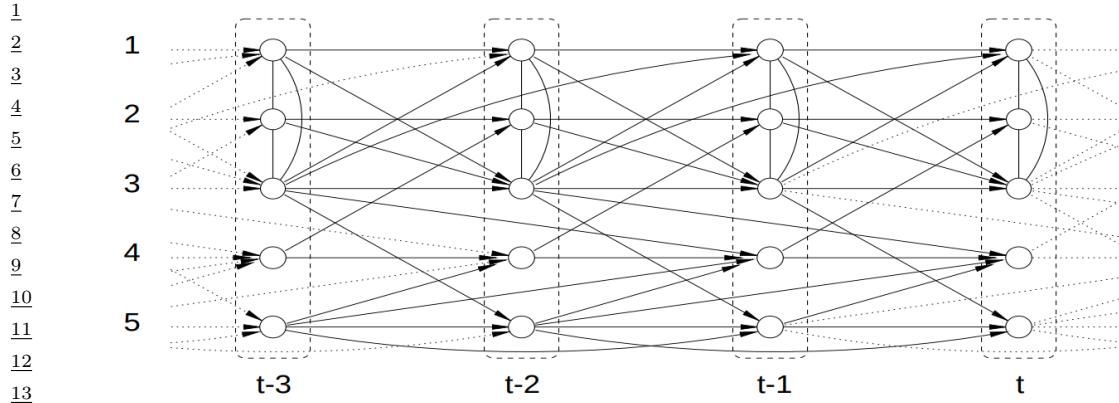
$$p(F, G|C, D) = \frac{1}{Z(C, D)}\phi(F, C)\phi(G, D)\phi(F, G) \quad (4.156)$$

For more details, see e.g., [KF09a, Sec 4.6.2].

#### 4.5.4.2 Acyclic directed mixed graphs

One can show [Pea09b, p51] that every latent variable DPGM can be rewritten in a way such that every latent variable is a root node with exactly two observed children. This is called the **projection** of the latent variable PGM, and is observationally indistinguishable from the original model.

Each such latent variable root node induces a dependence between its two children. We can represent this with a directed arc. The resulting graph is called an **acyclic directed mixed graph** or **ADMG**. See Figure 4.43 for an example. (A **mixed graph** is one with undirected, unidirected, and bidirected edges.)



*Figure 4.44: A VAR(2) process represented as a dynamic chain graph. From [DE00]. Used with kind permission of Rainer Dahlhaus.*

One can determine CI properties of ADMGs using a technique called **m-separation** [Ric03]. This is equivalent to d-separation in a graph where every bidirected edge  $Y_i \leftrightarrow Y_j$  is replaced by  $Y_i \leftarrow X_{ij} \rightarrow Y_j$ , where  $X_{ij}$  is a hidden variable for that edge.

The most common example of ADMGs is when everything is linear-Gaussian. This is known as a structural equation model and is discussed in Section 4.7.2.

#### 4.5.5 Comparing directed and undirected Gaussian PGMs

In this section, we compare directed and undirected Gaussian graphical models. In Section 4.2.3, we saw that directed GGMs correspond to sparse regression matrices. In Section 4.3.5, we saw that undirected GGMs correspond to sparse precision matrices.

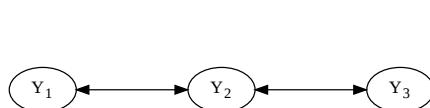
The advantage of the DAG formulation is that we can make the regression weights  $\mathbf{W}$ , and hence  $\Sigma$ , be conditional on covariate information [Pou04], without worrying about positive definite constraints. The disadvantage of the DAG formulation is its dependence on the order, although in certain domains, such as time series, there is already a natural ordering of the variables.

It is actually possible to combine both directed and undirected representations, resulting in a model known as a (Gaussian) **chain graph**. For example, consider a discrete-time, second-order Markov chain in which the observations are continuous,  $\mathbf{x}_t \in \mathbb{R}^D$ . The transition function can be represented as a (vector-valued) linear-Gaussian CPD:

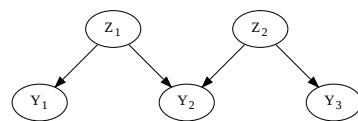
$$p(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{x}_{t-2}, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}_t | \mathbf{A}_1 \mathbf{x}_{t-1} + \mathbf{A}_2 \mathbf{x}_{t-2}, \Sigma) \quad (4.157)$$

This is called **vector auto-regressive** or **VAR** process of order 2. Such models are widely used in econometrics for time-series forecasting.

The time series aspect is most naturally modeled using a DPGM. However, if  $\Sigma^{-1}$  is sparse, then the correlation amongst the components within a time slice is most naturally modeled using a UPGM.



(a)



(b)

Figure 4.45: (a) A bi-directed graph. (b) The equivalent DAG. Here the  $z$  nodes are latent confounders. Adapted from Figures 5.12-5.13 of [Cho11].

For example, suppose we have

$$\mathbf{A}_1 = \begin{pmatrix} \frac{3}{5} & 0 & \frac{1}{5} & 0 & 0 \\ 0 & \frac{3}{5} & 0 & -\frac{1}{5} & 0 \\ \frac{2}{5} & \frac{1}{5} & \frac{3}{5} & 0 & 0 \\ 0 & 0 & 0 & -\frac{1}{2} & \frac{1}{5} \\ 0 & 0 & \frac{1}{5} & 0 & \frac{2}{5} \end{pmatrix}, \quad \mathbf{A}_2 = \begin{pmatrix} 0 & 0 & -\frac{1}{5} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & \frac{1}{5} & 0 & \frac{1}{3} \\ 0 & 0 & 0 & 0 & -\frac{1}{5} \end{pmatrix} \quad (4.158)$$

and

$$\boldsymbol{\Sigma} = \begin{pmatrix} 1 & \frac{1}{2} & \frac{1}{3} & 0 & 0 \\ \frac{1}{2} & 1 & -\frac{1}{3} & 0 & 0 \\ \frac{1}{3} & -\frac{1}{3} & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix}, \quad \boldsymbol{\Sigma}^{-1} = \begin{pmatrix} 2.13 & -1.47 & -1.2 & 0 & 0 \\ -1.47 & 2.13 & 1.2 & 0 & 0 \\ -1.2 & 1.2 & 1.8 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix} \quad (4.159)$$

The resulting graphical model is illustrated in Figure 4.44. Zeros in the transition matrices  $\mathbf{A}_1$  and  $\mathbf{A}_2$  correspond to absent directed arcs from  $x_{t-1}$  and  $x_{t-2}$  into  $x_t$ . Zeros in the precision matrix  $\boldsymbol{\Sigma}^{-1}$  correspond to absent undirected arcs between nodes in  $x_t$ .

#### 4.5.5.1 Covariance graphs

Sometimes we have a sparse covariance matrix rather than a sparse precision matrix. This can be represented using a **bi-directed graph**, where each edge has arrows in both directions, as in Figure 4.45(a). Here nodes that are not connected are unconditionally independent. For example in Figure 4.45(a) we see that  $Y_1 \perp Y_3$ . In the Gaussian case, this means  $\Sigma_{1,3} = \Sigma_{3,1} = 0$ . (A graph representing a sparse covariance matrix is called a **covariance graph**, see e.g., [Pen13]). By contrast, if this were an undirected model, we would have that  $Y_1 \perp Y_3 | Y_2$ , and  $\Lambda_{1,3} = \Lambda_{3,1} = 0$ , where  $\Lambda = \boldsymbol{\Sigma}^{-1}$ .

A bidirected graph can be converted to a DAG with latent variables, where each bidirected edge is replaced with a hidden variable representing a hidden common cause, or **confounder**, as illustrated in Figure 4.45(b). The relevant CI properties can then be determined using d-separation.

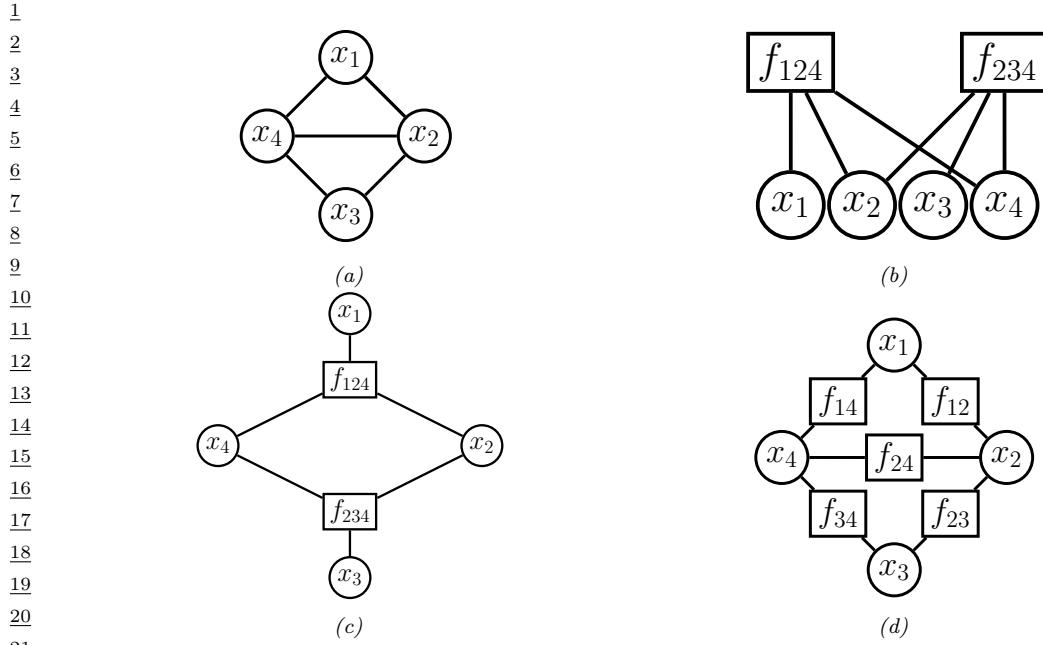


Figure 4.46: (a) A simple UPGM. (b) A factor graph representation assuming one potential per maximal clique. (c) Same as (b), but graph is visualized differently. (d) A factor graph representation assuming one potential per edge.

## 4.6 PGM extensions

In this section, we discuss some extensions of the basic PGM framework.

### 4.6.1 Factor graphs

A **factor graph** [KFL01; Loe04] is a graphical representation that unifies directed and undirected models. They come in two main “flavors”. The original version uses a bipartite graph, where we have nodes for random variables and nodes for factors, as we discuss in Section 4.6.1.1. An alternative form, known as a **Forney factor graphs** [For01] just has nodes for factors, and the variables are associated with edges, as we explain in Section 4.6.1.2.

#### 4.6.1.1 Bipartite factor graphs

A **factor graph** is an undirected bipartite graph with two kinds of nodes. Round nodes represent variables, square nodes represent factors, and there is an edge from each variable to every factor that mentions it. For example, consider the MRF in Figure 4.46(a). If we assume one potential per maximal clique, we get the factor graph in Figure 4.46(b), which represents the function

$$f(x_1, x_2, x_3, x_4) = f_{124}(x_1, x_2, x_4)f_{234}(x_2, x_3, x_4) \quad (4.160)$$

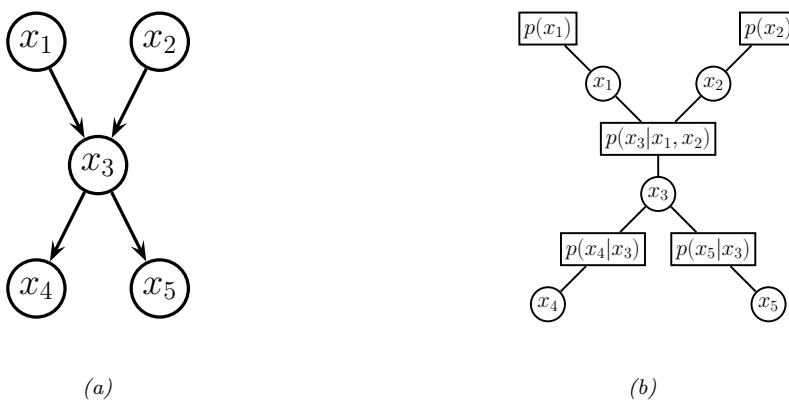


Figure 4.47: (a) A simple DPGM. (b) Its corresponding factor graph.

We can represent this in a topologically equivalent way as in Figure 4.46(c).

One advantage of factor graphs over UPGM diagrams is that they are more fine-grained. For example, suppose we associate one potential per edge, rather than per clique. In this case, we get the factor graph in Figure 4.46(d), which represents the function

$$f(x_1, x_2, x_3, x_4) = f_{14}(x_1, x_4)f_{12}(x_1, x_2)f_{34}(x_3, x_4)f_{23}(x_2, x_3)f_{24}(x_2, x_4) \quad (4.161)$$

We can also convert a DPGM to a factor graph: just create one factor per CPD, and connect that factor to all the variables that use that CPD. For example, Figure 4.47 represents the following factorization:

$$f(x_1, x_2, x_3, x_4, x_5) = f_1(x_1)f_2(x_2)f_{123}(x_1, x_2, x_3)f_{34}(x_3, x_4)f_{35}(x_3, x_5) \quad (4.162)$$

where we define  $f_{123}(x_1, x_2, x_3) = p(x_3|x_1, x_2)$ , etc. If each node has at most one parent (and hence the graph is a chain or simple tree), then there will be one factor per edge (root nodes can have their prior CPDs absorbed into their children's factors). Such models are equivalent to pairwise MRFs.

### 4.6.1.2 Forney factor graphs

A **Forney factor graph (FFG)**, also called a **normal factor graph**, is a graph in which nodes represent factors, and edges represent variables [For01; Loe04; Loe+07; CLV19]. This is more similar to standard neural network diagrams, and electrical engineering diagrams, where signals (represented as electronic pulses, or tensors, or probability distributions) propagate along wires and are modified by functions represented as nodes.

For example, consider the following factorized function:

$$f(x_1, \dots, x_5) = f_a(x_1)f_b(x_1, x_2)f_c(x_2, x_3, x_4)f_d(x_4, x_5) \quad (4.163)$$

We can visualize this as an FFG as in Figure 4.48a. The edge labeled  $x_3$  is called a **half-edge**, since it is only connected to one node; this is because  $x_3$  only participates in one factor. (Similarly for  $x_5$ .)

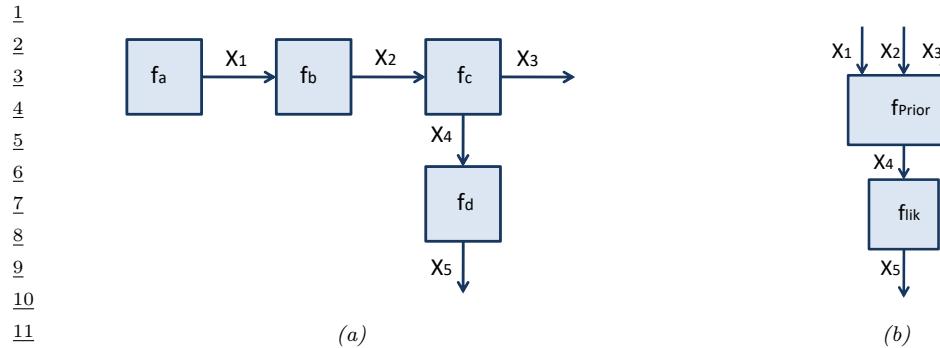


Figure 4.48: A Forney factor graph. (a) Directed version. (b) Hierarchical version.

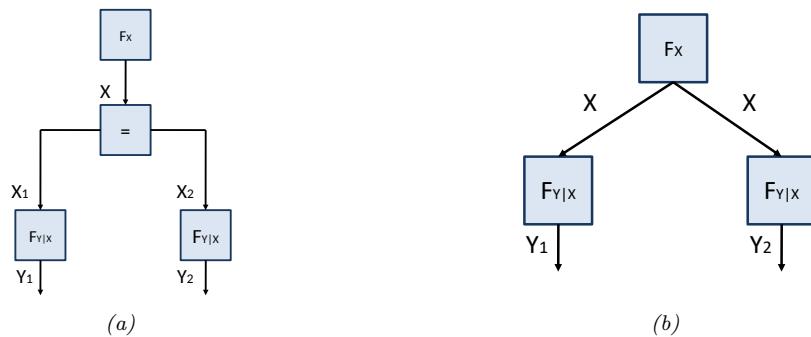


Figure 4.49: An FFG with an equality constraint node (left) and its corresponding simplified form (right).

The directionality associated with the edges is a useful mnemonic device if there is a natural order in which the variables are generated. In addition, associating directions with each edge allows us to uniquely name “messages” that are sent along each edge, which will prove useful when we discuss inference algorithms in Section 9.2.

In addition to being more similar to neural network diagrams, FFGs have the advantage over bipartite FGs in that they support hierarchical (compositional) construction, in which complex dependency structure between variables can be represented as a blackbox, with the input/output interface being represented by edges corresponding to the variables exposed by the blackbox. See Figure 4.48b for an example, which represents the function

$$f(x_1, \dots, x_5) = f_{\text{prior}}(x_1, x_2, x_3, x_4)f_{\text{lik}}(x_4, x_5) \quad (4.164)$$

The factor  $f_{\text{prior}}$  represents a (potentially complex) joint distribution  $p(x_1, x_2, x_3, x_4)$ , and the factor  $f_{\text{lik}}$  represents the likelihood term  $p(x_5|x_4)$ . Such models are widely used to build error-correcting codes (see Section 9.3.8).

To allow for variables to participate in more than 2 factors, equality constraint nodes are introduced, as illustrated in Figure 4.49(a). Formally, this is a factor defined as follows:

$$f_=(x, x_1, x_2) = \delta(x - x_1)\delta(x - x_2) \quad (4.165)$$

where  $\delta(u)$  is a Dirac delta if  $u$  is continuous, and a Kronecker delta if  $u$  is discrete. The effect of this factor is to ensure all the variables connected to the factor have the same value; intuitively, this factor acts like a “wire splitter”. Thus the function represented in Figure 4.49(a) is equivalent to the following:

$$f(x, y_1, y_2) = f_x(x) f_{y|x}(y_1, x) f_{y|x}(y_2, x) \quad (4.166)$$

This simplified form is represented in Figure 4.49(b), where we reuse the  $x$  variable across multiple edges. We have chosen the edge orientations to reflect our interpretation of the factors  $f_{y|x}(y, x)$  as likelihood terms,  $p(y|x)$ . We have also chosen to reuse the same  $f_{y|x}$  factor for both  $y$  variables; this is an example of **parameter tying**.

### 4.6.2 Probabilistic circuits

A **probabilistic circuit** is a kind of graphical model that supports efficient exact inference. It includes **arithmetic circuits** [Dar03; Dar09], **sum-product networks** (SPNs) [PD11; SCPD22]. and other kinds of model.

Here we briefly describe SPNs. An SPN is a probabilistic model, based on a directed tree-structured graph, in which terminal nodes represent univariate probability distributions and non-terminal nodes represent convex combinations (weighted sums) and products of probability functions. SPNs are similar to deep mixture models, in which we combine together dimensions. SPNs leverage context-specific independence to reduce the complexity of exact inference to time that is proportional to the number of links in the graph, as opposed to the treewidth of the graph (see Section 9.4.2).

SPNs are particularly useful for tasks such as missing data imputation of tabular data (see e.g., [Cla20; Ver+19]). A recent extension of SPNs, known as **einsum networks**, is proposed in [Peh+20] (see Section 9.6.1 for details on the connection between einstein summation and PGM inference).

### 4.6.3 Directed relational PGMs

A Bayesian network defines a joint probability distribution over a fixed number of random variables. By using plate notation (Section 4.2.8), we can define models with certain kinds of repetitive structure, and tied parameters, but many models are not expressible in this way. For example, it is not possible to represent even a simple HMM using plate notation (see Figure 29.12). Various notational extensions of plates have been proposed to handle repeated structure (see e.g., [HMK04; Die10]) but have not been widely adopted. The problem becomes worse when we have more complex domains, involving multiple objects which interact via multiple relationships.<sup>12</sup> Such models are called **relational probability models** or **RPMs**. In this section, we focus on directed RPMs; see Section 4.6.4 for the undirected case.

As in first order logic, RPMs have constant symbols (representing objects), function symbols (mapping one set of constants to another), and predicate symbols (representing relations between objects). We will assume that each function has a **type signature**. To illustrate this, consider an example from [RN19, Sec 15.1], which concerns online book reviews on sites such as Amazon. Suppose

<sup>12</sup> See e.g., this blog post from Rob Zinkov: <https://www.zinkov.com/posts/2013-07-28-stop-using-plates>.

1 there are two types of objects, Book and Customer, and the following functions and predicates:  
2

$$\text{Honest} : \text{Customer} \rightarrow \{\text{True}, \text{False}\} \quad (4.167)$$

$$\text{Kindess} : \text{Customer} \rightarrow \{1, 2, 3, 4, 5\} \quad (4.168)$$

$$\text{Quality} : \text{Book} \rightarrow \{1, 2, 3, 4, 5\} \quad (4.169)$$

$$\text{Recommendation} : \text{Customer} \times \text{Book} \rightarrow \{1, 2, 3, 4, 5\} \quad (4.170)$$

9 The constant symbols refer to specific objects. To keep things simple, we assume there are two  
10 books,  $B_1$  and  $B_2$ , and two customers,  $C_1$  and  $C_2$ . The **basic random variables** are obtained  
11 by instantiating each function with each possible combination of objects to create a set of **ground**  
12 **terms**. In this example, these variables are  $H(C_1)$ ,  $Q(B_1)$ ,  $R(C_1, B_2)$ , etc. (We use the abbreviations  
13  $H$ ,  $K$ ,  $Q$  and  $R$  for the functions Honest, Kindness, Quality and Recommendation.<sup>13</sup>)

14 We now need to specify the (conditional) distribution over these random variables. We define these  
15 distributions in terms of the generic indexed form of the variables, rather than the specific ground  
16 form. For example, we may use the following priors for the root nodes (variables with no parents):

$$H(c) \sim \text{Cat}(0.99, 0.01) \quad (4.171)$$

$$K(c) \sim \text{Cat}(0.1, 0.1, 0.2, 0.3, 0.3) \quad (4.172)$$

$$Q(b) \sim \text{Cat}(0.05, 0.2, 0.4, 0.2, 0.15) \quad (4.173)$$

22 For the recommendation nodes, we need to define a conditional distribution of the form

$$R(c, b) \sim \text{RecCPD}(H(c), K(c), Q(b)) \quad (4.174)$$

25 where RecCPD is the conditional probability distribution (CPD) for the recommendation node. If  
26 represented as a conditional probability table (CPT), this has  $2 \times 5 \times 5 = 50$  rows, each with 5  
27 entries. This table can encode our assumptions about what kind of ratings a book receives based on  
28 the quality of the book, but also properties of the reviewer, such as their honest and kindness. (More  
29 sophisticated models of human raters in the context of crowd-sourced data collection can be found in  
30 e.g., [LRC19].)

31 We can convert the above formulae into a graphical model “**template**”, as shown in Figure 4.50a.  
32 Given a set of objects, we can “**unroll**” the template to create a “**ground network**”, as shown in  
33 Figure 4.50b. There are  $C \times B + 2C + B$  random variables, with a corresponding joint state space  
34 (set of **possible worlds**) of size  $2^C 5^{C+B+BC}$ , which can get quite large. However, if we are only  
35 interested in answering specific queries, we can dynamically unroll small pieces of the network that  
36 are relevant to that query [GC90; Bre92].

37 Let us assume that only a subset of the  $R(c, b)$  entries are observed, and we would like to  
38 predict the missing entries of this matrix. This is essentially a simplified **recommender system**.  
39 (Unfortunately it ignores key aspects of the problem, such as the content/topic of the books, and  
40 the interests/preferences of the customers.) We can use standard probabilistic inference methods for  
41 graphical models (which we discuss in Chapter 9) to solve this problem.

43 Things get more interesting when we don’t know which objects are being referred to. For example,  
44 customer  $C_1$  might write a review of a book called “Probabilistic Machine Learning”, but do they

45 13. A unary function of an object that returns a basic type, such as Boolean or an integer, is often called an **attribute**  
46 of that object.

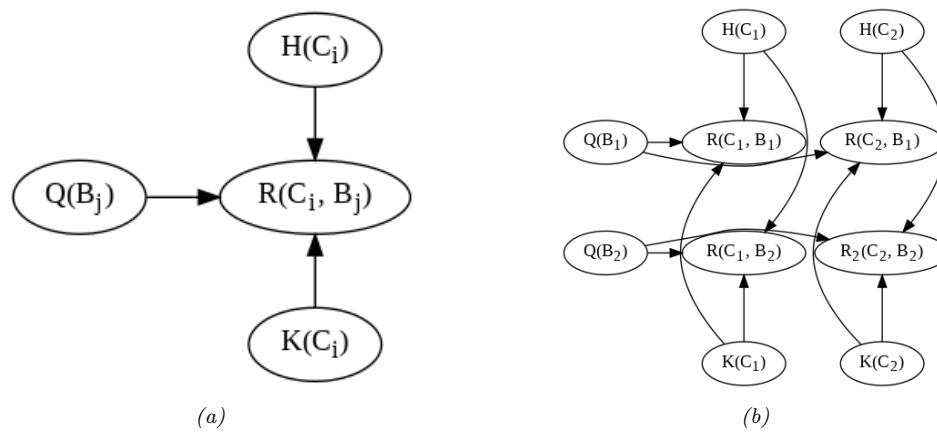


Figure 4.50: RPM for the book review domain. (a) Template for a generic customer  $C_i$  and book  $B_j$  pair.  $R$  is rating,  $Q$  is quality,  $H$  is honesty, and  $K$  is kindness. (b) Unrolled model for 2 books and 2 customers.

mean edition 1 ( $B_1$ ) or edition 2 ( $B_2$ )? To handle this kind of **relational uncertainty**, we can add all possible referents as parents to each relation. This is illustrated in Figure 4.51, where now  $Q(B_1)$  and  $Q(B_2)$  are both parents of  $R(C_1, B_1)$ . This is necessary because their review score might either depend on  $Q(B_1)$  or  $Q(B_2)$ , depending on which edition they are writing about. To disambiguate this, we create a new variable,  $L(C_i)$ , which specifies which version number of each book customer  $i$  is referring to. The new CPD for the recommendation node,  $p(R(c, b)|H(c), K(c), Q(1 : B), L(c))$ , has the form

$$R(c, b) \sim \text{RecCPT}(H(c), K(c), Q(b')) \text{ where } b' = L(c) \quad (4.175)$$

This CPD acts like a **multiplexer**, where the  $L(c)$  node specifies which of the parents  $Q(1 : B)$  to actually use.

Although the above problem may seem contrived, **identity uncertainty** is a widespread problem in many areas, such as citation analysis, credit card histories, and object tracking (see Section 4.6.5). In particular, the problem of **entity resolution** or **record linkage** — which refers to the task of mapping particular strings (such as names) to particular objects (such as people) — is a whole field of research (see e.g., [https://en.wikipedia.org/wiki/Record\\_linkage](https://en.wikipedia.org/wiki/Record_linkage) for an overview and [SHF15] for a Bayesian approach).

#### 4.6.4 Undirected relational PGMs

We can create **relational UGMs** in a manner which is analogous to relational DGMs (Section 4.6.3). This is particularly useful in the discriminative setting, for the same reasons that undirected CRFs are preferable to conditional DGMs (see Section 4.4).

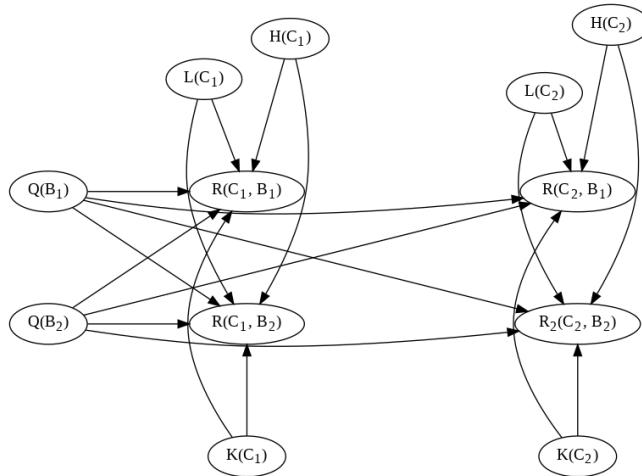


Figure 4.51: An extension of the book review RPM to handle identity uncertainty about which book a given customer is actually reviewing. The  $R(c, b)$  node now depends on all books, since we don't know which one is being referred to. We can select one of these parents based on the mapping specified by the user's library,  $L(c)$ .

#### 4.6.4.1 Collective classification

As an example of a relational UGM, suppose we are interested in the problem of classifying web pages of a university into types (e.g., student, professor, admin, etc.) Obviously we can do this based on the contents of the page (e.g., words, pictures, layout, etc.) However, we might also suppose there is information in the hyper-link structure itself. For example, it might be likely for students to cite professors, and professors to cite other professors, but there may be no links between admin pages and students / professors. When faced with a web page whose label is ambiguous, we can bias our estimate based on the estimated labels of its neighbors, as in a CRF. This process is known as **collective classification** (see e.g., [Sen+08]). To specify the CRF structure for a web-graph of arbitrary size and shape, we just specify a template graph and potential functions, and then unroll the template appropriately to match the topology of the web, making use of parameter tying.

#### 4.6.4.2 Markov logic networks

One particularly popular way of specifying relational UGMs is to use **first-order logic** rather than a graphical description of the template. The result is known as a **Markov logic network** [RD06; Dom+06; DL09].

For example, consider the sentences “Smoking causes cancer” and “If two people are friends, and one smokes, then so does the other”. We can write these sentences in first-order logic as follows:

$$\forall x. Sm(x) \implies Ca(x) \quad (4.176)$$

$$\forall x. \forall y. Fr(x, y) \wedge Sm(x) \implies Sm(y) \quad (4.177)$$

where  $Sm$  and  $Ca$  are predicates, and  $Fr$  is a relation.

	Fr(A,A)	Fr(B,B)	Fr(B,A)	Fr(A,B)	Sm(A)	Sm(B)	Ca(A)	Ca(B)
1	1	1	0	1	1	1	1	1
2	1	1	0	1	1	0	0	0
3	1	1	0	1	1	1	0	1
4								
5								

Table 4.5: Some possible joint instantiations of the 8 variables in the smoking example.

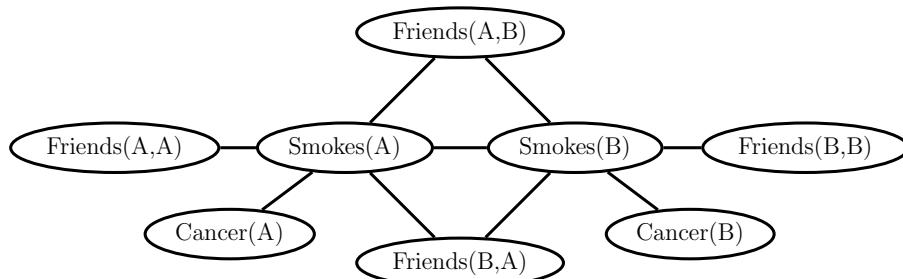


Figure 4.52: An example of a ground Markov logic network represented as a pairwise MRF for 2 people. Adapted from Figure 2.1 from [DL09]. Used with kind permission of Pedro Domingos.

It is convenient to write all formulas in **conjunctive normal form** (CNF), also known as **clausal form**. In this case, we get

$$\neg Sm(x) \vee Ca(x) \quad (4.178)$$

$$\neg Fr(x,y) \vee \neg Sm(x) \vee Sm(y) \quad (4.179)$$

The first clause can be read as “Either  $x$  does not smoke or he has cancer”, which is logically equivalent to Equation (4.176). (Note that in a clause, any unbound variable, such as  $x$ , is assumed to be universally quantified.)

Suppose there are just two objects (people) in the world, Anna and Bob, which we will denote by **constant symbols**  $A$  and  $B$ . We can then create 8 binary random variables  $Sm(x)$ ,  $Ca(x)$ , and  $Fr(x,y)$  for  $x, y \in \{A, B\}$ . This defines  $2^8$  **possible worlds**, some of which are shown in Table 4.5.<sup>14</sup>

Our goal is to define a probability distribution over these joint assignments. We can do this by creating a UGM with these variables, and adding a potential function to capture each logical rule or constraint. For example, we can encode the rule  $\neg Sm(x) \vee Ca(x)$  by creating a potential function

14. Note that we have not encoded the fact that  $Fr$  is a symmetric relation, so  $Fr(A,B)$  and  $Fr(B,A)$  might have different values. Similarly, we have the “degenerate” nodes  $Fr(A)$  and  $Fr(B)$ , since we did not enforce  $x \neq y$  in Equation (4.177). (If we add such constraints, then the model compiler, which generates the ground network, should avoid creating redundant nodes.)

1     $\Psi(Sm(x), Ca(x))$ , where we define  
2

$$\Psi(Sm(x), Ca(x)) = \begin{cases} 1 & \text{if } \neg Sm(x) \vee Ca(x) = T \\ 0 & \text{if } \neg Sm(x) \vee Ca(x) = F \end{cases} \quad (4.180)$$

3  
4    The result is the UGM in Figure 4.52.

5    The above approach will assign non-zero probability to all logically valid worlds. However, logical  
6 rules may not always be true. For example, smoking does not always cause cancer. We can relax the  
7 hard constraints by using non-zero potential functions. In particular, we can associate a weight with  
8 each rule, and thus get potentials such as  
9

$$\Psi(Sm(x), Ca(x)) = \begin{cases} e^w & \text{if } \neg Sm(x) \vee Ca(x) = T \\ e^0 & \text{if } \neg Sm(x) \vee Ca(x) = F \end{cases} \quad (4.181)$$

10  
11    where the value of  $w > 0$  controls strongly we want to enforce the corresponding rule.

12    The overall joint distribution has the form  
13

$$p(\mathbf{x}) = \frac{1}{Z(\mathbf{w})} \exp\left(\sum_i w_i n_i(\mathbf{x})\right) \quad (4.182)$$

14  
15    where  $n_i(\mathbf{x})$  is the number of instances of clause  $i$  which evaluate to true in assignment  $\mathbf{x}$ .

16    Given a grounded MLN model, we can then perform inference using standard methods. Of course,  
17 the ground models are often extremely large, so more efficient inference methods, which avoid creating  
18 the full ground model (known as **lifted inference**), must be used. See [DL09; KNP11] for details.

19    One way to gain tractability is to relax the discrete problem to a continuous one. This is the  
20 basic idea behind **hinge-loss MRFs** [Bac+15b], which support exact inference using scalable convex  
21 optimization. There is a template language for this model family known as **probabilistic soft logic**,  
22 which has a similar “flavor” to MLN, although it is not quite as expressive.

23    Recently MLNs have been combined with DL in various ways. For example, [Zha+20f] uses graph  
24 neural networks for inference. And [WP18] uses MLNs for evidence fusion, where the noisy predictions  
25 come from DNNs trained using weak supervision.

26    Finally, it is worth noting one subtlety which arises with undirected models, namely that the size  
27 of the unrolled model, which depends on the number of objects in the universe, can affect the results  
28 of inference, even if we have no data about the new objects. For example, consider an undirected  
29 chain of length  $T$ , with  $T$  hidden nodes  $z_t$  and  $T$  observed nodes  $y_t$ ; call this model  $M_1$ . Now suppose  
30 we double the length of the chain to  $2T$ , without adding more evidence; call this model  $M_2$ . We find  
31 that  $p(z_t|y_{1:T}, M_1) \neq p(z_t|y_{1:T}, M_2)$ , for  $t = 1 : T$ , even though we have not added new information,  
32 due to the different partition functions. This does not happen with a directed chain, because the  
33 newly added nodes can be marginalized out without affecting the original nodes, since the model is  
34 locally normalized and therefore modular. See [JBB09; Poo+12] for further discussion.

35

#### 36 4.6.5 Open-universe probability models

37    In Section 4.6.3, we discussed relational probability models, as well as the topic of identity uncertainty.  
38    However, we also implicitly made a **closed world assumption**, namely that the set of all objects is  
39 fixed and specified ahead of time. In many real world problems, this is an unrealistic assumption.

40

For example, in Section 29.9.3.5, we discuss the problem of tracking an unknown number of objects over time. As another example, consider the problem of enforcing the UN Comprehensive Nuclear Test Ban Treaty (CTBT). This requires monitoring seismic events, and determining if they were caused by nature or man-made explosions. Thus the number of objects of each type, as well as their source, is uncertain [ARS13].

As another (more peaceful) example, suppose we want to perform **citation matching**, in which we want to know whether to cite an arxiv version of a paper or the version on some conference website. Are these the same object? It is often hard to tell, since the titles and author might be the same, yet the content may have been updated. It is often necessary to use subtle cues, such as the date stored in the meta-data, to infer if the two “textual measurements” refer to the same underlying object (paper) or not [Pas+02].

In problems such as these, the number of objects of each type, as well as their relationships, is uncertain. This requires the use of **open-universe probability models** or **OUPM**, which can generate new objects as well as their properties [Rus15; MR10; LB19]. The first formal language for OUPMs was **BLOG** [Mil+05], which stands for “Bayesian LOGic”. This used a general purpose, but slow, MCMC inference scheme to sample over possible worlds of variable size and shape. [Las08; LLC20] describes another open-universe modeling language called **multi-entity Bayesian networks**.

Very recently, Facebook has released the **Bean Machine** library, available at <https://beanmachine.org/>, which supports more efficient inference in OUPMs. Details can be found in [Teh+20], as well as their blog post.<sup>15</sup>

#### 4.6.6 Programs as probability models

OUPMs, discussed in Section 4.6.5, let us define probability models over complex dynamic state spaces of unbounded and variable size. The set of possible worlds correspond to objects and their attributes and relationships. Another approach is to use a **probabilistic programming language** or **PPL**, in which we define the set of possible worlds as the set of **execution traces** generated by the program when it is endowed with a random choice mechanism. (This is a **procedural approach** to the problem, whereas OUPMs are a **declarative approach**.)

The difference between a probabilistic programming language and a standard one was described in [Gor+14] as follows: “Probabilistic programs are usual functional or imperative programs with two added constructs: (1) the ability to draw values at random from distributions, and (2) the ability to condition values of variables in a program via observation”. The former is a way to define  $p(\mathbf{z}, \mathbf{y})$ , and the latter is the same as standard Bayesian conditioning  $p(\mathbf{z}|\mathbf{y})$ .

Some recent examples of PPLs include **Gen** [CT+19], **Pyro** [Bin+19] and **Turing** [GXG18]. Inference in such models is often based on SMC, which we discuss in Chapter 13. For more details on PPLs, see e.g. [Mee+18].

### 4.7 Structural causal models

While probabilities encode our beliefs about a static world, causality tells us whether and how probabilities change when the world changes, be it by intervention or by act of imagination. —

<sup>15</sup> See <https://tinyurl.com/2svy5tmh>.

1 Judea Pearl [PM18b].  
2

3  
4 In this section, we discuss how we can use directed graphical model notation to represent **causal**  
5 **models**. We discuss causality in greater detail in Chapter 36, but we introduce some basic ideas and  
6 notation here, since it is foundational material that we will need in other parts of the book.

7 The core idea behind causal models is to create a mechanistic model of the world in which we  
8 can reason about the effects of local changes. The canonical example is an electronic circuit: we  
9 can predict the effects of any action, such as “knocking out” a particular transistor, or changing the  
10 resistance level of a wire, by modifying the circuit locally, and then “re-running” it from the same  
11 initial conditions.

12 We can generalize this idea to create a **structural causal models** or **SCM** [PGJ16], also called  
13 **functional causal model** [Sch19]. An SCM is a triple  $\mathcal{M} = (\mathcal{U}, \mathcal{V}, \mathcal{F})$ , where  $\mathcal{U} = \{U_i : i = 1 : N\}$   
14 is a set of unexplained or **exogenous** “noise” variables, which are passed as input to the model,  
15  $\mathcal{V} = \{V_i : i = 1 : N\}$  is a set of **endogeneous** variables that are part of the model itself, and  
16  $\mathcal{F} = \{f_i : i = 1 : N\}$  is a set of deterministic functions of the form  $V_i = f_i(V_{\text{pa}_i}, U_i)$ , where  $\text{pa}_i$  are the  
17 parents of variable  $i$ , and  $U_i \in \mathcal{U}$  are the external inputs. We assume the equations can be structured  
18 in a **recursive** way, so the dependency graph of nodes given their parents is a DAG. Finally, we  
19 assume our model is **causally sufficient**, which means that  $\mathcal{V}$  and  $\mathcal{U}$  are all of the causally relevant  
20 factors (although they may not all be observed). This is called the “**causal Markov assumption**”.

21 Of course, a model typically cannot represent all the variables that might influence observations or  
22 decisions. After all, models are *abstractions* of reality. The variables that we choose not to model  
23 explicitly in a functional way can be lumped into the unmodeled exogenous terms. To represent  
24 our ignorance about these terms, we can use a distribution  $p(\mathcal{U})$  over their values. By “pushing”  
25 this external noise through the deterministic part of the model, we induce a distribution over the  
26 endogeneous variables,  $p(\mathcal{V})$ , as in a probabilistic graphical model. However, SCMs make stronger  
27 assumptions than PGMs.

28 We usually assume  $p(\mathcal{U})$  is factorized (i.e., the  $U_i$  are independent); this is called a **Markovian**  
29 **SCM**. If the exogeneous noise terms are not independent, it would break the assumption that  
30 outcomes can be determined locally using deterministic functions. If there are believed to be  
31 dependencies between some of the  $U_i$ , we can add extra hidden parents to represent this; this is often  
32 depicted as a bidirected or undirected edge connecting the  $U_i$ , and is known as a **semi-Markovian**  
33 **SCM**.

34

35 **4.7.1 Example: causal impact of education on wealth**  
36

37 We now give a simple example of an SCM, based on [PM18b, p276]. Suppose we are interested in  
38 the causal effect of education on wealth. Let  $X$  represent the level of education of a person (on  
39 some numeric scale, say 0 = high school, 1 = college, 2 = graduate school), and  $Y$  represent their  
40 wealth (at some moment in time). In some cases we might expect that increasing  $X$  would increase  $Y$   
41 (although it of course depends on the nature of the degree, the nature of the job, etc). Thus we add  
42 an edge from  $X$  to  $Y$ . However, getting more education can cost a lot of money (in certain countries),  
43 which is a potentially confounding factor on wealth. Let  $Z$  be the debt incurred by a person based  
44 on their education. We add an edge from  $X$  to  $Z$  to reflect the fact that larger  $X$  means larger  $Z$  (in  
45 general), and we add an edge from  $Z$  to  $Y$  to reflect that larger  $Z$  means lower  $Y$  (in general).

46 We can represent our structural assumptions graphically as shown in Figure 4.53b(a). The  
47

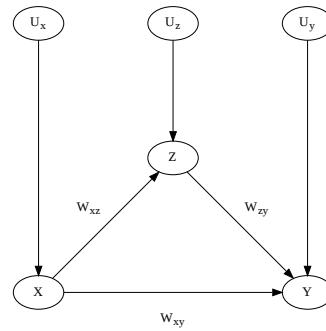
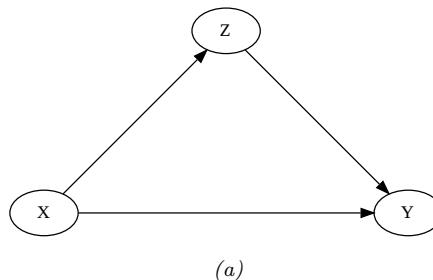


Figure 4.53: (a) PGM for modeling relationship between salary, education and debt. (b) Corresponding SCM.

corresponding SCM has the form:

$$X = f_X(U_x) \quad (4.183)$$

$$Z = f_Z(X, U_z) \quad (4.184)$$

$$Y = f_Y(X, Z, U_y) \quad (4.185)$$

for some set of functions  $f_x, f_y, f_z$ , and some prior distribution  $p(U_x, U_y, U_z)$ . We can also explicitly represent the exogeneous noise terms as shown in Figure 4.53b; this makes clear our assumption that the noise terms are a priori independent. (We return to this point later.)

## 4.7.2 Structural equation models

A **structural equation model** [Bol89; BP13], also known as a **path diagram**, is a special case of a structural causal model in which all the functional relationships are linear, and the prior on the noise terms is Gaussian. SEMs are widely used in economics and social science, due to the fact that they have a causal interpretation, yet they are computationally tractable.

For example, let us make an SEM version of our education example. We have

$$X = U_x \quad (4.186)$$

$$Z = c_z + w_{xz}X + U_z \quad (4.187)$$

$$Y = c_y + w_{xy}X + w_{zy}Z + U_y \quad (4.188)$$

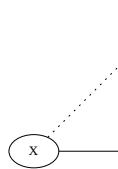
If we assume  $p(U_x) = \mathcal{N}(U_x|0, \sigma_x^2)$ ,  $p(U_z) = \mathcal{N}(U_z|0, \sigma_z^2)$ , and  $p(U_y) = \mathcal{N}(U_y|0, \sigma_y^2)$ , then the model can be converted to the following Gaussian DGM:

$$p(X) = \mathcal{N}(X|\mu_x, \sigma_x^2) \quad (4.189)$$

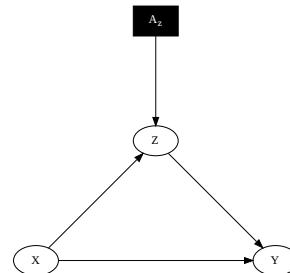
$$p(Z|X) = \mathcal{N}(Z|c_z + w_{xz}X, \sigma_z^2) \quad (4.190)$$

$$p(Y|X, Z) = \mathcal{N}(Y|c_y + w_{xy}X + w_{zy}Z, \sigma_y^2) \quad (4.191)$$

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13



(a)



(b)

14 *Figure 4.54: An SCM in which we intervene on Z. (a) Hard intervention, in which we clamp Z and thus cut  
15 its incoming edges (shown as dotted). (b) Soft intervention, in which we change Z’s mechanism. The square  
16 node is an “action” node, using the influence diagram notation from Section 34.2.*

17  
18  
19  
20 We can relax the linearity assumption, to allow arbitrarily flexible functions, and relax the Gaussian  
21 assumption, to allow any noise distribution. The resulting “nonparametric SEMs” are equivalent to  
22 structural causal models. (For a more detailed comparison between SEMs and SCMs, see [Pea12;  
23 BP13; Shi00].)

### 24 25 4.7.3 Do operator and augmented DAGs 26

27 One of the main advantages of SCMs is that they let us predict the effect of **interventions**, which  
28 are actions that change one or more local mechanisms. A simple intervention is to force a variable to  
29 have a given value, e.g., we can force a gene to be “on” or “off”. This is called a **perfect intervention**  
30 and is written as  $\text{do}(X_i = x_i)$ , where we have introduced new notation for the “**do**” operator (as  
31 in the verb “to do”). This notation means we actively clamp variable  $X_i$  to value  $x_i$  (as opposed to  
32 just observing that it has this value). Since the value of  $X_i$  is now independent of its usual parents,  
33 we should “cut” the incoming edges to node  $X_i$  in the graph. This is called the “**graph surgery**”  
34 operation.

35 In Figure 4.54a we illustrate this for our education SCM, where we force  $Z$  to have a given value.  
36 For example, we may set  $Z = 0$ , by paying off everyone’s student debt. Note that  $p(X|\text{do}(Z = z)) \neq$   
37  $p(X|Z = z)$ , since the intervention changes the model. For example, if we see someone with a debt of  
38 0, we may infer that they probably did not get higher education, i.e.,  $p(X \geq 1|Z = 0)$  is small; but if  
39 we pay off everyone’s college loans, then observing someone with no debt in this modified world should  
40 not change our beliefs about whether they got higher education, i.e.,  $p(X \geq 1|\text{do}(Z = 0)) = p(X \geq 1)$ .

41 In more realistic scenarios, we may not be able to set a variable to a specific value, but we may  
42 be able to change it from its current value in some way. For example, we may be able to reduce  
43 everyone’s debt by some fixed amount, say  $\Delta = -10,000$ . Thus we replace  $Z = f_Z(X, U_z)$  with  
44  $Z = f'_z(Z, U_z)$ , where  $f'_z(Z, U_z) = f_z(Z, U_z) + \Delta$ . This is called an **additive intervention**.

45 To model this kind of scenario, we can add create an **augmented DAG**, in which every variable is  
46 augmented with an additional parent node, representing whether or not the variable’s mechanism is  
47

Level	Activity	Questions	Examples
1:Association. $p(Y a)$	Seeing	How would seeing $A$ change my belief in $Y$ ?	Someone took aspirin, how likely is it their headache will be cured?
2:Intervention. $p(Y \text{do}(a))$	Doing	What if I do $A$ ?	If I take aspirin, will my headache be cured?
3:Counterfactuals. $p(Y^a \text{do}(a'), y')$	Imagining	Was it $A$ that caused $Y$ ?	Would my headache be cured had I not taken aspirin?

Table 4.6: Pearl's causal hierarchy. Adapted from Table 1 of [Pea19].

changed in some way [Daw02; Daw15; CPD17]. These extra variables are represented by square nodes, and correspond to decision variables or actions, as in the influence diagram formalism (Section 34.2). The same formalism is used in MDPs for reinforcement learning (see Section 34.5).

We give an example of this in Figure 4.54b, where we add the  $A_z \in \{0, 1\}$  node to specify whether we use the debt reduction policy or not. The modified mechanism for  $Z$  becomes

$$Z = f'_Z(X, U_x, A_z) = \begin{cases} f_Z(X, U_x) & \text{if } A_z = 0 \\ f_Z(X, U_x) + \Delta & \text{if } A_z = 1 \end{cases} \quad (4.192)$$

With this new definition, conditioning on the effects of an action can be performed using standard probabilistic inference. That is,  $p(Q|\text{do}(A_z = a), E = e) = p(Q|A_z = a, E = e)$ , where  $Q$  is the query (e.g., the event  $X \geq 1$ ) and  $E$  are the (possibly empty) evidence variables. This is because the  $A_z$  node has no parents, so it has no incoming edges to cut when we clamp it.

Although the augmented DAG allows us to use standard notation (no explicit do operators) and inference machinery, the use of “surgical” interventions, which delete incoming edges to a node that is set to a value, results in a simpler graph, which can simplify many calculations, particularly in the non-parametric setting (see [Pea09b, p361] for a discussion). It is therefore a useful abstraction, even if it is less general than the augmented DAG approach.

#### 4.7.4 Counterfactuals

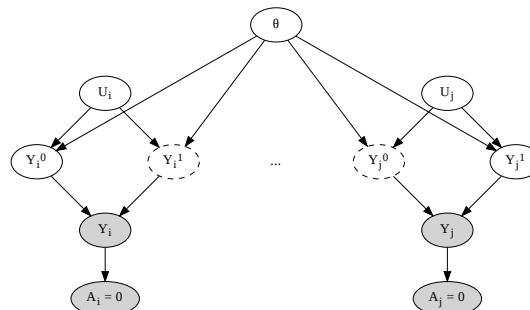
So far we have been focused on predicting the **effects of causes**, so we can choose the optimal action (e.g., if I have a headache, I have to decide should I take an aspirin or not). This can be tackled using standard techniques from Bayesian decision theory, as we have seen (see [Daw00; Daw15; LR19; Roh21; DM22] for more details).

Now suppose we are interested in the **causes of effects**. For example, suppose I took the aspirin and my headache did go away. I might be interested in the **counterfactual question** “if I had not taken the aspirin, would my headache have gone away anyway?”. This kind of reasoning is crucial for legal reasoning (see e.g., [DMM17]), as well as for tasks like explainability and fairness.

Counterfactual reasoning requires strictly more assumptions than reasoning about interventions (see e.g., [DM22]). Indeed, Judea Pearl has proposed what he calls the **causal hierarchy** [Pea09b; PGJ16; PM18b], which has three levels of analysis, each more powerful than the last, but each making stronger assumptions. See Table 4.6 for a summary.

In counterfactual reasoning, we want to answer questions of the type  $p(Y^{a'}|\text{do}(a), y)$ , which is read as: “what is the probability distribution over outcomes  $Y$  if I were to do  $a'$ , given that I have already

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13



14 *Figure 4.55: Illustration of the potential outcomes framework as a SCM. The nodes with dashed edges are*  
 15 *unobserved. In this example, for unit 1, we select action  $A_1 = 0$  and observe  $Y_1 = Y_1^0 = y_1$ , whereas for unit*  
 16 *2, we select action  $A_2 = 1$  and observe  $Y_2 = Y_2^1 = y_2$ .*

17  
18  
19

20 done  $a$  and observed outcome  $y$ . (We can also condition on any other evidencee that was observed,  
 21 such as covariates  $\mathbf{x}$ .) The quantity  $Y^{a'}$  is often called a **potential outcome** [Rub74], since it is  
 22 the outcome that would occur in a hypothetical world in which you did  $a'$  instead of  $a$ . (Note that  
 23  $p(Y^{a'} = y)$  is equivalent to  $p(Y = y|\text{do}(a'))$ , and is an interventional prediction, not a counterfactual  
 24 one.)

25 The assumptions behind the potential outcomes framework can be clearly expressed using a  
 26 structural causal model. We illustrate this in Figure 4.55 for a simple case where there are two  
 27 possible actions. We see that we have a set of “**units**”, such as individual patients, indexed by  
 28 subscripts. Each unit is associated with a hidden exogeneous random noise source,  $U_i$ , that captures  
 29 everything that is unique about that unit. This noise gets deterministically mapped to two potential  
 30 outcomes,  $Y_i^0$  and  $Y_i^1$ , depending on which action is taken. For any given unit, we only get to observe  
 31 one of the outcomes, namely the one corresponding to the action that was actually chosen. In the  
 32 figure, for unit 1, we chose action  $A_1 = 0$ , so we get to see  $Y_1^0 = y_1$ , whereas for unit 2, we chose  
 33 action  $A_2 = 1$ , so we get to see  $Y_2^1 = y_2$ . The fact that we cannot simultaneously see both outcomes  
 34 for the same unit is called the “**fundamental problem of causal inference**” [Hol86].

35 We will assume the noise sources are independent, which is known as the “stable unit treatment  
 36 value assumption” or **SUTVA**. (This would not be true if the treatment on person  $j$  could somehow  
 37 affect the outcome of person  $i$ , e.g., due to spreading disease or information between  $i$  and  $j$ .) We  
 38 also assume that the deterministic mechanisms that map noise to outcomes are the same across  
 39 all units (represented by the shared parameter vector  $\theta$  in Figure 4.55). We need to make one final  
 40 assumption, namely that the exogeneous noise is not affected by our actions. (This is a formalization  
 41 of the assumption known as “all else being equal”, or (in legal terms) “**ceteris paribus**”).

42 With the above assumptions, we can predict what the outcome *for an individual unit* would have  
 43 been in the alternative universe where we picked the other action. The procedure is as follows. First  
 44 we perform **abduction** using SCM  $G$ , to infer  $p(U_i|A_i = a, Y_i = y_i)$ , which is the posterior over  
 45 the latent factors for unit  $i$  given the observed evidence in the actual world. Second we perform  
 46 **intervention**, in which we modify the causal mechanisms of  $G$  by replacing  $A_i = a$  with  $A_i = a'$  to  
 47

1 get  $G_{a'}$ . Third we perform **prediction**, in which we propagate the distribution of the latent factors,  
2  $p(U_i|A_i = a, Y_i = y_i)$ , through the modified SCM  $G_{a'}$  to get  $p(Y_i^{a'}|A_i = a, Y_i = y_i)$ .

3 In Figure 4.55, we see that we have two copies of every possible outcome variable, to represent  
4 the set of possible worlds. Of course, we only get to see one such world, based on the actions that  
5 we actually took. More generally, a model in which we “clone” all the deterministic variables, with  
6 the noise being held constant between the two branches of the graph for the same unit, is called a  
7 **twin network** [Pea09b]. We will see a more practical example in Section 29.12.7, where we discuss  
8 assessing the counterfactual causal impact of an intervention in a time series. (See also [RR11; RR13],  
9 who propose a related formalism known as **single world intervention graph** or **SWIG**.)

10 We see from the above that the potential outcomes framework is mathematically equivalent to  
11 structural causal models, but does not use graphical model notation. This has led to heated debate  
12 between the founders of the two schools of thought.<sup>16</sup> The SCM approach is more popular in  
13 computer science (see e.g., [PJS17; Sch19; Sch+21b]), and the PO approach is more popular in  
14 economics (see e.g. [AP09; Imb19]). Modern textbooks on causality usually use both formalisms (see  
15 e.g., [HR20a; Nea20]).

171819202122232425262728293031323334353637383940414243

44 16. The potential outcomes framework is based on the work of Donald Rubin, and others, and is therefore sometimes  
45 called the **Rubin Causal Model** (see e.g., [https://en.wikipedia.org/wiki/Rubin\\_causal\\_model](https://en.wikipedia.org/wiki/Rubin_causal_model)). The structural  
46 causal models framework is based on the work of Judea Pearl and others. See e.g., <http://causality.cs.ucla.edu/blog/index.php/2012/12/03/judea-pearl-on-potential-outcomes/> for a discussion of the two.

47



# 5 Information theory

Machine learning is fundamentally about **information processing**. But what do we mean by “information”? We discuss this in Section 5.1–Section 5.3. We then go on to briefly discuss two main applications of information theory. The first application is **data compression** or **source coding**, which is the problem of removing redundancy from data so it can be represented more compactly, either in a lossless way (e.g., ZIP files) or a lossy way (e.g., MP3 files). See Section 5.4 for details. The second application is **error correction** or **channel coding**, which means encoding data in such a way that it is robust to errors when sent over a noisy channel, such as a telephone line or a satellite link. See Section 5.5 for details.

It turns out that methods for data compression and error correction both rely on having an accurate probabilistic model of the data. For compression, a probabilistic model is needed so the sender can assign shorter **codewords** to data vectors which occur most often, and hence save space. For error correction, a probabilistic model is needed so the receiver can infer the most likely source message by combining the received noisy message with a prior over possible messages.

It is clear that probabilistic machine learning is useful for information theory. However, information theory is also useful for machine learning. Indeed, we have seen that Bayesian machine learning is about representing and reducing our uncertainty, and so is fundamentally about information. In Section 5.6.2, we explore this direction in more detail, where we discuss the information bottleneck.

For more information on information theory, see e.g., [Mac03; CT06].

## 5.1 KL divergence

*This section is written with Alex Alemi.*

To discuss information theory, we need some way to measure or quantify information itself. Let’s say we start with some distribution describing our degrees of belief about a random variable, call it  $q(x)$ . We then want to update our degrees of belief to some new distribution  $p(x)$ , perhaps because we’ve taken some new measurements or merely thought about the problem a bit longer. What we seek is a mathematical way to quantify the magnitude of this update, which we’ll denote  $I[p\|q]$ . What sort of criteria would be reasonable for such a measure? We discuss this issue below, and then define a quantity that satisfies these criteria.

1 2 **5.1.1 Desiderata**

3 For simplicity, imagine we are describing a distribution over  $N$  possible events. In this case, the  
4 probability distribution  $q(\mathbf{x})$  consists of  $N$  non-negative real numbers that add up to 1. To be even  
5 more concrete, imagine we are describing the random variable representing the suit of the next card  
6 we'll draw from a deck:  $S \in \{\clubsuit, \spadesuit, \heartsuit, \diamondsuit\}$ . Imagine we initially believe the distributions over suits to  
7 be uniform:  $q = [\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4}]$ . If our friend told us they removed all of the red cards we could update  
8 to:  $q' = [\frac{1}{2}, \frac{1}{2}, 0, 0]$ . Alternatively, we might believe some diamonds changed into clubs and want to  
9 update to  $q'' = [\frac{3}{8}, \frac{2}{8}, \frac{2}{8}, \frac{1}{8}]$ . Is there a good way to quantify *how much* we've updated our beliefs?  
10 Which is a larger update:  $q \rightarrow q'$  or  $q \rightarrow q''$ ?

11 It seems desireable that any useful such measure would satisfy the following properties:  
12

- 13 1. *continuous* in its arguments: If we slightly perturb either our starting or ending distribution,  
14 it should similarly have a small effect on the magnitude of the update. For example:  $I[p \parallel \frac{1}{4} +$   
15  $\epsilon, \frac{1}{4}, \frac{1}{4}, \frac{1}{4} - \epsilon]$  should be close to  $I[p \parallel q]$  for small  $\epsilon$ , where  $q = [\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4}]$ .
- 16 2. *non-negative*:  $I[p \parallel q] \geq 0$  for all  $p(\mathbf{x})$  and  $q(\mathbf{x})$ . The magnitude of our updates are non-negative.  
17
- 18 3. *permutation invariant*: The magnitude of the update should not depend on the order we choose for  
19 the elements of  $\mathbf{x}$ . For example, it shouldn't matter if I list my probabilities for the suits of cards  
20 in the order  $\clubsuit, \spadesuit, \heartsuit, \diamondsuit$  or  $\clubsuit, \diamondsuit, \heartsuit, \spadesuit$ , if I keep the order consistent across all of the distributions,  
21 I should get the same answer. For example:  $I[a, b, c, d \parallel e, f, g, h] = I[a, d, c, b \parallel e, h, g, f]$ .  
22
- 23 4. *monotonic* for uniform distributions: While it's hard to say how large the updates in our beliefs  
24 are in general, there are some special cases for which we have a strong intuition. If our beliefs  
25 update from a uniform distribution on  $N$  elements to one that is uniform in  $N'$  elements, the  
26 information gain should be an increasing function of  $N$  and a decreasing function of  $N'$ . For  
27 instance changing from a uniform distribution on all four suits  $[\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4}]$  (so  $N = 4$ ) to only one  
28 suit, such as all clubs,  $[1, 0, 0, 0]$  where  $N' = 1$ , is a larger update than if I only updated to the  
29 card being black,  $[\frac{1}{2}, \frac{1}{2}, 0, 0]$  where  $N' = 2$ .
- 30 5. satisfy a natural *chain rule*: So far we've been describing our beliefs in what will happen on the next  
31 card draw as a single random variable representing the suit of the next card ( $S \in \{\clubsuit, \spadesuit, \heartsuit, \diamondsuit\}$ ).  
32 We could equivalently describe the same physical process in two steps. First we consider the  
33 random variable representing the color of the card ( $C \in \{\blacksquare, \square\}$ ), which could be either black  
34 ( $\blacksquare = \{\clubsuit, \spadesuit\}$ ) or red ( $\square = \{\heartsuit, \diamondsuit\}$ ). Then, if we draw a red card we describe our belief that it is  $\heartsuit$   
35 versus  $\diamondsuit$ . If it was instead black we would assign beliefs to it being  $\clubsuit$  versus  $\spadesuit$ . We can convert  
36 any distribution over the four suits into this conditional factorization, for example:  
37

$$\text{p}(S) = \left[ \frac{3}{8}, \frac{2}{8}, \frac{2}{8}, \frac{1}{8} \right] \quad (5.1)$$

40 becomes

$$\text{p}(C) = \left[ \frac{5}{8}, \frac{3}{8} \right] \quad \text{p}(\{\clubsuit, \spadesuit\} | C = \blacksquare) = \left[ \frac{3}{5}, \frac{2}{5} \right] \quad \text{p}(\{\heartsuit, \diamondsuit\} | C = \square) = \left[ \frac{2}{3}, \frac{1}{3} \right]. \quad (5.2)$$

45 In the same way we could decompose our uniform distribution  $q$ . Obviously, for our measure of  
46 information to be of use the magnitude of the update needs to be the same regardless of how we  
47

choose to describe what is ultimately the same physical process. What we need is some way to relate what would be four different invocations of our information function:

$$I_S \equiv I[p(S)\|q(S)] \quad (5.3)$$

$$I_C \equiv I[p(C)\|q(C)] \quad (5.4)$$

$$I_{\blacksquare} \equiv I[p(\{\clubsuit, \spadesuit\}|C = \blacksquare)\|q(\{\clubsuit, \spadesuit\}|C = \blacksquare)] \quad (5.5)$$

$$I_{\square} \equiv I[p(\{\heartsuit, \diamondsuit\}|C = \square)\|q(\{\heartsuit, \diamondsuit\}|C = \square)]. \quad (5.6)$$

Clearly  $I_S$  should be some function of  $\{I_C, I_{\blacksquare}, I_{\square}\}$ . Our last desiderata is that the way we measure the magnitude of our updates will have  $I_S$  be a linear combination of  $I_C, I_{\blacksquare}, I_{\square}$ . In particular, we will require that they combine as a weighted linear combinations, with weights set by the probability that we would find ourselves in that branch according to the distribution  $p$ :

$$I_S = I_C + p(C = \blacksquare)I_{\blacksquare} + p(C = \square)I_{\square} = I_C + \frac{5}{8}I_{\blacksquare} + \frac{3}{5}I_{\square} \quad (5.7)$$

Stating this requirement more generally: If we partition  $\mathbf{x}$  into two pieces  $[\mathbf{x}_L, \mathbf{x}_R]$ , so that we can write  $p(\mathbf{x}) = p(\mathbf{x}_L)p(\mathbf{x}_R|\mathbf{x}_L)$  and similarly for  $q$ , the magnitude of the update should be

$$I[p(\mathbf{x})\|q(\mathbf{x})] = I[p(\mathbf{x}_L)\|q(\mathbf{x}_L)] + \mathbb{E}_{p(\mathbf{x}_L)} [I[p(\mathbf{x}_R|\mathbf{x}_L)\|q(\mathbf{x}_R|\mathbf{x}_L)]] . \quad (5.8)$$

Notice that this requirement *breaks the symmetry between our two distributions*: The right hand side asks us to take the expected conditional information gain with respect to the marginal, but we need to decide which of two marginals to take the expectation with respect to.

### 5.1.2 The KL divergence uniquely satisfies the desiderata

We will now define a quantity that is the only measure (up to a multiplicative constant) that satisfies the above desiderata. The **Kullback-Leibler divergence** or **KL divergence**, also known as the **information gain** or **relative entropy**, is defined as follows:

$$D_{\text{KL}}(p \parallel q) \triangleq \sum_{k=1}^K p_k \log \frac{p_k}{q_k}. \quad (5.9)$$

This naturally extends to continuous distributions:

$$D_{\text{KL}}(p \parallel q) \triangleq \int dx p(x) \log \frac{p(x)}{q(x)}. \quad (5.10)$$

Next we will verify that this definition satisfies all of our desiderata. (The proof that it is the unique measure which captures these properties can be found in e.g., [Hob69; Rén61].)

#### 5.1.2.1 Continuity of KL

One of our desiderata was that our measure of information gain should be continuous. The KL divergence is manifestly continuous in its arguments except potentially when  $p_k$  or  $q_k$  is zero. In the first case, notice that the limit as  $p \rightarrow 0$  is well behaved:

$$\lim_{p \rightarrow 0} p \log \frac{p}{q} = 0. \quad (5.11)$$

1 Taking this as the definition of the value of the integrand when  $p = 0$  will make it continuous there.  
2 Notice that we do have a problem however if  $q = 0$  in some place that  $p \neq 0$ . Our information  
3 gain requires that our original distribution of beliefs  $q$  has some support everywhere the updated  
4 distribution does. Intuitively it would require an infinite amount of information for us to update our  
5 beliefs in some outcome to change from being exactly 0 to some positive value.  
6

### 7 5.1.2.2 Non-negativity of KL divergence

8 In this section, we prove that the KL divergence as defined is always non-negative. We will make use  
9 of **Jensen's inequality**, which states that for any convex function  $f$ , we have that  
10

$$\underline{12} \quad f\left(\sum_{i=1}^n \lambda_i \mathbf{x}_i\right) \leq \sum_{i=1}^n \lambda_i f(\mathbf{x}_i) \quad (5.12)$$

13 where  $\lambda_i \geq 0$  and  $\sum_{i=1}^n \lambda_i = 1$ . This can be proved by induction, where the base case with  $n = 2$   
14 follows by definition of convexity.

15  
16 **Theorem 5.1.1.** (*Information inequality*)  $D_{\text{KL}}(p \parallel q) \geq 0$  with equality iff  $p = q$ .

17 *Proof.* We now prove the theorem, following [CT06, p28]. As we noted in the previous section, the  
18 KL divergence requires special consideration when  $p(x)$  or  $q(x) = 0$ , the same is true here. Let  
19  $A = \{x : p(x) > 0\}$  be the support of  $p(x)$ . Using the convexity of the log function and Jensen's  
20 inequality, we have that  
21

$$\underline{22} \quad -D_{\text{KL}}(p \parallel q) = -\sum_{x \in A} p(x) \log \frac{p(x)}{q(x)} = \sum_{x \in A} p(x) \log \frac{q(x)}{p(x)} \quad (5.13)$$

$$\underline{23} \quad \leq \log \sum_{x \in A} p(x) \frac{q(x)}{p(x)} = \log \sum_{x \in A} q(x) \quad (5.14)$$

$$\underline{24} \quad \leq \log \sum_{x \in \mathcal{X}} q(x) = \log 1 = 0 \quad (5.15)$$

25 Since  $\log(x)$  is a strictly concave function ( $-\log(x)$  is convex), we have equality in Equation (5.14) iff  
26  $p(x) = cq(x)$  for some  $c$  that tracks the fraction of the whole space  $\mathcal{X}$  contained in  $A$ . We have equality  
27 in Equation (5.15) iff  $\sum_{x \in A} q(x) = \sum_{x \in \mathcal{X}} q(x) = 1$ , which implies  $c = 1$ . Hence  $D_{\text{KL}}(p \parallel q) = 0$  iff  
28  $p(x) = q(x)$  for all  $x$ .  $\square$   
29

30 The non-negativity of KL divergence often feels as though it's one of the most useful results in  
31 Information Theory. It is a good result to keep in your back pocket. Anytime you can rearrange an  
32 expression in terms of KL divergence terms, since those are guaranteed to be non-negative, dropping  
33 them immediately generates a bound.  
34

### 35 5.1.2.3 KL divergence is invariant to reparameterizations

36 We wanted our measure of information to be invariant to permutations of the labels. The discrete  
37 form is manifestly permutation invariant as summations are. The KL divergence actually satisfies a  
38

1 much stronger property of reparameterization invariance. Namely, we can transform our random  
2 variable through an arbitrary invertible map and it won't change the value of the KL divergence.  
3

4 If we transform our random variable from  $x$  to some  $y = f(x)$  we know that  $p(x) dx = p(y) dy$  and  
5  $q(x) dx = q(y) dy$ . Hence the KL divergence remains the same for both random variables:

$$\underline{6} \quad D_{\text{KL}}(p(x) \parallel q(x)) = \int dx p(x) \log \frac{p(x)}{q(x)} = \int dy p(y) \log \left( \frac{p(y)}{q(y)} \left| \frac{dy}{dx} \right| \right) = D_{\text{KL}}(p(y) \parallel q(y)). \quad (5.16)$$

10 Because of this reparameterization invariance we can rest assured that when we measure the KL  
11 divergence between two distributions we are measuring something about the distributions and not the  
12 way we choose to represent the space in which they are defined. We are therefore free to transform  
13 our data into a convenient basis of our choosing, such as a Fourier bases for images, without affecting  
14 the result.

#### 16 5.1.2.4 Monotonicity for uniform distributions

17 Consider updating a probability distribution from a uniform distribution on  $N$  elements to a uniform  
18 distribution on  $N'$  elements. The KL divergence is:

$$\underline{20} \quad D_{\text{KL}}(p \parallel q) = \sum_k \frac{1}{N'} \log \frac{\frac{1}{N'}}{\frac{1}{N}} = \log \frac{N}{N'}, \quad (5.17)$$

23 or the log of the ratio of the elements before and after the update. This satisfies our monotonicity  
24 requirement.

25 We can interpret this result as follows: Consider finding an element of a sorted array by means of  
26 bisection. A well designed yes/no question can cut the search space in half. Measured in bits, the  
27 KL divergence tells us how many well designed yes/no questions are required on average to move  
28 from  $q$  to  $p$ .

#### 30 5.1.2.5 Chain rule for KL divergence

31 Here we show that the KL divergence satisfies a natural chain rule:

$$\underline{33} \quad D_{\text{KL}}(p(x, y) \parallel q(x, y)) = \int dx dy p(x, y) \log \frac{p(x, y)}{q(x, y)} \quad (5.18)$$

$$\underline{35} \quad = \int dx dy p(x, y) \left[ \log \frac{p(x)}{q(x)} + \log \frac{p(y|x)}{q(y|x)} \right] \quad (5.19)$$

$$\underline{38} \quad = D_{\text{KL}}(p(x) \parallel q(x)) + \mathbb{E}_{p(x)} [D_{\text{KL}}(p(y|x) \parallel q(y|x))]. \quad (5.20)$$

39 We can rest assured that we can decompose our distributions into their conditionals and the KL  
40 divergences will just add.

41 As a notational convenience, the **conditional KL divergence** is defined to be the expected value  
42 of the KL divergence between two conditional distributions:

$$\underline{44} \quad D_{\text{KL}}(p(y|x) \parallel q(y|x)) \triangleq \int dx p(x) \int dy p(y|x) \log \frac{p(y|x)}{q(y|x)}. \quad (5.21)$$

46 This allows us to drop many expectation symbols.

1 **5.1.3 Thinking about KL**

3 In this section, we discuss some qualitative properties of the KL divergence.

5 **5.1.3.1 Units of KL**

7 Above we said that the desiderata we listed determined the KL divergence up to a multiplicative  
8 constant. Because the KL divergence is logarithmic, and logarithms in different bases are the same  
9 up to a multiplicative constant, our choice of the base of the logarithm when we compute the KL  
10 divergence is a choice akin to choosing which units to measure the information in.

11 If the KL divergence is measured with the base-2 logarithm, it is said to have units of **bits**, short  
12 for “binary digits”. If measured using the natural logarithm as we normally do for mathematical  
13 convenience, it is said to be measured in **nats** for “natural units”.

14 To convert between the systems, we use  $\log_2 y = \frac{\log y}{\log 2}$ . Hence

$$\underline{16} \quad 1 \text{ bit} = \log 2 \text{ nats} \sim 0.693 \text{ nats} \tag{5.22}$$

$$\underline{17} \quad 1 \text{ nat} = \frac{1}{\log 2} \text{ bits} \sim 1.44 \text{ bits.} \tag{5.23}$$

20 **5.1.3.2 Asymmetry of the KL divergence**

22 The KL divergence is *not* symmetric in its two arguments. While many find this asymmetry confusing  
23 at first, we can see that the asymmetry stems from our requirement that we have a natural chain  
24 rule. When we decompose the distribution into its conditional, we need to take an expectation with  
25 respect to the variables being conditioned on. In the KL divergence we take this expectation with  
26 respect to the first argument  $p(x)$ . This breaks the symmetry between the two distributions.

27 At a more intuitive level, we can see that the information required to move from  $q$  to  $p$  is in general  
28 different than the information required to move from  $p$  to  $q$ . For example, consider the KL divergence  
29 between two Bernoulli distributions, the first with the probability of success given by 0.443 and the  
30 second with 0.975:

$$\underline{31} \quad D_{\text{KL}} = 0.975 \log \frac{0.975}{0.443} + 0.025 \log \frac{0.025}{0.557} = 0.692 \text{ nats} \sim 1.0 \text{ bits.} \tag{5.24}$$

33 So it takes 1 bit of information to update from a [0.443, 0.557] distribution to a [0.975, 0.025] Bernoulli  
34 distribution. What about the reverse?

$$\underline{36} \quad D_{\text{KL}} = 0.443 \log \frac{0.443}{0.975} + 0.557 \log \frac{0.557}{0.025} = 1.38 \text{ nats} \sim 2.0 \text{ bits,} \tag{5.25}$$

38 so it takes two bits, or twice as much information to move the other way. Thus we see that starting  
39 with a distribution that is nearly even and moving to one that is nearly certain takes about 1 bit of  
40 information, or one well designed yes/no question. To instead move us from near certainty in an  
41 outcome to something that is akin to the flip of a coin requires more persuasion.

43 **5.1.3.3 Minimizing forwards vs reverse KL**

45 The asymmetry of KL means that finding a  $p$  that is close to  $q$  by minimizing  $D_{\text{KL}}(p \parallel q)$  (also called  
46 the **inclusive KL**) gives different behavior than minimizing  $D_{\text{KL}}(q \parallel p)$  (also called the **exclusive**  
47

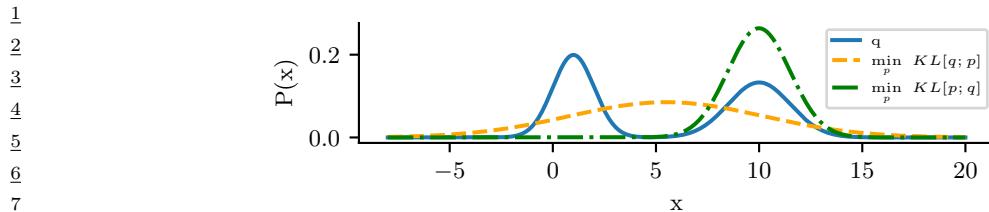


Figure 5.1: Demonstration of the mode-covering or mode-seeking behavior of KL divergence. The original distribution  $q$  is bimodal. When we minimize  $D_{\text{KL}}(q \parallel p)$ , then  $p$  covers the modes of  $q$  (orange). When we minimize  $D_{\text{KL}}(p \parallel q)$ , then  $p$  ignores some of the modes of  $q$  (green). Generated by [minimize\\_kl\\_divergence.ipynb](#).

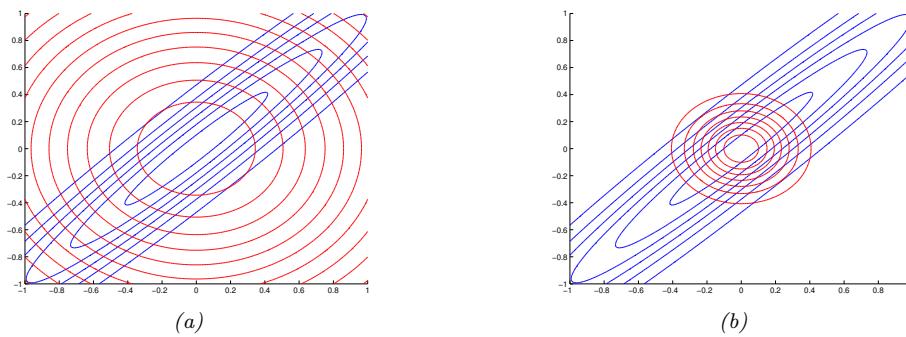


Figure 5.2: Illustrating forwards vs reverse KL on a symmetric Gaussian. The blue curves are the contours of the true distribution  $p$ . The red curves are the contours of a factorized approximation  $q$ . (a) Minimizing  $D_{\text{KL}}(p \parallel q)$ . (b) Minimizing  $D_{\text{KL}}(q \parallel p)$ . Adapted from Figure 10.2 of [Bis06]. Generated by [kl\\_pq\\_gauss.ipynb](#).

For example, consider the bimodal distribution  $q$  shown in blue in Figure 5.1, which we approximate with a unimodal Gaussian. To prevent  $D_{\text{KL}}(q \parallel p)$  from becoming infinite, we must have  $p > 0$  whenever  $q > 0$  (i.e.,  $p$  must have support everywhere  $q$  does), so  $p$  tends to **cover** both modes as it must be nonvanishing everywhere  $q$  is; this is called **mode-covering** or **zero-avoiding** behavior (orange curve).

By contrast, to prevent  $D_{\text{KL}}(p \parallel q)$  from becoming infinite, we must have  $p = 0$  whenever  $q = 0$ , which creates **mode-seeking** or **zero-forcing** behavior (green curve).

#### 5.1.3.4 Moment projection

Suppose we compute  $q$  by minimizing the forwards KL:

$$q = \underset{q}{\operatorname{argmin}} D_{\text{KL}}(p \parallel q) \quad (5.26)$$

This is called **M-projection**, or **moment projection** since the optimal  $q$  matches the moments of  $p$ ; this is called **moment matching**.

To see why, let us assume that  $q$  is an exponential family distribution of the form

$$q(\mathbf{x}) = h(\mathbf{x}) \exp[\boldsymbol{\eta}^T \mathcal{T}(\mathbf{x}) - \log Z(\boldsymbol{\eta})] \quad (5.27)$$

where  $\mathcal{T}(\mathbf{x})$  is the vector of sufficient statistics, and  $\boldsymbol{\eta}$  are the natural parameters. The first order optimality conditions are as follows:

$$\partial_{\eta_i} D_{\text{KL}}(p \parallel q) = -\partial_{\eta_i} \int_{\mathbf{x}} p(\mathbf{x}) \log q(\mathbf{x}) \quad (5.28)$$

$$= -\partial_{\eta_i} \int_{\mathbf{x}} p(\mathbf{x}) \log (h(\mathbf{x}) \exp[\boldsymbol{\eta}^T \mathcal{T}(\mathbf{x}) - \log Z(\boldsymbol{\eta})]) \quad (5.29)$$

$$= -\partial_{\eta_i} \int_{\mathbf{x}} p(\mathbf{x}) (\boldsymbol{\eta}^T \mathcal{T}(\mathbf{x}) - \log Z(\boldsymbol{\eta})) \quad (5.30)$$

$$= - \int_{\mathbf{x}} p(\mathbf{x}) \mathcal{T}_i(\mathbf{x}) + \mathbb{E}_{q(\mathbf{x})} [\mathcal{T}_i(\mathbf{x})] \quad (5.31)$$

$$= -\mathbb{E}_{p(\mathbf{x})} [\mathcal{T}_i(\mathbf{x})] + \mathbb{E}_{q(\mathbf{x})} [\mathcal{T}_i(\mathbf{x})] = 0 \quad (5.32)$$

where in the penultimate line we used the fact that the derivative of the log partition function yields the expected sufficient statistics, as shown in Equation (2.199). Hence the expected sufficient statistics (moments of the distribution) must match.

As an example, suppose the true target distribution  $p$  is a correlated 2d Gaussian,  $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$ , where

$$\boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \quad \boldsymbol{\Sigma} = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{12}^T & \Sigma_{22} \end{pmatrix} \quad \boldsymbol{\Lambda} = \begin{pmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{12}^T & \Lambda_{22} \end{pmatrix} \quad (5.33)$$

We will approximate this with a distribution  $q$  which is a product of two 1d Gaussians, i.e., a Gaussian with a diagonal covariance matrix:

$$q(\mathbf{x}|\mathbf{m}, \mathbf{V}) = \mathcal{N}(x_1|m_1, v_1)\mathcal{N}(x_2|m_2, v_2) \quad (5.34)$$

If we perform moment matching, the optimal  $q$  must therefore have the following form:

$$q(\mathbf{x}) = \mathcal{N}(x_1|\mu_1, \Sigma_{11})\mathcal{N}(x_2|\mu_2, \Sigma_{22}) \quad (5.35)$$

In Figure 5.2(a), we show the resulting distribution. We see that  $q$  covers (includes)  $p$ , but its support is too broad (under-confidence).

### 5.1.3.5 Information projection

Now suppose we compute  $q$  by minimizing the reverse KL:

$$q = \underset{q}{\operatorname{argmin}} D_{\text{KL}}(q \parallel p) \quad (5.36)$$

This is called **I-projection**, or **information projection**.

In the case of a 2d Gaussian, one can show [Mac03, p435] that the optimal solution has the form

$$q(\mathbf{x}) = \mathcal{N}(x_1|m_1, \Lambda_{11}^{-1})\mathcal{N}(x_2|m_2, \Lambda_{22}^{-1}) \quad (5.37)$$

$$m_1 = \mu_1 - \Lambda_{11}^{-1}\Lambda_{12}(m_2 - \mu_2) \quad (5.38)$$

$$m_2 = \mu_2 - \Lambda_{22}^{-1}\Lambda_{21}(m_1 - \mu_1) \quad (5.39)$$

To solve this fixed point equation, we can set  $m_i = \mu_i$ , so we have again captured the mean exactly. However, the posterior variance is too narrow, i.e., the approximate posterior is overconfident. This is shown in Figure 5.2(b). (Note, however, that minimizing the reverse KL does not always result in an overly compact approximation, as explained in [Tur+08].)

### 5.1.3.6 KL as expected weight of evidence

Imagine you have two different hypotheses you wish to select between, which we'll label  $P$  and  $Q$ . You collect some data  $D$ . Bayes' rule tells us how to update our beliefs in the hypotheses being correct:

$$\Pr(P|D) = \frac{\Pr(D|P)}{\Pr(D)} \Pr(P). \quad (5.40)$$

Normally this requires being able to evaluate the marginal likelihood  $\Pr(D)$ , which is difficult. If we instead consider the ratio of the probabilities for the two hypotheses:

$$\frac{\Pr(P|D)}{\Pr(Q|D)} = \frac{\Pr(D|P)}{\Pr(D|Q)} \frac{\Pr(P)}{\Pr(Q)}, \quad (5.41)$$

the marginal likelihood drops out. Taking the logarithm of both sides, and identifying the probability of the data under the model as the likelihood we find:

$$\log \frac{\Pr(P|D)}{\Pr(Q|D)} = \log \frac{p(D)}{q(D)} + \log \frac{\Pr(P)}{\Pr(Q)}. \quad (5.42)$$

The posterior log probability ratio for one hypothesis over the other is just our prior log probability ratio plus a term that I. J. Good called the **weight of evidence** [Goo85]  $D$  for hypothesis  $P$  over  $Q$ :

$$w[P/Q; D] \triangleq \log \frac{p(D)}{q(D)}. \quad (5.43)$$

With this interpretation, the KL divergence is the expected weight of evidence for  $P$  over  $Q$  given by each observation, provided  $P$  were correct. Thus we see that data will (on average) add rather than subtract evidence towards the correct hypothesis, since KL divergence is always non-negative in expectation (see Section 5.1.2.2).

### 5.1.4 Properties of KL

Below are some other useful properties of the KL divergence.

#### 5.1.4.1 Compression Lemma

An important general purpose result for the KL divergence is the Compression Lemma:

**Theorem 5.1.2.** *For any distributions  $P$  and  $Q$  with a well-defined KL divergence, and for any scalar function  $\phi$  defined on the domain of the distributions we have that:*

$$\mathbb{E}_P [\phi] \leq \log \mathbb{E}_Q [e^\phi] + D_{\text{KL}}(P \| Q). \quad (5.44)$$

1 2 3 *Proof.* We know that the KL divergence between any two distributions is non-negative. Consider a distribution of the form:

4  
5 
$$g(x) = \frac{q(x)}{\mathcal{Z}} e^{\phi(x)}. \quad (5.45)$$
  
6

7 where the *partition function* is given by:  
8

9  
10 
$$\mathcal{Z} = \int dx q(x) e^{\phi(x)}. \quad (5.46)$$
  
11

12 Taking the KL divergence between  $p(x)$  and  $g(x)$  and rearranging gives the bound:  
13

14 
$$D_{\text{KL}}(P \parallel G) = D_{\text{KL}}(P \parallel Q) - \mathbb{E}_P[\phi(x)] + \log(\mathcal{Z}) \geq 0. \quad (5.47)$$
  
15

□

18 One way to view the compression lemma is that it provides what is termed the Donsker-Varadhan  
19 variational representation of the KL divergence:  
20

21  
22 
$$D_{\text{KL}}(P \parallel Q) = \sup_{\phi} \mathbb{E}_P[\phi(x)] - \log \mathbb{E}_Q[e^{\phi(x)}]. \quad (5.48)$$
  
23

24 In the space of all possible functions  $\phi$  defined on the same domain as the distributions, assuming all  
25 of the values above are finite, the KL divergence is the supremum achieved. For any fixed function  
26  $\phi(x)$ , the right hand side provides a lower bound on the true KL divergence.  
27

28 Another use of the compression lemma is that it provides a way to estimate the expectation of  
29 some function with respect to an unknown distribution  $P$ . In this spirit, the Compression Lemma  
30 can be used to power a set of what are known as PAC Bayes bounds of losses with respect to the  
31 true distribution in terms of measured losses with respect to a finite training set. See for example  
32 Section 17.4.5 or Banerjee [Ban06].  
33

#### 34 5.1.4.2 Data processing inequality for KL 35

36 We now show that any processing we do on samples from two different distributions makes their  
37 samples approach one another. This is called the **data processing inequality**, since it shows that  
38 we cannot increase the information gain from  $q$  to  $p$  by processing our data and then measuring it.  
39

40 **Theorem 5.1.3.** *Consider two different distributions  $p(x)$  and  $q(x)$  combined with a probabilistic  
41 channel  $t(y|x)$ . If  $p(y)$  is the distribution that results from sending samples from  $p(x)$  through the  
42 channel  $t(y|x)$  and similarly for  $q(y)$  we have that:*  
43

44 
$$D_{\text{KL}}(p(x) \parallel q(x)) \geq D_{\text{KL}}(p(y) \parallel q(y)) \quad (5.49)$$
  
45

4647

1 *Proof.* The proof uses Jensen's inequality from Section 5.1.2.2 again. Call  $p(x, y) = p(x)t(y|x)$  and  
2  $q(x, y) = q(x)t(y|x)$ .

3

$$D_{\text{KL}}(p(x) \parallel q(x)) = \int dx p(x) \log \frac{p(x)}{q(x)} \quad (5.50)$$

4

$$= \int dx \int dy p(x)t(y|x) \log \frac{p(x)t(y|x)}{q(x)t(y|x)} \quad (5.51)$$

5

$$= \int dx \int dy p(x, y) \log \frac{p(x, y)}{q(x, y)} \quad (5.52)$$

6

$$= - \int dy p(y) \int dx p(x|y) \log \frac{q(x, y)}{p(x, y)} \quad (5.53)$$

7

$$\geq - \int dy p(y) \log \left( \int dx p(x|y) \frac{q(x, y)}{p(x, y)} \right) \quad (5.54)$$

8

$$= - \int dy p(y) \log \left( \frac{q(y)}{p(y)} \int dx q(x|y) \right) \quad (5.55)$$

9

$$= \int dy p(y) \log \frac{p(y)}{q(y)} = D_{\text{KL}}(p(y) \parallel q(y)) \quad (5.56)$$

□

10 One way to interpret this result is that any processing done to random samples makes it harder to  
11 tell two distributions apart.

12 As a special form of processing, we can simply marginalize out a subset of random variables.

13 **Corollary 5.1.1.** (*Monotonicity of KL divergence*)

14

$$D_{\text{KL}}(p(x, y) \parallel q(x, y)) \geq D_{\text{KL}}(p(x) \parallel q(x)) \quad (5.57)$$

15 *Proof.* The proof is essentially the same as the one above.

16

$$D_{\text{KL}}(p(x, y) \parallel q(x, y)) = \int dx \int dy p(x, y) \log \frac{p(x, y)}{q(x, y)} \quad (5.58)$$

17

$$= - \int dy p(y) \int dx p(x|y) \log \left( \frac{q(y)}{p(y)} \frac{q(x|y)}{p(x|y)} \right) \quad (5.59)$$

18

$$\geq - \int dy p(y) \log \left( \frac{q(y)}{p(y)} \int dx q(x|y) \right) \quad (5.60)$$

19

$$= \int dy p(y) \log \frac{p(y)}{q(y)} = D_{\text{KL}}(p(y) \parallel q(y)) \quad (5.61)$$

(5.62) □

20 One intuitive interpretation of this result is that if you only partially observe random variables, it  
21 is harder to distinguish between two candidate distributions than if you observed all of them.

1 2 **5.1.5 KL divergence and MLE**

3 Suppose we want to find the distribution  $q$  that is as close as possible to  $p$ , as measured by KL  
4 divergence:

5

$$\underline{6} \quad q^* = \arg \min_q D_{\text{KL}}(p \parallel q) = \arg \min_q \int p(x) \log p(x) dx - \int p(x) \log q(x) dx \quad (5.63)$$

7

8 Now suppose  $p$  is the empirical distribution, which puts a probability atom on the observed training  
9 data and zero mass everywhere else:

10

$$\underline{11} \quad p_{\mathcal{D}}(x) = \frac{1}{N_{\mathcal{D}}} \sum_{n=1}^{N_{\mathcal{D}}} \delta(x - x_n) \quad (5.64)$$

12

13 Using the sifting property of delta functions we get

14

$$\underline{15} \quad D_{\text{KL}}(p_{\mathcal{D}} \parallel q) = - \int p_{\mathcal{D}}(x) \log q(x) dx + C \quad (5.65)$$

16

17

$$\underline{18} \quad = - \int \left[ \frac{1}{N_{\mathcal{D}}} \sum_n \delta(x - x_n) \right] \log q(x) dx + C \quad (5.66)$$

19

20

$$\underline{21} \quad = - \frac{1}{N_{\mathcal{D}}} \sum_n \log q(x_n) + C \quad (5.67)$$

22

23 where  $C = \int p_{\mathcal{D}}(x) \log p_{\mathcal{D}}(x)$  is a constant independent of  $q$ .

24 We can rewrite the above as follows

25

$$\underline{26} \quad D_{\text{KL}}(p_{\mathcal{D}} \parallel q) = \mathbb{H}_{ce}(p_{\mathcal{D}}, q) - \mathbb{H}(p_{\mathcal{D}}) \quad (5.68)$$

27

28 where

29

$$\underline{30} \quad \mathbb{H}_{ce}(p, q) \triangleq - \sum_k p_k \log q_k \quad (5.69)$$

31

32 is known as the **cross entropy**. The quantity  $\mathbb{H}_{ce}(p_{\mathcal{D}}, q)$  is the average negative log likelihood  
33 of  $q$  evaluated on the training set. Thus we see that minimizing KL divergence to the empirical  
34 distribution is equivalent to maximizing likelihood.

35 This perspective points out the flaw with likelihood-based training, namely that it puts too  
36 much weight on the training set. In most applications, we do not really believe that the empirical  
37 distribution is a good representation of the true distribution, since it just puts “spikes” on a finite  
38 set of points, and zero density everywhere else. Even if the dataset is large (say 1M images), the  
39 universe from which the data is sampled is usually even larger (e.g., the set of “all natural images”  
40 is much larger than 1M). Thus we need to somehow smooth the empirical distribution by sharing  
41 probability mass between “similar” inputs.

42 43 **5.1.6 KL divergence and Bayesian Inference**

44 Bayesian inference itself can be motivated as the solution to a particular minimization problem of  
45 KL.

46

Consider a prior set of beliefs described by a joint distribution  $q(\theta, D) = q(\theta)q(D|\theta)$ , involving some *prior*  $q(\theta)$  and some *likelihood*  $q(D|\theta)$ . If we happen to observe some particular dataset  $D_0$ , how should we update our beliefs? We could search for the joint distribution that is as close as possible to our prior beliefs but that respects the constraint that we now know the value of the data:

$$p(\theta, D) = \operatorname{argmin} D_{\text{KL}}(p(\theta, D) \parallel q(\theta, D)) \text{ such that } p(D) = \delta(D - D_0). \quad (5.70)$$

where  $\delta(D - D_0)$  is a degenerate distribution that puts all its mass on the dataset  $D$  that is identically equal to  $D_0$ . Writing the KL out in its chain rule form:

$$D_{\text{KL}}(p(\theta, D) \parallel q(\theta, D)) = D_{\text{KL}}(p(D) \parallel q(D)) + D_{\text{KL}}(p(\theta|D) \parallel q(\theta|D)), \quad (5.71)$$

makes clear that the solution is given by the joint distribution:

$$p(\theta, D) = p(D)p(\theta|D) = \delta(D - D_0)q(\theta|D). \quad (5.72)$$

Our updated beliefs have a marginal over the  $\theta$

$$p(\theta) = \int dD p(\theta, D) = \int dD \delta(D - D_0)q(\theta|D) = q(\theta|D = D_0), \quad (5.73)$$

which is just the usual Bayesian posterior from our prior beliefs evaluated at the data we observed.

By contrast, the usual statement of Bayes' rule is just a trivial observation about the chain rule of probabilities:

$$q(\theta, D) = q(D)q(\theta|D) = q(\theta)q(D|\theta) \implies q(\theta|D) = \frac{q(D|\theta)}{q(D)}q(\theta). \quad (5.74)$$

Notice that this relates the conditional distribution  $q(\theta|D)$  in terms of  $q(D|\theta)$ ,  $q(\theta)$  and  $q(D)$ , but that these are all different ways to write the same distribution. Bayes rule does not tell us how we ought to *update* our beliefs in light of evidence, for that we need some other principle [Cat+11].

One of the nice things about this interpretation of Bayesian inference is that it naturally generalizes to other forms of constraints rather than assuming we have observed the data exactly.

If there was some additional measurement error that was well understood, we ought to instead of pegging out updated beliefs to be a delta function on the observed data, simply peg it to be the well understood distribution  $p(D)$ . For example, we might not know the precise value the data takes, but believe after measuring things that it is a Gaussian distribution with a certain mean and standard deviation.

Because of the chain rule of KL, this has no effect on our updated conditional distribution over parameters, which remains the Bayesian posterior:  $p(\theta|D) = q(\theta|D)$ . However, this does change our marginal beliefs about the parameters, which are now:

$$p(\theta) = \int dD p(D)q(\theta|D). \quad (5.75)$$

This generalization of Bayes' rule is sometimes called **Jeffrey's conditionalization rule** [Cat08].

1 **5.1.7 KL divergence and Exponential Families**

3 The KL divergence between two exponential family distributions from the same family has a nice  
4 closed form, as we explain below.

5 Consider  $p(\mathbf{x})$  with natural parameter  $\boldsymbol{\eta}$ , base measure  $h(\mathbf{x})$  and sufficient statistics  $\mathcal{T}(\mathbf{x})$ :

$$\underline{7} \quad p(\mathbf{x}) = h(\mathbf{x}) \exp[\boldsymbol{\eta}^\top \mathcal{T}(\mathbf{x}) - A(\boldsymbol{\eta})] \quad (5.76)$$

8 where

$$\underline{10} \quad A(\boldsymbol{\eta}) = \log \int h(\mathbf{x}) \exp(\boldsymbol{\eta}^\top \mathcal{T}(\mathbf{x})) d\mathbf{x} \quad (5.77)$$

12 is the *log partition function*, a convex function of  $\boldsymbol{\eta}$ .

14 The KL divergence between two exponential family distributions from the same family is as follows:

$$\underline{15} \quad D_{\text{KL}}(p(\mathbf{x}|\boldsymbol{\eta}_1) \parallel p(\mathbf{x}|\boldsymbol{\eta}_2)) = \mathbb{E}_{\boldsymbol{\eta}_1} [(\boldsymbol{\eta}_1 - \boldsymbol{\eta}_2)^\top \mathcal{T}(\mathbf{x}) - A(\boldsymbol{\eta}_1) + A(\boldsymbol{\eta}_2)] \quad (5.78)$$

$$\underline{17} \quad = (\boldsymbol{\eta}_1 - \boldsymbol{\eta}_2)^\top \boldsymbol{\mu}_1 - A(\boldsymbol{\eta}_1) + A(\boldsymbol{\eta}_2) \quad (5.79)$$

18 where  $\boldsymbol{\mu}_j \triangleq \mathbb{E}_{\boldsymbol{\eta}_j} [\mathcal{T}(\mathbf{x})]$ .

21 **5.1.7.1 Example: KL divergence between two Gaussians**

22 An important example is the KL divergence between two multivariate Gaussian distributions, which  
23 is given by

$$\begin{aligned} \underline{25} \quad & D_{\text{KL}}(\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_1, \Sigma_1) \parallel \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_2, \Sigma_2)) \\ \underline{26} \quad &= \frac{1}{2} \left[ \text{tr}(\Sigma_2^{-1} \Sigma_1) + (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)^\top \Sigma_2^{-1} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) - D + \log \left( \frac{\det(\Sigma_2)}{\det(\Sigma_1)} \right) \right] \end{aligned} \quad (5.80)$$

29 In the scalar case, this becomes

$$\begin{aligned} \underline{31} \quad & D_{\text{KL}}(\mathcal{N}(x|\mu_1, \sigma_1) \parallel \mathcal{N}(x|\mu_2, \sigma_2)) = \log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2} \\ \underline{32} \end{aligned} \quad (5.81)$$

33 **5.1.8 Bregman divergence**

35 Let  $f : \Omega \rightarrow \mathbb{R}$  be a continuously differentiable, strictly convex function defined on a closed convex  
36 set  $\Omega$ . We define the **Bregman divergence** associated with  $f$  as follows [Bre67]:

$$\underline{38} \quad B_f(\mathbf{w} \parallel \mathbf{v}) = f(\mathbf{w}) - f(\mathbf{v}) - (\mathbf{w} - \mathbf{v})^\top \nabla f(\mathbf{v}) \quad (5.82)$$

40 To understand this, let

$$\underline{41} \quad \hat{f}_\mathbf{v}(\mathbf{w}) = f(\mathbf{v}) + (\mathbf{w} - \mathbf{v})^\top \nabla f(\mathbf{v}) \quad (5.83)$$

43 be a first order Taylor series approximation to  $f$  centered at  $\mathbf{v}$ . Then the Bregman divergence is the  
44 difference from this linear approximation:

$$\underline{46} \quad B_f(\mathbf{w} \parallel \mathbf{v}) = f(\mathbf{w}) - \hat{f}_\mathbf{v}(\mathbf{w}) \quad (5.84)$$

47

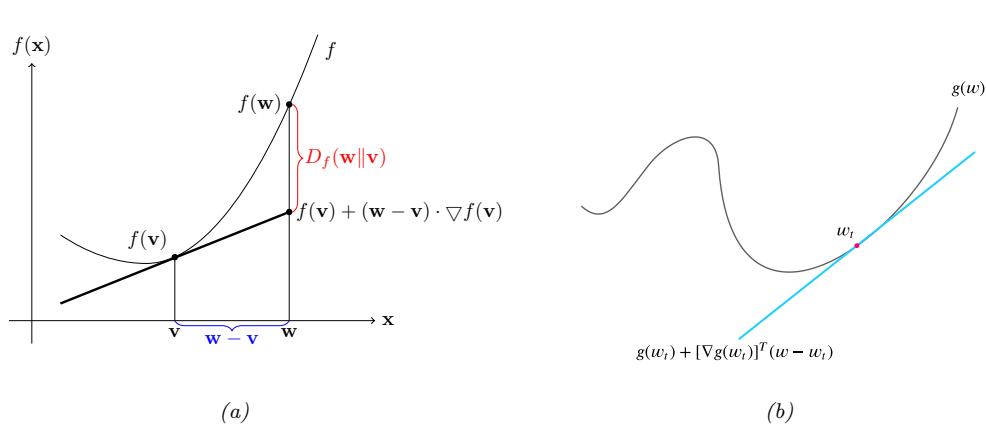


Figure 5.3: (a) Illustration of Bregman divergence. (b) A locally linear approximation to a non-convex function.

See Figure 5.3a for an illustration. Since  $f$  is convex, we have  $B_f(\mathbf{w}||\mathbf{v}) \geq 0$ , since  $\hat{f}_v$  is a linear lower bound on  $f$ .

Below we mention some important special cases of Bregman divergences.

- If  $f(\mathbf{w}) = \|\mathbf{w}\|^2$ , then  $B_f(\mathbf{w}||\mathbf{v}) = \|\mathbf{w} - \mathbf{v}\|^2$  is the squared Euclidean distance.
- If  $f(\mathbf{w}) = \mathbf{w}^\top \mathbf{Q} \mathbf{w}$ , then  $B_f(\mathbf{w}||\mathbf{v})$  is the squared Mahalanobis distance.
- If  $\mathbf{w}$  are the natural parameters of an exponential family distribution, and  $f(\mathbf{w}) = \log Z(\mathbf{w})$  is the log normalizer, then the Bregman divergence is the same as the Kullback Leibler divergence, as we show in Section 5.1.8.1.

### 5.1.8.1 KL is a Bregman divergence

Recall that the log partition function  $A(\boldsymbol{\eta})$  is a convex function. We can therefore use it to define the Bregman divergence (Section 5.1.8) between the two distributions,  $p$  and  $q$ , as follows:

$$B_f(\boldsymbol{\eta}_q||\boldsymbol{\eta}_p) = A(\boldsymbol{\eta}_q) - A(\boldsymbol{\eta}_p) - (\boldsymbol{\eta}_q - \boldsymbol{\eta}_p)^\top \nabla_{\boldsymbol{\eta}_p} A(\boldsymbol{\eta}_p) \quad (5.85)$$

$$= A(\boldsymbol{\eta}_q) - A(\boldsymbol{\eta}_p) - (\boldsymbol{\eta}_q - \boldsymbol{\eta}_p)^\top \mathbb{E}_p [\mathcal{T}(\mathbf{x})] \quad (5.86)$$

$$= D_{\text{KL}}(p \parallel q) \quad (5.87)$$

where we exploited the fact that the gradient of the log partition function computes the expected sufficient statistics as shown in Section 2.3.3.

In fact, the KL divergence is the only divergence that is both a Bregman divergence and an  $f$ -divergence (Section 2.7.1) [Ama09].

1 **5.2 Entropy**

3 In this section, we discuss the **entropy** of a distribution  $p$ , which is just a shifted and scaled version  
4 of the KL divergence between the probability distribution and the uniform distribution, as we will  
5 see.  
6

7 **5.2.1 Definition**

9 The entropy of a discrete random variable  $X$  with distribution  $p$  over  $K$  states is defined by  
10

$$\begin{aligned} \text{H}(X) &\triangleq -\sum_{k=1}^K p(X=k) \log p(X=k) = -\mathbb{E}_X [\log p(X)] \end{aligned} \quad (5.88)$$

14 We can use logarithms to any base, but we commonly use log base 2, in which case the units are  
15 called bits, or log base e, in which case the units are called nats, as we explained in Section 5.1.3.1.

16 The entropy is equivalent to a constant minus the KL divergence from the uniform distribution:  
17

$$\text{H}(X) = \log K - D_{\text{KL}}(p(X) \parallel u(X)) \quad (5.89)$$

$$\begin{aligned} D_{\text{KL}}(p(X) \parallel u(X)) &= \sum_{k=1}^K p(X=k) \log \frac{p(X=k)}{\frac{1}{K}} \end{aligned} \quad (5.90)$$

$$= \log K + \sum_{k=1}^K p(X=k) \log p(X=k) \quad (5.91)$$

26 If  $p$  is uniform, the KL is zero, and we see that the entropy achieves its maximal value of  $\log K$ .

27 For the special case of binary random variables,  $X \in \{0, 1\}$ , we can write  $p(X=1) = \theta$  and  
28  $p(X=0) = 1 - \theta$ . Hence the entropy becomes

$$\text{H}(X) = -[p(X=1) \log p(X=1) + p(X=0) \log p(X=0)] \quad (5.92)$$

$$= -[\theta \log \theta + (1 - \theta) \log(1 - \theta)] \quad (5.93)$$

32 This is called the **binary entropy function**, and is also written  $\text{H}(\theta)$ . We plot this in Figure 5.4.  
33 We see that the maximum value of 1 bit occurs when the distribution is uniform,  $\theta = 0.5$ . A fair coin  
34 requires a single yes/no question to determine its state.  
35

36 **5.2.2 Differential entropy for continuous random variables**

38 If  $X$  is a continuous random variable with pdf  $p(x)$ , we define the **differential entropy** as  
39

$$h(X) \triangleq - \int_{\mathcal{X}} dx p(x) \log p(x) \quad (5.94)$$

42 assuming this integral exists.

44 For example, one can show that the entropy of a  $d$ -dimensional Gaussian is

$$h(\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})) = \frac{1}{2} \log |2\pi e \boldsymbol{\Sigma}| = \frac{1}{2} \log[(2\pi e)^d |\boldsymbol{\Sigma}|] = \frac{d}{2} + \frac{d}{2} \log(2\pi) + \frac{1}{2} \log |\boldsymbol{\Sigma}| \quad (5.95)$$

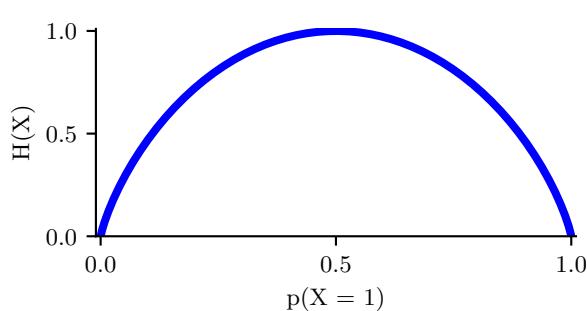


Figure 5.4: Entropy of a Bernoulli random variable as a function of  $\theta$ . The maximum entropy is  $\log_2 2 = 1$ . Generated by [bernpoulli\\_entropy\\_fig.ipynb](#).

In the 1d case, this becomes

$$h(\mathcal{N}(\mu, \sigma^2)) = \frac{1}{2} \log [2\pi e \sigma^2] \quad (5.96)$$

Note that, unlike the discrete case, *differential entropy can be negative*. This is because pdf's can be bigger than 1. For example, suppose  $X \sim U(0, a)$ . Then

$$h(X) = - \int_0^a dx \frac{1}{a} \log \frac{1}{a} = \log a \quad (5.97)$$

If we set  $a = 1/8$ , we have  $h(X) = \log_2(1/8) = -3$  bits.

One way to understand differential entropy is to realize that all real-valued quantities can only be represented to finite precision. It can be shown [CT91, p228] that the entropy of an  $n$ -bit quantization of a continuous random variable  $X$  is approximately  $h(X) + n$ . For example, suppose  $X \sim U(0, \frac{1}{8})$ . Then in a binary representation of  $X$ , the first 3 bits to the right of the binary point must be 0 (since the number is  $\leq 1/8$ ). So to describe  $X$  to  $n$  bits of accuracy only requires  $n - 3$  bits, which agrees with  $h(X) = -3$  calculated above.

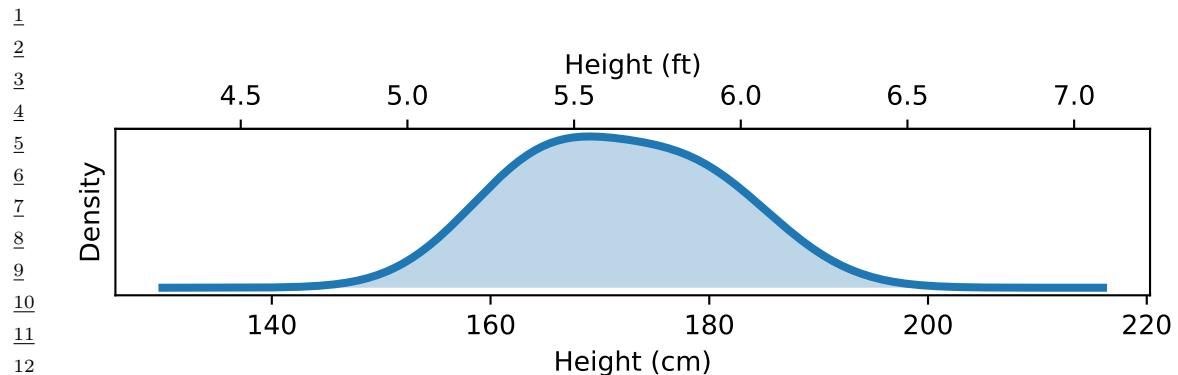
The continuous entropy also lacks the reparameterization independence of KL divergence Section 5.1.2.3. In particular, if we transform our random variable  $y = f(x)$ , the entropy transforms. To see this, note that the change of variables tells us that

$$p(y) dy = p(x) dx \implies p(y) = p(x) \left| \frac{dy}{dx} \right|^{-1}, \quad (5.98)$$

Thus the continuous entropy transforms as follows:

$$h(X) = - \int dx p(x) \log p(x) = h(Y) - \int dy p(y) \log \left| \frac{dy}{dx} \right|. \quad (5.99)$$

We pick up a factor in the continuous entropy of the log of the determinant of the Jacobian of the transformation. This changes the value for the continuous entropy even for simply rescaling the random variable such as when we change units. For example in Figure 5.5 we show the distribution



*Figure 5.5: Distribution of adult heights. The continuous entropy of the distribution depends on its units of measurement. If heights are measured in feet, this distribution has a continuous entropy of 0.43 bits. If measured in centimeters it's 5.4 bits. If measured in meters it's -1.3 bits. Data taken from <https://ourworldindata.org/human-height>.*

of adult human heights (it is bimodal because while both male and female heights are normally distributed, they differ noticeably). The continuous entropy of this distribution depends on the units it is measured in. If measured in feet, the continuous entropy is 0.43 bits. Intuitively this is because human heights mostly span less than a foot. If measured in centimeters it is instead 5.4 bits. There are 30.48 centimeters in a foot,  $\log_2 30.48 = 4.9$  explaining the difference. If we measured the continuous entropy of the same distribution measured in meters we would obtain -1.3 bits!

### 5.2.3 Typical sets

The **typical set** of a probability distribution is the set whose elements have an information content that is close to that of the expected information content from random samples from the distribution. More precisely, for a distribution  $p(\mathbf{x})$  with support  $\mathbf{x} \in \mathcal{X}$ , the  $\epsilon$ -typical set  $\mathcal{A}_\epsilon^N \in \mathcal{X}^N$  for  $p(\mathbf{x})$  is the set of all length  $N$  sequences such that

$$\mathbb{H}(p(\mathbf{x})) - \epsilon \leq -\frac{1}{N} \log p(\mathbf{x}_1, \dots, \mathbf{x}_N) \leq \mathbb{H}(p(\mathbf{x})) + \epsilon \quad (5.100)$$

If we assume  $p(\mathbf{x}_1, \dots, \mathbf{x}_N) = \prod_{n=1}^N p(\mathbf{x}_n)$ , then we can interpret the term in the middle as the  $N$ -sample empirical estimate of the entropy. The **asymptotic equipartition property** or **AEP** states that this will converge (in probability) to the true entropy as  $N \rightarrow \infty$  [CT06]. Thus the typical set has probability close to 1, and is thus a compact summary of what we can expect to be generated by  $p(\mathbf{x})$ .

1 **5.2.4 Cross entropy and perplexity**

2 A standard way to measure how close a model  $q$  is to a true distribution  $p$  is in terms of the KL  
3 divergence (Section 5.1), given by

4

$$D_{\text{KL}}(p \parallel q) = \sum_x p(x) \log \frac{p(x)}{q(x)} = \mathbb{H}_{ce}(p, q) - \mathbb{H}(p) \quad (5.101)$$

5

6 where  $\mathbb{H}_{ce}(p, q)$  is the **cross entropy**

7

8

$$\mathbb{H}_{ce}(p, q) = - \sum_x p(x) \log q(x) \quad (5.102)$$

9

10 and  $\mathbb{H}(p) = \mathbb{H}_{ce}(p, p)$  is the entropy, which is a constant independent of the model.

11 In language modeling, it is common to report an alternative performance measure known as the  
12 **perplexity**. This is defined as

13

$$\text{perplexity}(p, q) \triangleq 2^{\mathbb{H}_{ce}(p, q)} \quad (5.103)$$

14

15 We can compute an empirical approximation to the cross entropy as follows. Suppose we approxi-  
16 mate the true distribution with an empirical distribution based on data sampled from  $p$ :

17

$$p_{\mathcal{D}}(x|\mathcal{D}) = \frac{1}{N} \sum_{n=1}^N \mathbb{I}(x = x_n) \quad (5.104)$$

18

19 In this case, the cross entropy is given by

20

$$H = - \frac{1}{N} \sum_{n=1}^N \log p(x_n) = - \frac{1}{N} \log \prod_{n=1}^N p(x_n) \quad (5.105)$$

21

22 The corresponding perplexity is given by

23

$$\text{perplexity}(p_{\mathcal{D}}, p) = 2^{-\frac{1}{N} \log(\prod_{n=1}^N p(x_n))} = 2^{\log(\prod_{n=1}^N p(x_n))^{-\frac{1}{N}}} \quad (5.106)$$

24

25

$$= \left( \prod_{n=1}^N p(x_n) \right)^{-1/N} = \sqrt[N]{\prod_{n=1}^N \frac{1}{p(x_n)}} \quad (5.107)$$

26

27 In the case of language models, we usually condition on previous words when predicting the next  
28 word. For example, in a bigram model, we use a second order Markov model of the form  $p(x_n|x_{n-1})$ .  
29 We define the **branching factor** of a language model as the number of possible words that can  
30 follow any given word. For example, suppose the model predicts that each word is equally likely,  
31 regardless of context, so  $p(x_n|x_{n-1}) = 1/K$ , where  $K$  is the number of words in the vocabulary. Then  
32 the perplexity is  $((1/K)^N)^{-1/N} = K$ . If some symbols are more likely than others, and the model  
33 correctly reflects this, its perplexity will be lower than  $K$ . However, we have  $\mathbb{H}(p^*) \leq \mathbb{H}_{ce}(p^*, p)$ , so  
34 we can never reduce the perplexity below  $2^{-\mathbb{H}(p^*)}$ .

1 **5.3 Mutual information**

2 The KL divergence gave us a way to measure how similar two distributions were. How should we  
3 measure how dependant two random variables are? One thing we could do is turn the question  
4 of measuring the dependence of two random variables into a question about the similarity of their  
5 distributions. This gives rise to the notion of **mutual information** (MI) between two random  
6 variables, which we define below.

7 **5.3.1 Definition**

8 The mutual information between rv's  $X$  and  $Y$  is defined as follows:

$$\underline{13} \quad \mathbb{I}(X; Y) \triangleq D_{\text{KL}}(p(x, y) \parallel p(x)p(y)) = \sum_{y \in Y} \sum_{x \in X} p(x, y) \log \frac{p(x, y)}{p(x)p(y)} \quad (5.108)$$

16 (We write  $\mathbb{I}(X; Y)$  instead of  $\mathbb{I}(X, Y)$ , in case  $X$  and/or  $Y$  represent sets of variables; for example, we  
17 can write  $\mathbb{I}(X; Y, Z)$  to represent the MI between  $X$  and  $(Y, Z)$ .) For continuous random variables,  
18 we just replace sums with integrals.

19 It is easy to see that MI is always non-negative, even for continuous random variables, since

$$\underline{21} \quad \mathbb{I}(X; Y) = D_{\text{KL}}(p(x, y) \parallel p(x)p(y)) \geq 0 \quad (5.109)$$

22 We achieve the bound of 0 iff  $p(x, y) = p(x)p(y)$ .

24 **5.3.2 Interpretation**

26 Knowing that the mutual information is a KL divergence between the joint and factored marginal  
27 distributions tells us that the MI measures the information gain if we update from a model that treats  
28 the two variables as independent  $p(x)p(y)$  to one that models their true joint density  $p(x, y)$ .

29 To gain further insight into the meaning of MI, it helps to re-express it in terms of joint and  
30 conditional entropies, as follows:

$$\underline{32} \quad \mathbb{I}(X; Y) = \mathbb{H}(X) - \mathbb{H}(X|Y) = \mathbb{H}(Y) - \mathbb{H}(Y|X) \quad (5.110)$$

34 Thus we can interpret the MI between  $X$  and  $Y$  as the reduction in uncertainty about  $X$  after  
35 observing  $Y$ , or, by symmetry, the reduction in uncertainty about  $Y$  after observing  $X$ . Incidentally,  
36 this result gives an alternative proof that conditioning, on average, reduces entropy. In particular, we  
37 have  $0 \leq \mathbb{I}(X; Y) = \mathbb{H}(X) - \mathbb{H}(X|Y)$ , and hence  $\mathbb{H}(X|Y) \leq \mathbb{H}(X)$ .

38 We can also obtain a different interpretation. One can show that

$$\underline{40} \quad \mathbb{I}(X; Y) = \mathbb{H}(X, Y) - \mathbb{H}(X|Y) - \mathbb{H}(Y|X) \quad (5.111)$$

41 Finally, one can show that

$$\underline{43} \quad \mathbb{I}(X; Y) = \mathbb{H}(X) + \mathbb{H}(Y) - \mathbb{H}(X, Y) \quad (5.112)$$

45 See Figure 5.6 for a summary of these equations in terms of an **information diagram**. (Formally,  
46 this is a signed measure mapping set expressions to their information-theoretic counterparts [Yeu91a].)

47

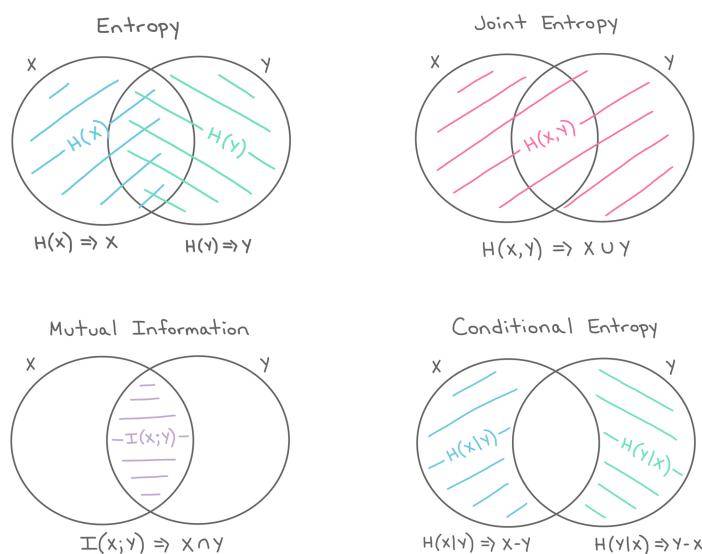


Figure 5.6: The marginal entropy, joint entropy, conditional entropy and mutual information represented as information diagrams. Used with kind permission of Katie Everett.

### 5.3.3 Data processing inequality

Suppose we have an unknown variable  $X$ , and we observe a noisy function of it, call it  $Y$ . If we process the noisy observations in some way to create a new variable  $Z$ , it should be intuitively obvious that we cannot increase the amount of information we have about the unknown quantity,  $X$ . This is known as the **data processing inequality**. We now state this more formally, and then prove it.

**Theorem 5.3.1.** Suppose  $X \rightarrow Y \rightarrow Z$  forms a Markov chain, so that  $X \perp Z|Y$ . Then  $\mathbb{I}(X;Y) \geq \mathbb{I}(X;Z)$ .

*Proof.* By the chain rule for mutual information we can expand the mutual information in two different ways:

$$\mathbb{I}(X;Y,Z) = \mathbb{I}(X;Z) + \mathbb{I}(X;Y|Z) \quad (5.113)$$

$$= \mathbb{I}(X;Y) + \mathbb{I}(X;Z|Y) \quad (5.114)$$

Since  $X \perp Z|Y$ , we have  $\mathbb{I}(X;Z|Y) = 0$ , so

$$\mathbb{I}(X;Z) + \mathbb{I}(X;Y|Z) = \mathbb{I}(X;Y) \quad (5.115)$$

Since  $\mathbb{I}(X;Y|Z) \geq 0$ , we have  $\mathbb{I}(X;Y) \geq \mathbb{I}(X;Z)$ . Similarly one can prove that  $\mathbb{I}(Y;Z) \geq \mathbb{I}(X;Z)$ .  $\square$

1    **5.3.4 Sufficient Statistics**

3 An important consequence of the DPI is the following. Suppose we have the chain  $\theta \rightarrow X \rightarrow s(X)$ .  
4 Then

5     $\mathbb{I}(\theta; s(X)) \leq \mathbb{I}(\theta; X) \tag{5.116}$

7 If this holds with equality, then we say that  $s(X)$  is a **sufficient statistic** of the data  $X$  for the  
8 purposes of inferring  $\theta$ . In this case, we can equivalently write  $\theta \rightarrow s(X) \rightarrow X$ , since we can  
9 reconstruct the data from knowing  $s(X)$  just as accurately as from knowing  $\theta$ .

10 An example of a sufficient statistic is the data itself,  $s(X) = X$ , but this is not very useful, since it  
11 doesn't summarize the data at all. Hence we define a **minimal sufficient statistic**  $s'(X)$  as one  
12 which is sufficient, and which contains no extra information about  $\theta$ ; thus  $s'(X)$  maximally compresses  
13 the data  $X$  without losing information which is relevant to predicting  $\theta$ . More formally, we say  $s$  is a  
14 minimal sufficient statistic for  $X$  if  $s(X) = f(s'(X))$  for some function  $f$  and all sufficient statistics  
15  $s'(X)$ . We can summarize the situation as follows:

16     $\theta \rightarrow s(X) \rightarrow s'(X) \rightarrow X \tag{5.117}$

18 Here  $s'(X)$  takes  $s(X)$  and adds redundant information to it, thus creating a one-to-many mapping.

19 For example, a minimal sufficient statistic for a set of  $N$  Bernoulli trials is simply  $N$  and  $N_1 =$   
20  $\sum_n \mathbb{I}(X_n = 1)$ , i.e., the number of successes. In other words, we don't need to keep track of the  
21 entire sequence of heads and tails and their ordering, we only need to keep track of the total number  
22 of heads and tails. Similarly, for inferring the mean of a Gaussian distribution with known variance  
23 we only need to know the empirical mean and number of samples.

24 Earlier in Section 5.1.7 we motivated the exponential family of distributions as being the ones that  
25 are minimal in the sense that they contain no other information than constraints on some statistics of  
26 the data. It makes sense then that the statistics used to generate exponential family distributions are  
27 sufficient. It also hints at the more remarkable fact of the **Pitman-Koopman-Darmois theorem**,  
28 which says that for any distribution whose domain is fixed, it is only the exponential family that  
29 admits sufficient statistics with bounded dimensionality as the number of samples increases Diaconis  
30 [Dia88].

32    **5.3.5 Multivariate mutual information**

34 There are several ways to generalize the idea of mutual information to a set of random variables as  
35 we discuss below.

36

37    **5.3.5.1 Total correlation**

38

39 The simplest way to define multivariate MI is to use the **total correlation** [Wat60] or **multi-**  
40 **information** [SV98], defined as

41     $\mathbb{T}\mathbb{C}(\{X_1, \dots, X_D\}) \triangleq D_{\mathbb{KL}} \left( p(\mathbf{x}) \parallel \prod_d p(x_d) \right) \tag{5.118}$

44     $= \sum_{\mathbf{x}} p(\mathbf{x}) \log \frac{p(\mathbf{x})}{\prod_{d=1}^D p(x_d)} = \sum_d \mathbb{H}(x_d) - \mathbb{H}(\mathbf{x}) \tag{5.119}$

45

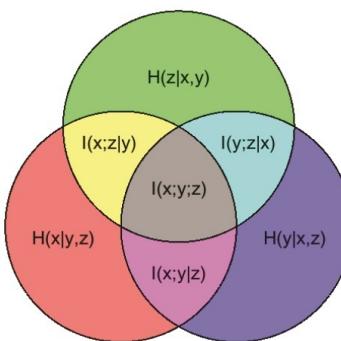


Figure 5.7: Illustration of multivariate mutual information between three random variables. From [https://en.wikipedia.org/wiki/Mutual\\_information](https://en.wikipedia.org/wiki/Mutual_information). Used with kind permission of Wikipedia author PAR.

For example, for 3 variables, this becomes

$$\mathbb{TC}(X, Y, Z) = \mathbb{H}(X) + \mathbb{H}(Y) + \mathbb{H}(Z) - \mathbb{H}(X, Y, Z) \quad (5.120)$$

where  $\mathbb{H}(X, Y, Z)$  is the joint entropy

$$\mathbb{H}(X, Y, Z) = - \sum_x \sum_y \sum_z p(x, y, z) \log p(x, y, z) \quad (5.121)$$

One can show that the multi-information is always non-negative, and is zero iff  $p(\mathbf{x}) = \prod_d p(x_d)$ .

However, this means the quantity is non-zero even if only a pair of variables interact. For example, if  $p(X, Y, Z) = p(X, Y)p(Z)$ , then the total correlation will be non-zero, even though there is no 3 way interaction. This motivates the alternative definition in Section 5.3.5.2.

### 5.3.5.2 Interaction information (co-information)

The conditional mutual information can be used to give an inductive definition of the **multivariate mutual information (MMI)** as follows:

$$\mathbb{I}(X_1; \dots; X_D) = \mathbb{I}(X_1; \dots; X_{D-1}) - \mathbb{I}(X_1; \dots; X_{D-1}|X_D) \quad (5.122)$$

This is called the **multiple mutual information** [Yeu91b], or the **co-information** [Bel03]. This definition is equivalent, up to a sign change, to the **interaction information** [McG54; Han80; JB03; Bro09].

For 3 variables, the MMI is given by

$$\mathbb{I}(X; Y; Z) = \mathbb{I}(X; Y) - \mathbb{I}(X; Y|Z) \quad (5.123)$$

$$= \mathbb{I}(X; Z) - \mathbb{I}(X; Z|Y) \quad (5.124)$$

$$= \mathbb{I}(Y; Z) - \mathbb{I}(Y; Z|X) \quad (5.125)$$

1 This can be interpreted as the change in mutual information between two pairs of variables when  
2 conditioning on the third. Note that this quantity is symmetric in its arguments.

3 By the definition of conditional mutual information, we have

4

$$\mathbb{I}(X; Z|Y) = \mathbb{I}(Z; X, Y) - \mathbb{I}(Y; Z) \quad (5.126)$$

5

6 Hence we can rewrite Equation (5.124) as follows:

7

$$\mathbb{I}(X; Y; Z) = \mathbb{I}(X; Z) + \mathbb{I}(Y; Z) - \mathbb{I}(X, Y; Z) \quad (5.127)$$

8 This tells us that the MMI is the difference between how much we learn about  $Z$  given  $X$  and  $Y$   
9 individually vs jointly (see also Section 5.3.5.3).

10 The 3-way MMI is illustrated in the information diagram in Figure 5.7. The way to interpret such  
11 diagrams when we have multiple variables is as follows: the area of a shaded area that includes circles  
12  $A, B, C, \dots$  and excludes circles  $F, G, H, \dots$  represents  $\mathbb{I}(A; B; C; \dots | F, G, H, \dots)$ ; if  $B = C = \emptyset$ , this  
13 is just  $\mathbb{H}(A|F, G, H, \dots)$ ; if  $F = G = H = \emptyset$ , this is just  $\mathbb{I}(A; B; C, \dots)$ .

14

15

### 16 5.3.5.3 Synergy and redundancy

17 The MMI is  $\mathbb{I}(X; Y; Z) = \mathbb{I}(X; Z) + \mathbb{I}(Y; Z) - \mathbb{I}(X, Y; Z)$ . We see that this can be positive, zero or  
18 negative. If some of the information about  $Z$  that is provided by  $X$  is also provided by  $Y$ , then  
19 there is some **redundancy** between  $X$  and  $Y$  (wrt  $Z$ ). In this case,  $\mathbb{I}(X; Z) + \mathbb{I}(Y; Z) > \mathbb{I}(X, Y; Z)$ ,  
20 so (from Equation (5.127)) we see that the MMI will be positive. If, by contrast, we learn more  
21 about  $Z$  when we see  $X$  and  $Y$  together, we say there is some **synergy** between them. In this case,  
22  $\mathbb{I}(X; Z) + \mathbb{I}(Y; Z) < \mathbb{I}(X, Y; Z)$ , so the MMI will be negative.

23

### 24 5.3.5.4 MMI and causality

25

26 The sign of the MMI can be used to distinguish between different kinds of directed graphical models,  
27 which can sometimes be interpreted causally (see Chapter 36 for a general discussion of causality).  
28 For example, consider a model of the form  $X \leftarrow Z \rightarrow Y$ , where  $Z$  is a “cause” of  $X$  and  $Y$ . For  
29 example, suppose  $X$  represents the event it is raining,  $Y$  represents the event that the sky is dark,  
30 and  $Z$  represents the event that the sky is cloudy. Conditioning on the common cause  $Z$  renders  
31 the children  $X$  and  $Y$  independent, since if I know it is cloudy, noticing that the sky is dark does  
32 not change my beliefs about whether it will rain or not. Consequently  $\mathbb{I}(X; Y|Z) \leq \mathbb{I}(X; Y)$ , so  
33  $\mathbb{I}(X; Y; Z) \geq 0$ .

34 Now consider the case where  $Z$  is a common effect,  $X \rightarrow Z \leftarrow Y$ . In this case, conditioning on  
35  $Z$  makes  $X$  and  $Y$  dependent, due to the explaining away phenomenon (see Section 4.2.4.2). For  
36 example, if  $X$  and  $Y$  are independent random bits, and  $Z$  is the XOR of  $X$  and  $Y$ , then observing  
37  $Z = 1$  means that  $p(X \neq Y|Z = 1) = 1$ , so  $X$  and  $Y$  are now dependent (information-theoretically,  
38 not causally), even though they were a priori independent. Consequently  $\mathbb{I}(X; Y|Z) \geq \mathbb{I}(X; Y)$ , so  
39  $\mathbb{I}(X; Y; Z) \leq 0$ .

40 Finally, consider a Markov chain,  $X \rightarrow Y \rightarrow Z$ . We have  $\mathbb{I}(X; Z|Y) \leq \mathbb{I}(X; Z)$  and so the MMI  
41 must be positive.

42

1 **5.3.5.5 MMI and entropy**

2 We can also write the MMI in terms of entropies. Specifically, we know that

3  $\mathbb{I}(X; Y) = \mathbb{H}(X) + \mathbb{H}(Y) - \mathbb{H}(X, Y)$  (5.128)

4 and

5  $\mathbb{I}(X; Y|Z) = \mathbb{H}(X, Z) + \mathbb{H}(Y, Z) - \mathbb{H}(Z) - \mathbb{H}(X, Y, Z)$  (5.129)

6 Hence we can rewrite Equation (5.123) as follows:

7  $\mathbb{I}(X; Y; Z) = [\mathbb{H}(X) + \mathbb{H}(Y) + \mathbb{H}(Z)] - [\mathbb{H}(X, Y) + \mathbb{H}(X, Z) + \mathbb{H}(Y, Z)] + \mathbb{H}(X, Y, Z)$  (5.130)

8 Contrast this to Equation (5.120).

9 More generally, we have

10  $\mathbb{I}(X_1, \dots, X_D) = - \sum_{\mathcal{T} \subseteq \{1, \dots, D\}} (-1)^{|\mathcal{T}|} \mathbb{H}(\mathcal{T})$  (5.131)

11 For sets of size 1, 2 and 3 this expands as follows:

12  $I_1 = H_1$  (5.132)

13  $I_{12} = H_1 + H_2 - H_{12}$  (5.133)

14  $I_{123} = H_1 + H_2 + H_3 - H_{12} - H_{13} - H_{23} + H_{123}$  (5.134)

15 We can use the **Mobius inversion formula** to derive the following dual relationship:

16  $\mathbb{H}(\mathcal{S}) = - \sum_{\mathcal{T} \subseteq \mathcal{S}} (-1)^{|\mathcal{T}|} \mathbb{I}(\mathcal{T})$  (5.135)

17 for sets of variables  $\mathcal{S}$ .

18 Using the chain rule for entropy, we can also derive the following expression for the 3-way MMI:

19  $\mathbb{I}(X; Y; Z) = \mathbb{H}(Z) - \mathbb{H}(Z|X) - \mathbb{H}(Z|Y) + \mathbb{H}(Z|X, Y)$  (5.136)

20 **5.3.6 Variational bounds on mutual information**

21 In this section, we discuss methods for computing upper and lower bounds on MI that use variational approximations to the intractable distributions. This can be useful for representation learning (Chapter 32). This approach was first suggested in [BA03]. For a more detailed overview of variational bounds on mutual information, see Poole et al. [Poo+19b].

22 **5.3.6.1 Upper bound**

23 Suppose that the joint  $p(\mathbf{x}, \mathbf{y})$  is intractable to evaluate, but that we can sample from  $p(\mathbf{x})$  and evaluate the conditional distribution  $p(\mathbf{y}|\mathbf{x})$ . Furthermore, suppose we approximate  $p(\mathbf{y})$  by  $q(\mathbf{y})$ .

1 Then we can compute an upper bound on the MI as follows:  
2

$$\begin{aligned} \underline{4} \quad \mathbb{I}(\mathbf{x}; \mathbf{y}) &= \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} \left[ \log \frac{p(\mathbf{y}|\mathbf{x})q(\mathbf{y})}{p(\mathbf{y})q(\mathbf{y})} \right] \\ \underline{5} \end{aligned} \tag{5.137}$$

$$\begin{aligned} \underline{6} \quad &= \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} \left[ \log \frac{p(\mathbf{y}|\mathbf{x})}{q(\mathbf{y})} \right] - D_{\text{KL}}(p(\mathbf{y}) \parallel q(\mathbf{y})) \\ \underline{7} \end{aligned} \tag{5.138}$$

$$\begin{aligned} \underline{8} \quad &\leq \mathbb{E}_{p(\mathbf{x})} \left[ \mathbb{E}_{p(\mathbf{y}|\mathbf{x})} \left[ \log \frac{p(\mathbf{y}|\mathbf{x})}{q(\mathbf{y})} \right] \right] \\ \underline{9} \end{aligned} \tag{5.139}$$

$$\begin{aligned} \underline{10} \quad &= \mathbb{E}_{p(\mathbf{x})} [D_{\text{KL}}(p(\mathbf{y}|\mathbf{x}) \parallel q(\mathbf{y}))] \\ \underline{11} \end{aligned} \tag{5.140}$$

12

13 This bound is tight if  $q(\mathbf{y}) = p(\mathbf{y})$ .

14 What's happening here is that  $\mathbb{I}(Y; X) = \mathbb{H}(Y) - \mathbb{H}(Y|X)$  and we've assumed we know  $p(\mathbf{y}|\mathbf{x})$   
15 and so can estimate  $\mathbb{H}(Y|X)$  well. While we don't know  $\mathbb{H}(Y)$ , we can upper bound it using some  
16 model  $q(\mathbf{y})$ . Our model can never do better than  $p(\mathbf{y})$  itself (the non-negativity of KL), so our  
17 entropy estimate errs too large, and hence our MI estimate will be an upper bound.

181920

### 21 5.3.6.2 BA lower bound

22 Suppose that the joint  $p(\mathbf{x}, \mathbf{y})$  is intractable to evaluate, but that we can evaluate  $p(\mathbf{x})$ . Furthermore,  
23 suppose we approximate  $p(\mathbf{x}|\mathbf{y})$  by  $q(\mathbf{x}|\mathbf{y})$ . Then we can derive the following variational lower bound  
24 on the mutual information:

25

$$\begin{aligned} \underline{26} \quad \mathbb{I}(\mathbf{x}; \mathbf{y}) &= \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} \left[ \log \frac{p(\mathbf{x}|\mathbf{y})}{p(\mathbf{x})} \right] \\ \underline{27} \end{aligned} \tag{5.141}$$

$$\begin{aligned} \underline{28} \quad &= \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} \left[ \log \frac{q(\mathbf{x}|\mathbf{y})}{p(\mathbf{x})} \right] + \mathbb{E}_{p(\mathbf{y})} [D_{\text{KL}}(p(\mathbf{x}|\mathbf{y}) \parallel q(\mathbf{x}|\mathbf{y}))] \\ \underline{29} \end{aligned} \tag{5.142}$$

$$\begin{aligned} \underline{30} \quad &\geq \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} \left[ \log \frac{q(\mathbf{x}|\mathbf{y})}{p(\mathbf{x})} \right] = \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [\log q(\mathbf{x}|\mathbf{y})] + h(\mathbf{x}) \\ \underline{31} \end{aligned} \tag{5.143}$$

32

33 where  $h(\mathbf{x})$  is the differential entropy of  $\mathbf{x}$ . This is called the **BA lower bound**, after the authors  
34 Barber and Agakov [BA03].

3536

### 37 5.3.6.3 NWJ lower bound

38

39 The BA lower bound requires a tractable normalized distribution  $q(\mathbf{x}|\mathbf{y})$  that we can evaluate  
40 pointwise. If we reparameterize this distribution in a clever way, we can generate a lower bound that  
41 does not require a normalized distribution. Let's write:

42

$$\begin{aligned} \underline{43} \quad q(\mathbf{x}|\mathbf{y}) &= \frac{p(\mathbf{x})e^{f(\mathbf{x}, \mathbf{y})}}{Z(\mathbf{y})} \\ \underline{44} \end{aligned} \tag{5.144}$$

45

1 with  $Z(\mathbf{y}) = \mathbb{E}_{p(\mathbf{x})} [e^{f(\mathbf{x}, \mathbf{y})}]$  the normalization constant or partition function. Plugging this into the  
2 BA lower bound above we obtain:

3

$$\mathbb{E}_{p(\mathbf{x}, \mathbf{y})} \left[ \log \frac{p(\mathbf{x}) e^{f(\mathbf{x}, \mathbf{y})}}{p(\mathbf{x}) Z(\mathbf{y})} \right] = \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [f(\mathbf{x}, \mathbf{y})] - \mathbb{E}_{p(\mathbf{y})} [Z(\mathbf{y})] \quad (5.145)$$

4

$$= \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [f(\mathbf{x}, \mathbf{y})] - \mathbb{E}_{p(\mathbf{y})} \left[ \log \mathbb{E}_{p(\mathbf{x})} \left[ e^{f(\mathbf{x}, \mathbf{y})} \right] \right] \quad (5.146)$$

5

$$\triangleq I_{DV}(X; Y). \quad (5.147)$$

6 This is the **Donsker Varadhan lower bound** [DV75].

7 We can construct a more tractable version of this by using the fact that the log function can be  
8 upper bounded by a straight line using

9

$$\log x \leq \frac{x}{a} + \log a - 1 \quad (5.148)$$

10 If we set  $a = e$ , we get

11

$$\mathbb{I}(X; Y) \geq \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [f(\mathbf{x}, \mathbf{y})] - e^{-1} \mathbb{E}_{p(\mathbf{y})} Z(\mathbf{y}) \triangleq I_{NWJ}(X; Y) \quad (5.149)$$

12 This is called the **NWJ lower bound** (after the authors of Nguyen, Wainwright, and Jordan  
13 [NWJ10a]), or the f-GAN KL [NCT16a], or the MINE-f score [Bel+18].

#### 24 5.3.6.4 InfoNCE lower bound

25 If we instead explore a multi-sample extension to the DV bound above, we can generate the following  
26 lower bound (see [Poo+19b] for the derivation):

27

$$\mathbb{I}_{\text{NCE}} = \mathbb{E} \left[ \frac{1}{K} \sum_{i=1}^K \log \frac{e^{f(\mathbf{x}_i, \mathbf{y}_i)}}{\frac{1}{K} \sum_{j=1}^K e^{f(\mathbf{x}_i, \mathbf{y}_j)}} \right] \quad (5.150)$$

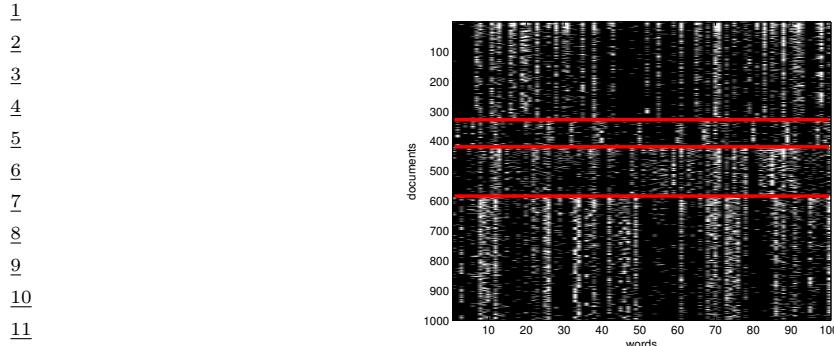
28

$$= \log K - \mathbb{E} \left[ \frac{1}{K} \sum_{i=1}^K \log \left( 1 + \sum_{j \neq i}^K e^{f(\mathbf{x}_i, \mathbf{y}_j) - f(\mathbf{x}_i, \mathbf{y}_i)} \right) \right] \quad (5.151)$$

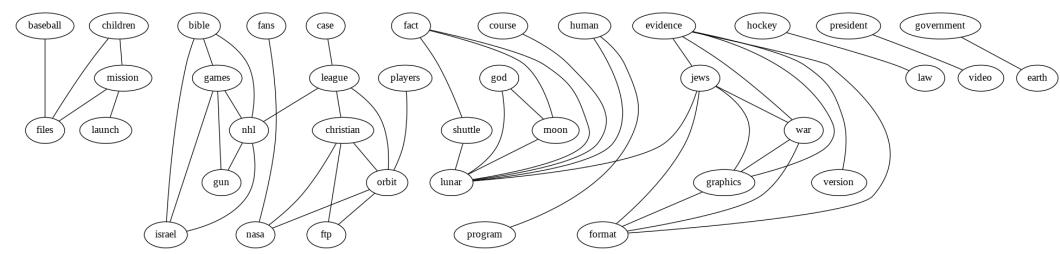
29 where the expectation is over paired samples from the joint  $p(X, Y)$ . The quantity in Equation (5.151)  
30 is called the **InfoNCE** estimate, and was proposed in [OLV18a; Hen+19a]. (NCE stands for “noise  
31 contrastive estimation”, and is discussed in Section 24.4.)

32 The intuition here is that mutual information is a divergence between the joint  $p(\mathbf{x}, \mathbf{y})$  and the  
33 product of the marginals,  $p(\mathbf{x})p(\mathbf{y})$ . In other words, mutual information is a measurement of how  
34 distinct sampling pairs jointly is from sampling  $\mathbf{x}$ s and  $\mathbf{y}$ s independently. The InfoNCE bound  
35 provides a lower bound on the true mutual information by attempting to train a model to distinguish  
36 between these two situations.

37 Although this is a valid lower bound, we may need to use a large batch size  $K$  to estimate the  
38 MI if the MI is large, since  $\mathbb{I}_{\text{NCE}} \leq \log K$ . (Recently [SE20a] proposed to use a multi-label classifier,  
39 rather than a multi-class classifier, to overcome this limitation.)



13 Figure 5.8: Subset of size  $16242 \times 100$  of the 20-newsgroups data. We only show 1000 rows, for clarity. Each  
14 row is a document (represented as a bag-of-words bit vector), each column is a word. The red lines separate  
15 the 4 classes, which are (in descending order) comp, rec, sci, talk (these are the titles of USENET groups).  
16 We can see that there are subsets of words whose presence or absence is indicative of the class. The data is  
17 available from <http://cs.nyu.edu/~rweis/data.html>. Generated by newsgroups\_visualize.ipynb.



27 Figure 5.9: Part of a relevance network constructed from the 20 newsgroup data. data shown in Figure 5.8.  
28 We show edges whose mutual information is greater than or equal to 20% of the maximum pairwise MI.  
29 For clarity, the graph has been cropped, so we only show a subset of the nodes and edges. Generated by  
30 relevance\_network\_newsgroup\_demo.ipynb.

### 33 5.3.7 Relevance networks

34 If we have a set of related variables, we can compute a **relevance network**, in which we add an  $i - j$   
35 edge if the pairwise mutual information  $\mathbb{I}(X_i; X_j)$  is above some threshold. In the Gaussian case,  
36  $\mathbb{I}(X_i; X_j) = -\frac{1}{2} \log(1 - \rho_{ij}^2)$ , where  $\rho_{ij}$  is the correlation coefficient, and the resulting graph is called  
37 a **covariance graph** (Section 4.5.5.1). However, we can also apply it to discrete random variables.  
38 Relevance networks are quite popular in systems biology [Mar+06], where they are used to visualize  
39 the interaction between genes. But they can also be applied to other kinds of datasets. For example,  
40 Figure 5.9 visualizes the MI between words in the 20 newsgroup dataset shown in Figure 5.8. The  
41 results seem intuitively reasonable.  
42

43 However, relevance networks suffer from a major problem: the graphs are usually very dense, since  
44 most variables are dependent on most other variables, even after thresholding the MIs. For example,  
45 suppose  $X_1$  directly influences  $X_2$  which directly influences  $X_3$  (e.g., these form components of a  
46 signalling cascade,  $X_1 - X_2 - X_3$ ). Then  $X_1$  has non-zero MI with  $X_3$  (and vice versa), so there  
47

1 will be a  $1 - 3$  edge as well as the  $1 - 2$  and  $2 - 3$  edges; thus the graph may be fully connected,  
2 depending on the threshold.

3 A solution to this is to learn a probabilistic graphical model, which represents conditional *in-*  
4 *dependence*, rather than *dependence*. In the chain example, there will not be a  $1 - 3$  edge, since  
5  $X_1 \perp X_3 | X_2$ . Consequently graphical models are usually much sparser than relevance networks. See  
6 Chapter 30 for details.  
7

## 8 5.4 Data compression (source coding)

9 **Data compression**, also known as **source coding**, is at the heart of information theory. It is also  
10 related to probabilistic machine learning. The reason for this is as follows: if we can model the  
11 probability of different kinds of data samples, then we can assign short **code words** to the most  
12 frequently occurring ones, reserving longer encodings for the less frequent ones. This is similar to  
13 the situation in natural language, where common words (such as “a”, “the”, “and”) are generally  
14 much shorter than rare words. Thus the ability to compress data requires an ability to discover  
15 the underlying patterns, and their relative frequencies, in the data. This has led Marcus Hutter  
16 to propose that compression be used as an objective way to measure performance towards general  
17 purpose AI. More precisely, he is offering 50,000 Euros to anyone who can compress the first 100MB  
18 of (English) Wikipedia better than some baseline. This is known as the **Hutter prize**.<sup>1</sup>  
19

20 In this section, we give a brief summary of some of the key ideas in data compression. For details,  
21 see e.g., [Mac03; CT06; YMT22].  
22

### 23 5.4.1 Lossless compression

24 Discrete data, such as natural language, can always be compressed in such a way that we can uniquely  
25 recover the original data. This is called **lossless compression**.

26 Claude Shannon proved that the expected number of bits needed to losslessly encode some data  
27 coming from distribution  $p$  is at least  $\mathbb{H}(p)$ . This is known as the **source coding theorem**. Achieving  
28 this lower bound requires coming up with good probability models, as well as good ways to design  
29 codes based on those models. Because of the non-negativity of the KL divergence,  $\mathbb{H}_{ce}(p, q) \geq \mathbb{H}(p)$ ,  
30 so if we use any model  $q$  other than the true model  $p$  to compress the data, it will take some excess  
31 bits. The number of excess bits is exactly  $D_{\text{KL}}(p \parallel q)$ .

32 Common techniques for realizing lossless codes include Huffman coding, arithmetic coding and  
33 asymmetric numeral systems [Dud13]. The input to these algorithms is a probability distribution  
34 over strings (which is where ML comes in). This distribution is often represented using a latent  
35 variable model (see e.g., [TBB19; KAH19]).  
36

### 37 5.4.2 Lossy compression and the rate-distortion tradeoff

38 To encode real-valued signals, such as images and sound, as a digital signal, we first have to quantize  
39 the signal into a sequence of symbols. A simple way to do this is to use vector quantization. We can  
40 then compress this discrete sequence of symbols using lossless coding methods. However, when we  
41 uncompress, we lose some information. Hence this approach is called **lossy compression**.  
42

43 1. For details, see <http://prize.hutter1.net>.

In this section, we quantify this tradeoff between the size of the representation (number of symbols we use), and the resulting error. We will use the terminology of the variational information bottleneck discussed in Section 5.6.2 (except here we are in the unsupervised setting). In particular, we assume we have a stochastic encoder  $p(\mathbf{z}|\mathbf{x})$ , a stochastic decoder  $d(\mathbf{x}|\mathbf{z})$  and a prior marginal  $m(\mathbf{z})$ .

We define the **distortion** of an encoder-decoder pair (as in Section 5.6.2) as follows:

$$D = - \int d\mathbf{x} p(\mathbf{x}) \int d\mathbf{z} e(\mathbf{z}|\mathbf{x}) \log d(\mathbf{x}|\mathbf{z}) \quad (5.152)$$

If the decoder is a deterministic model plus Gaussian noise,  $d(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|f_d(\mathbf{z}), \sigma^2)$ , and the encoder is deterministic,  $e(\mathbf{z}|\mathbf{x}) = \delta(\mathbf{z} - f_e(\mathbf{x}))$ , then this becomes

$$D = \frac{1}{\sigma^2} \mathbb{E}_{p(\mathbf{x})} [| | f_d(f_e(\mathbf{x})) - \mathbf{x} | |^2] \quad (5.153)$$

This is just the expected **reconstruction error** that occurs if we (deterministically) encode and then decode the data using  $f_e$  and  $f_d$ .

We define the **rate** of our model as follows:

$$R = \int d\mathbf{x} p(\mathbf{x}) \int d\mathbf{z} e(\mathbf{z}|\mathbf{x}) \log \frac{e(\mathbf{z}|\mathbf{x})}{m(\mathbf{z})} \quad (5.154)$$

$$= \mathbb{E}_{p(\mathbf{x})} [D_{\text{KL}}(e(\mathbf{z}|\mathbf{x}) \parallel m(\mathbf{z}))] \quad (5.155)$$

$$= \int d\mathbf{x} \int d\mathbf{z} p(\mathbf{x}, \mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{x})m(\mathbf{z})} \geq \mathbb{I}(\mathbf{x}, \mathbf{z}) \quad (5.156)$$

This is just the average KL between our encoding distribution and the marginal. If we use  $m(\mathbf{z})$  to design an optimal code, then the rate is the *excess* number of bits we need to pay to encode our data using  $m(\mathbf{z})$  rather than the true **aggregate posterior**  $p(\mathbf{z}) = \int d\mathbf{x} p(\mathbf{x})e(\mathbf{z}|\mathbf{x})$ .

There is a fundamental tradeoff between the rate and distortion. To see why, note that a trivial encoding scheme would set  $e(\mathbf{z}|\mathbf{x}) = \delta(\mathbf{z} - \mathbf{x})$ , which simply uses  $\mathbf{x}$  as its own best representation. This would incur 0 distortion (and hence maximize the likelihood), but it would incur a high rate, since each  $e(\mathbf{z}|\mathbf{x})$  distribution would be unique, and far from  $m(\mathbf{z})$ . In other words, there would be no compression. Conversely, if  $e(\mathbf{z}|\mathbf{x}) = \delta(\mathbf{z} - \mathbf{0})$ , the encoder would ignore the input. In this case, the rate would be 0, but the distortion would be high.

We can characterize the tradeoff more precisely using the variational lower and upper bounds on the mutual information from Section 5.3.6. From that section, we know that

$$H - D \leq \mathbb{I}(\mathbf{x}; \mathbf{z}) \leq R \quad (5.157)$$

where  $H$  is the (differential) entropy

$$H = - \int d\mathbf{x} p(\mathbf{x}) \log p(\mathbf{x}) \quad (5.158)$$

For discrete data, all probabilities are bounded above by 1, and hence  $H \geq 0$  and  $D \geq 0$ . In addition, the rate is always non-negative,  $R \geq 0$ , since it is the average of a KL divergence. (This is true for either discrete or continuous encodings  $\mathbf{z}$ .) Consequently, we can plot the set of achievable values of  $R$  and  $D$  as shown in Figure 5.10. This is known as a **rate distortion curve**.

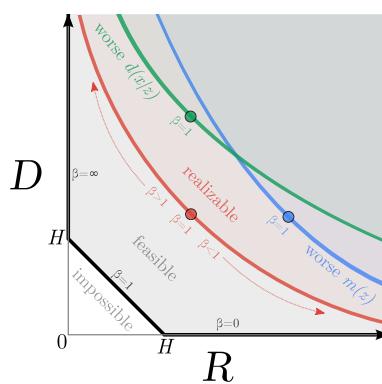


Figure 5.10: Illustration of the rate-distortion tradeoff. See text for details. From Figure 1 of [Ale+18]. Used with kind permission of Alex Alemi.

The bottom horizontal line corresponds to the zero distortion setting,  $D = 0$ , in which we can perfectly encode and decode our data. This can be achieved by using the trivial encoder where  $e(\mathbf{z}|\mathbf{x}) = \delta(\mathbf{z} - \mathbf{x})$ . Shannon's source coding theorem tells us that the minimum number of bits we need to use to encode data in this setting is the entropy of the data, so  $R \geq H$  when  $D = 0$ . If we use a suboptimal marginal distribution  $m(\mathbf{z})$  for coding, we will increase the rate without affecting the distortion.

The left vertical line corresponds to the zero rate setting,  $R = 0$ , in which the latent code is independent of  $\mathbf{z}$ . In this case, the decoder  $d(\mathbf{x}|\mathbf{z})$  is independent of  $\mathbf{z}$ . However, we can still learn a joint probability model  $p(\mathbf{x})$  which does not use latent variables, e.g., this could be an autoregressive model. The minimal distortion such a model could achieve is again the entropy of the data,  $D \geq H$ .

The black diagonal line illustrates solutions that satisfy  $D = H - R$ , where the upper and lower bounds are tight. In practice, we cannot achieve points on the diagonal, since that requires the bounds to be tight, and therefore assumes our models  $e(\mathbf{z}|\mathbf{x})$  and  $d(\mathbf{x}|\mathbf{z})$  are perfect. This is called the “non-parametric limit”. In the finite data setting, we will always incur additional error, so the RD plot will trace a curve which is shifted up, as shown in Figure 5.10.

We can generate different solutions along this curve by minimizing the following objective:

$$J = D + \beta R = \int d\mathbf{x} p(\mathbf{x}) \int d\mathbf{z} e(\mathbf{z}|\mathbf{x}) \left[ -\log d(\mathbf{x}|\mathbf{z}) + \beta \log \frac{e(\mathbf{z}|\mathbf{x})}{m(\mathbf{z})} \right] \quad (5.159)$$

If we set  $\beta = 1$ , and define  $q(\mathbf{z}|\mathbf{x}) = e(\mathbf{z}|\mathbf{x})$ ,  $p(\mathbf{x}|\mathbf{z}) = d(\mathbf{x}|\mathbf{z})$ , and  $p(\mathbf{z}) = m(\mathbf{z})$ , this exactly matches the VAE objective in Section 21.2. To see this, note that the ELBO from Section 10.1.2 can be written as

$$\mathcal{L} = -(D + R) = \mathbb{E}_{p(\mathbf{x})} \left[ \mathbb{E}_{e(\mathbf{z}|\mathbf{x})} [\log d(\mathbf{x}|\mathbf{z})] - \mathbb{E}_{e(\mathbf{z}|\mathbf{x})} \left[ \log \frac{e(\mathbf{z}|\mathbf{x})}{m(\mathbf{z})} \right] \right] \quad (5.160)$$

which we recognize as the expected reconstruction error minus the KL term  $D_{\text{KL}}(e(\mathbf{z}|\mathbf{x}) \parallel m(\mathbf{z}))$ .

If we allow  $\beta \neq 1$ , we recover the  $\beta$ -VAE objective discussed in Section 21.3.1. Note, however, that the  $\beta$ -VAE model cannot distinguish between different solutions on the diagonal line, all of which have

1  $\beta = 1$ . This is because all such models have the same marginal likelihood (and hence same ELBO),  
2 although they differ radically in terms of whether they learn an interesting latent representation or  
3 not. Thus likelihood is not a sufficient metric for comparing the quality of unsupervised representation  
4 learning methods, as discussed in Section 21.3.1.  
5

6 For further discussion on the inherent conflict between rate, distortion and *perception*, see Blau  
7 and Michaeli [BM19]. For techniques for evaluating rate distortion curves for models see Huang, Cao,  
8 and Grosse [HCG20].  
9

### 10 5.4.3 Bits back coding

11 In the previous section we penalized the rate of our code using the average KL divergence,  $\mathbb{E}_{p(\mathbf{x})} [R(\mathbf{x})]$ ,  
12 where  
13

$$\begin{aligned} \text{14} \quad R(\mathbf{x}) &\triangleq \int d\mathbf{z} p(\mathbf{z}|\mathbf{x}) \log \frac{p(\mathbf{z}|\mathbf{x})}{m(\mathbf{z})} = \mathbb{H}_{ce}(p(\mathbf{z}|\mathbf{x}), m(\mathbf{z})) - \mathbb{H}(p(\mathbf{z}|\mathbf{x})). \end{aligned} \quad (5.161) \quad \text{15}$$

16 The first term is the cross entropy, which is the expected number of bits we need to encode  $\mathbf{x}$ ; the  
17 second term is the entropy, which is the minimum number of bits. Thus we are penalizing the *excess*  
18 number of bits required to communicate the code to a receiver. How come we don't have to "pay for"  
19 the actual (total) number of bits we use, which is the cross entropy?  
20

21 The reason is that we could in principle get the bits needed by the optimal code given back to  
22 us; this is called **bits back coding** [HC93; FH97]. The argument goes as follows. Imagine Alice is  
23 trying to (losslessly) communicate some data, such as an image  $\mathbf{x}$ , to Bob. Before they went their  
24 separate ways, both Alice and Bob decided to share their encoder  $p(\mathbf{z}|\mathbf{x})$ , marginal  $m(\mathbf{z})$  and decoder  
25 distributions  $d(\mathbf{x}|\mathbf{z})$ . To communicate an image, Alice will use a **two part code**. First, she will  
26 sample a code  $\mathbf{z} \sim p(\mathbf{z}|\mathbf{x})$  from her encoder, and communicate that to Bob over a channel designed  
27 to efficiently encode samples from the marginal  $m(\mathbf{z})$ ; this costs  $-\log_2 m(\mathbf{z})$  bits. Next Alice will  
28 use her decoder  $d(\mathbf{x}|\mathbf{z})$  to compute the residual error, and losslessly send that to Bob at the cost of  
29  $-\log_2 d(\mathbf{x}|\mathbf{z})$  bits. The expected total number of bits required here is what we naively expected:  
30

$$\text{31} \quad \mathbb{E}_{p(\mathbf{z}|\mathbf{x})} [-\log_2 d(\mathbf{x}|\mathbf{z}) - \log_2 m(\mathbf{z})] = D + \mathbb{H}_{ce}(p(\mathbf{z}|\mathbf{x}), m(\mathbf{z})). \quad (5.162) \quad \text{32}$$

33 We see that this is the distortion plus cross entropy, not distortion plus rate. So how do we get the  
34 bits back, to convert the cross entropy to a rate term?

35 The trick is that Bob actually receives more information than we suspected. Bob can use the code  
36  $\mathbf{z}$  and the residual error to perfectly reconstruct  $\mathbf{x}$ . However, Bob also knows what specific code Alice  
37 sent,  $\mathbf{z}$ , as well as what encoder she used,  $p(\mathbf{z}|\mathbf{x})$ . When Alice drew the sample code  $\mathbf{z} \sim p(\mathbf{z}|\mathbf{x})$ , she  
38 had to use some kind of entropy source in order to generate the random sample. Suppose she did it  
39 by picking words sequentially from a compressed copy of Moby Dick, in order to generate a stream  
40 of random bits. On Bob's end, he can reverse engineer all of the sampling bits, and thus recover the  
41 compressed copy of Moby Dick! Thus Alice can use the extra randomness in the choice of  $\mathbf{z}$  to share  
42 more information.

43 While in the original formulation the bits-back argument was largely theoretical, offering a thought  
44 experiment for why we should penalize our models with the KL instead of the cross entropy, recently  
45 several practical real world algorithms have been developed that actually achieve the bits-back goal.  
46 These include [HHLMF18; AT20; TBB19; YBM20; HLA19].  
47

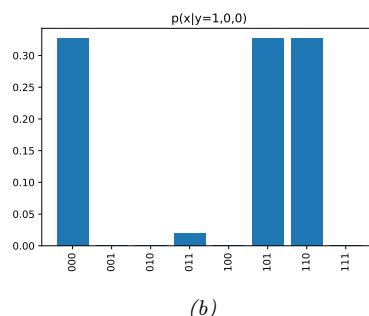
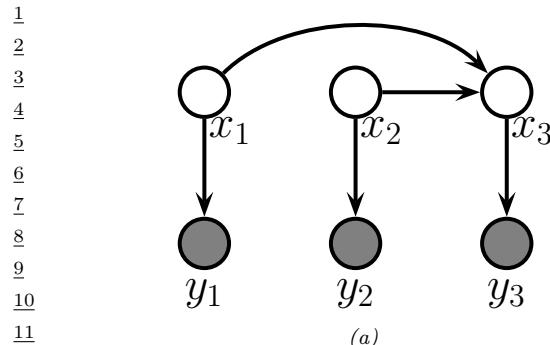


Figure 5.11: (a) A simple error-correcting code DPGM.  $x_i$  are the sent bits,  $y_i$  are the received bits.  $x_3$  is an even parity check bit computed from  $x_1$  and  $x_2$ . (b) Posterior over codewords given that  $\mathbf{y} = (1, 0, 0)$ ; the probability of a bit flip is 0.2. Generated by [error\\_correcting\\_code\\_demo.ipynb](#).

## 5.5 Error-correcting codes (channel coding)

The idea behind **error correcting codes** is to add redundancy to a signal  $\mathbf{x}$  (which is the result of encoding the original data), such that when it is sent over to the receiver via a noisy transmission line (such as a cell phone connection), the receiver can recover from any corruptions that might occur to the signal. This is called **channel coding**.

In more detail, let  $\mathbf{x} \in \{0, 1\}^m$  be the source message, where  $m$  is called the **block length**. Let  $\mathbf{y}$  be the result of sending  $\mathbf{x}$  over a **noisy channel**. This is a corrupted version of the message. For example, each message bit may get flipped independently with probability  $\alpha$ , in which case  $p(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^m p(y_i|x_i)$ , where  $p(y_i|x_i=0) = [1 - \alpha, \alpha]$  and  $p(y_i|x_i=1) = [\alpha, 1 - \alpha]$ . Alternatively, we may add Gaussian noise, so  $p(y_i|x_i=b) = \mathcal{N}(y_i|\mu_b, \sigma^2)$ . The receiver's goal is to infer the true message from the noisy observations, i.e., to compute  $\operatorname{argmax}_{\mathbf{x}} p(\mathbf{x}|\mathbf{y})$ .

A common way to increase the chance of being able to recover the original signal is to add **parity check bits** to it before sending it. These are deterministic functions of the original signal, which specify if the sum of the input bits is odd or even. This provides a form of **redundancy**, so that if one bit is corrupted, we can still infer its value, assuming the other bits are not flipped. (This is reasonable since we assume the bits are corrupted independently at random, so it is less likely that multiple bits are flipped than just one bit.)

For example, suppose we have two original message bits, and we add one parity bit. This can be modeled using a directed graphical model as shown in Figure 5.11(a). This graph encodes the following joint probability distribution:

$$p(\mathbf{x}, \mathbf{y}) = p(x_1)p(x_2)p(x_3|x_1, x_2) \prod_{i=1}^3 p(y_i|x_i) \quad (5.163)$$

The priors  $p(x_1)$  and  $p(x_2)$  are uniform. The conditional term  $p(x_3|x_1, x_2)$  is deterministic, and computes the parity of  $(x_1, x_2)$ . In particular, we have  $p(x_3 = 1|x_1, x_2) = 1$  if the total number of 1s in the block  $x_{1:2}$  is odd. The likelihood terms  $p(y_i|x_i)$  represent a bit flipping noisy channel model, with noise level  $\alpha = 0.2$ .

1 Suppose we observe  $\mathbf{y} = (1, 0, 0)$ . We know that this cannot be what the sender sent, since this  
 2 violates the parity constraint (if  $x_1 = 1$  then we know  $x_3 = 1$ ). Instead, the 3 posterior modes for  $\mathbf{x}$   
 3 are 000 (first bit was flipped), 110 (second bit was flipped), and 101 (third bit was flipped). The only  
 4 other configuration with non-zero support in the posterior is 011, which corresponds to the much less  
 5 likely hypothesis that three bits were flipped (see Figure 5.11(b)). All other hypotheses (001, 010 and  
 6 100) are inconsistent with the deterministic method used to create codewords. (See Section 9.2.3.2  
 7 for further discussion of this point.)  
 8

9 In practice, we use more complex coding schemes that are more efficient, in the sense that they  
 10 add less redundant bits to the message, but still guarantee that errors can be corrected. For details,  
 11 see Section 9.3.8.  
 12

## 13 5.6 The information bottleneck

14 In this section, we discuss discriminative models  $p(\mathbf{y}|\mathbf{x})$  that use a *stochastic bottleneck* between the  
 15 input  $\mathbf{x}$  and the output  $\mathbf{y}$  to prevent overfitting, and improve robustness and calibration.  
 16

### 18 5.6.1 Vanilla IB

19 We say that  $\mathbf{z}$  is a **representation** of  $\mathbf{x}$  if  $\mathbf{z}$  is a (possibly stochastic) function of  $\mathbf{x}$ , and hence can  
 20 be described by the conditional  $p(\mathbf{z}|\mathbf{x})$ . We say that a representation  $\mathbf{z}$  of  $\mathbf{x}$  is **sufficient** for task  $\mathbf{y}$   
 21 if  $\mathbf{y} \perp \mathbf{x} | \mathbf{z}$ , or equivalently, if  $\mathbb{I}(\mathbf{z}; \mathbf{y}) = \mathbb{I}(\mathbf{x}; \mathbf{y})$ , i.e.,  $\mathbb{H}(\mathbf{y}|\mathbf{z}) = \mathbb{H}(\mathbf{y}|\mathbf{x})$ . We say that a representation  
 22 is a **minimal sufficient statistic** if  $\mathbf{z}$  is sufficient and there is no other  $\mathbf{z}$  with smaller  $\mathbb{I}(\mathbf{z}; \mathbf{x})$  value.  
 23 Thus we would like to find a representation  $\mathbf{z}$  that maximizes  $\mathbb{I}(\mathbf{z}; \mathbf{y})$  while minimizing  $\mathbb{I}(\mathbf{z}; \mathbf{x})$ . That  
 24 is, we would like to optimize the following objective:  
 25

$$26 \quad \min \beta \mathbb{I}(\mathbf{z}; \mathbf{x}) - \mathbb{I}(\mathbf{z}; \mathbf{y}) \quad (5.164)$$

27 where  $\beta \geq 0$ , and we optimize wrt the distributions  $p(\mathbf{z}|\mathbf{x})$  and  $p(\mathbf{y}|\mathbf{z})$ . This is called the **information**  
 28 **bottleneck principle** [TPB99]. This generalizes the concept of minimal sufficient statistic to take  
 29 into account that there is a tradeoff between sufficiency and minimality, which is captured by the  
 30 Lagrange multiplier  $\beta > 0$ .  
 31

32 This principle is illustrated in Figure 5.12. We assume  $Z$  is a function of  $X$ , but is independent  
 33 of  $Y$ , i.e., we assume the graphical model  $Z \leftarrow X \leftrightarrow Y$ . This corresponds to the following joint  
 34 distribution:  
 35

$$36 \quad p(\mathbf{x}, \mathbf{y}, \mathbf{z}) = p(\mathbf{z}|\mathbf{x})p(\mathbf{y}|\mathbf{z})p(\mathbf{x}) \quad (5.165)$$

37 Thus  $Z$  can capture any amount of information about  $X$  that it wants, but cannot contain information  
 38 that is unique to  $Y$ , as illustrated in Figure 5.12a. The optimal representation only captures  
 39 information about  $X$  that is useful for  $Y$ ; to prevent us “wasting capacity” and fitting irrelevant  
 40 details of the input,  $Z$  should also minimize information about  $X$ , as shown in Figure 5.12b.  
 41

42 If all the random variables are discrete, and  $\mathbf{z} = e(\mathbf{x})$  is a deterministic function of  $\mathbf{x}$ , then the  
 43 algorithm of [TPB99] can be used to minimize the IB objective in Section 5.6. The objective can  
 44 also be solved analytically if all variables are jointly Gaussian [Che+05] (the resulting method can be  
 45 viewed as a form of supervised PCA). But in general, it is intractable to solve this problem exactly.  
 46 We discuss a tractable approximation in Section 5.6.2. (More details can be found in e.g., [SZ22].)  
 47

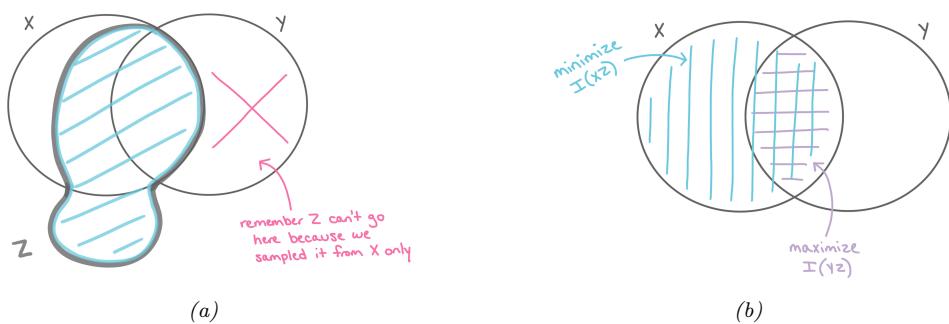


Figure 5.12: Information diagrams for information bottleneck. (a)  $Z$  can contain any amount of information about  $X$  (whether it useful for predicting  $Y$  or not), but it cannot contain information about  $Y$  that is not shared with  $X$ . (b) The optimal representation for  $Z$  maximizes  $\mathbb{I}(Z, Y)$  and minimizes  $\mathbb{I}(Z, X)$ . Used with kind permission of Katie Everett.

### 5.6.2 Variational IB

In this section, we derive a variational upper bound on Equation (5.164), leveraging ideas from Section 5.3.6. This is called the **variational IB** or **VIB** method [Ale+16]. The key trick will be to use the non-negativity of the KL divergence to write

$$\int d\mathbf{x} p(\mathbf{x}) \log p(\mathbf{x}) \geq \int d\mathbf{x} p(\mathbf{x}) \log q(\mathbf{x}) \quad (5.166)$$

for any distribution  $q$ . (Note that both  $p$  and  $q$  may be conditioned on other variables.)

To explain the method in more detail, let us define the following notation. Let  $e(\mathbf{z}|\mathbf{x}) = p(\mathbf{z}|\mathbf{x})$  represent the encoder,  $b(\mathbf{z}|\mathbf{y}) \approx p(\mathbf{z}|\mathbf{y})$  represent the backwards encoder,  $d(\mathbf{z}|\mathbf{y}) \approx p(\mathbf{z}|\mathbf{y})$  represent the classifier (decoder), and  $m(\mathbf{z}) \approx p(\mathbf{z})$  represent the marginal. (Note that we get to choose  $p(\mathbf{z}|\mathbf{x})$ , but the other distributions are derived by approximations of the corresponding marginals and conditionals of the exact joint  $p(\mathbf{x}, \mathbf{y}, \mathbf{z})$ .) Also, let  $\langle \cdot \rangle$  represent expectations wrt the relevant terms from the  $p(\mathbf{x}, \mathbf{y}, \mathbf{z})$  joint.

With this notation, we can derive a lower bound on  $\mathbb{I}(\mathbf{z}; \mathbf{y})$  as follows:

$$\mathbb{I}(\mathbf{z}; \mathbf{y}) = \int d\mathbf{y} d\mathbf{z} p(\mathbf{y}, \mathbf{z}) \log \frac{p(\mathbf{y}, \mathbf{z})}{p(\mathbf{y})p(\mathbf{z})} \quad (5.167)$$

$$= \int d\mathbf{y} d\mathbf{z} p(\mathbf{y}, \mathbf{z}) \log p(\mathbf{y}|\mathbf{z}) - \int d\mathbf{y} d\mathbf{z} p(\mathbf{y}, \mathbf{z}) \log p(\mathbf{y}) \quad (5.168)$$

$$= \int d\mathbf{y} d\mathbf{z} p(\mathbf{z}) p(\mathbf{y}|\mathbf{z}) \log p(\mathbf{y}|\mathbf{z}) - \text{const} \quad (5.169)$$

$$\geq \int d\mathbf{y} d\mathbf{z} p(\mathbf{y}, \mathbf{z}) \log d(\mathbf{y}|\mathbf{z}) \quad (5.170)$$

$$= \langle \log d(\mathbf{y}|\mathbf{z}) \rangle \quad (5.171)$$

where we exploited the fact that  $\mathbb{H}(p(\mathbf{y}))$  is a constant that is independent of our representation.

1 Note that we can approximate the expectations by sampling from  
2

3  $p(\mathbf{y}, \mathbf{z}) = \int d\mathbf{x} p(\mathbf{x})p(\mathbf{y}|\mathbf{x})p(\mathbf{z}|\mathbf{x}) = \int d\mathbf{x} p(\mathbf{x}, \mathbf{y})e(\mathbf{z}|\mathbf{x})$  (5.172)  
4

5 This is just the empirical distribution “pushed through” the encoder.  
6

7 Similarly, we can derive an upper bound on  $\mathbb{I}(\mathbf{z}; \mathbf{x})$  as follows:

8  $\mathbb{I}(\mathbf{z}; \mathbf{x}) = \int dz dx p(\mathbf{x}, \mathbf{z}) \log \frac{p(\mathbf{z}, \mathbf{x})}{p(\mathbf{x})p(\mathbf{z})}$  (5.173)  
9

10  $= \int dz dx p(\mathbf{x}, \mathbf{z}) \log p(\mathbf{z}|\mathbf{x}) - \int dz p(\mathbf{z}) \log p(\mathbf{z})$  (5.174)  
11

12  $\leq \int dz dx p(\mathbf{x}, \mathbf{z}) \log p(\mathbf{z}|\mathbf{x}) - \int dz p(\mathbf{z}) \log m(\mathbf{z})$  (5.175)  
13

14  $= \int dz dx p(\mathbf{x}, \mathbf{z}) \log \frac{e(\mathbf{z}|\mathbf{x})}{m(\mathbf{z})}$  (5.176)  
15

16  $= \langle \log e(\mathbf{z}|\mathbf{x}) \rangle - \langle m(\mathbf{z}) \rangle$  (5.177)  
17

18 Note that we can approximate the expectations by sampling from  $p(\mathbf{x}, \mathbf{z}) = p(\mathbf{x})p(\mathbf{z}|\mathbf{x})$ .  
19

20 Putting it altogether, we get the following upper bound on the IB objective:

21  $\beta \mathbb{I}(\mathbf{x}; \mathbf{z}) - \mathbb{I}(\mathbf{z}; \mathbf{y}) \leq \beta (\langle \log e(\mathbf{z}|\mathbf{x}) \rangle - \langle \log m(\mathbf{z}) \rangle) - \langle \log d(\mathbf{y}|\mathbf{z}) \rangle$  (5.178)  
22

23 Thus the VIB objective is

24  $\mathcal{L}_{\text{VIB}} = \beta (\mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})e(\mathbf{z}|\mathbf{x})} [\log e(\mathbf{z}|\mathbf{x}) - \log m(\mathbf{z})] - \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})e(\mathbf{z}|\mathbf{x})d(\mathbf{y}|\mathbf{z})} [\log d(\mathbf{y}|\mathbf{z})])$  (5.179)  
25

26  $= -\mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})e(\mathbf{z}|\mathbf{x})d(\mathbf{y}|\mathbf{z})} [\log d(\mathbf{y}|\mathbf{z})] + \beta \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [D_{\text{KL}}(e(\mathbf{z}|\mathbf{x}) \parallel m(\mathbf{z}))]$  (5.180)  
27

28 We can now take stochastic gradients of this objective and minimize it (wrt the parameters of  
29 the encoder, decoder and marginal) using SGD. (We assume the distributions are reparameterizable,  
30 as discussed in Section 6.5.4.) For the encoder  $e(\mathbf{z}|\mathbf{x})$ , we often use a conditional Gaussian, and  
31 for the decoder  $d(\mathbf{y}|\mathbf{z})$ , we often use a softmax classifier. For the marginal,  $m(\mathbf{z})$ , we should use  
32 a flexible model, such as a mixture of Gaussians, since it needs to approximate the **aggregated**  
33 **posterior**  $p(\mathbf{z}) = \int dz p(\mathbf{x})e(\mathbf{z}|\mathbf{x})$ , which is a mixture of  $N$  Gaussians (assuming  $p(\mathbf{x})$  is an empirical  
34 distribution with  $N$  samples, and  $e(\mathbf{z}|\mathbf{x})$  is a Gaussian).

35 We illustrate this in Figure 5.13, where we fit an MLP model to MNIST. We use a 2d bottleneck  
36 layer before passing to the softmax. In panel a, we show the embedding learned by a deterministic  
37 encoder. We see that each image gets mapped to a point, and there is little overlap between classes,  
38 or between instances. In panels b-c, we show the embedding learned by a stochastic encoder. Each  
39 image gets mapped to a Gaussian distribution, we show the mean and the covariance separately. The  
40 classes are still well separated, but individual instances of a class are no longer distinguishable, since  
41 such information is not relevant for prediction purposes.

### 42 5.6.3 Conditional entropy bottleneck

43 The IB tries to maximize  $\mathbb{I}(Z; Y)$  while minimizing  $\mathbb{I}(Z; X)$ . We can write this objective as  
44

45  $\min \mathbb{I}(\mathbf{x}; \mathbf{z}) - \lambda \mathbb{I}(\mathbf{y}; \mathbf{z})$  (5.181)  
46

47

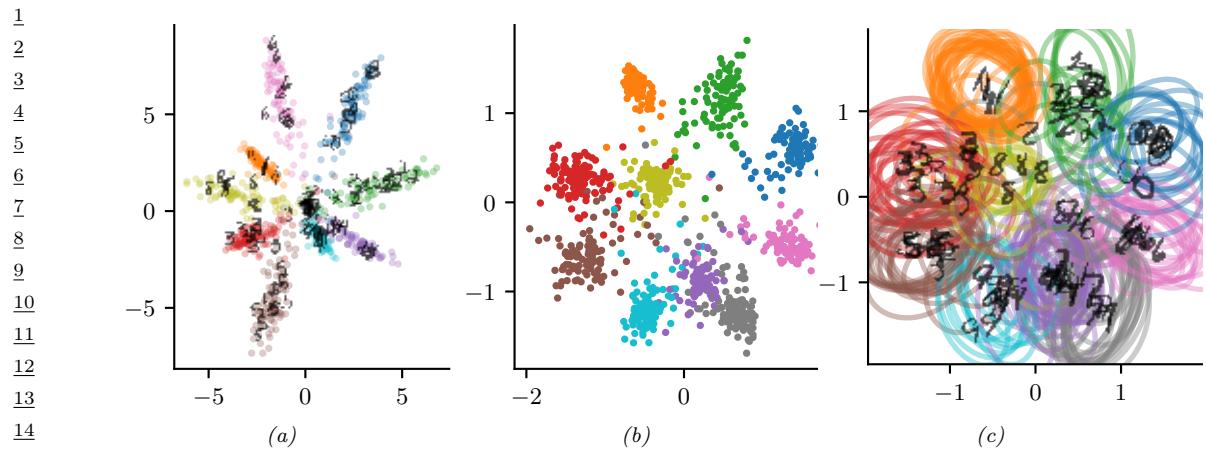


Figure 5.13: 2d embeddings of MNIST digits created by an MLP classifier. (a) Deterministic model. (b-c) VIB model, means and covariances. Generated by [vib\\_demo.ipynb](#). Used with kind permission of Alex Alemi.

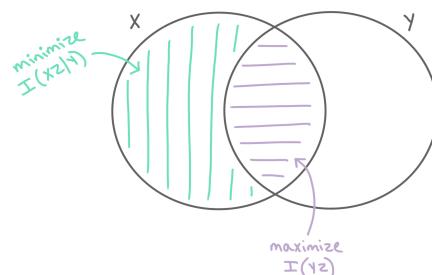


Figure 5.14: Conditional entropy bottleneck (CEB) chooses a representation  $Z$  that maximizes  $I(Z; Y)$  and minimizes  $I(X, Z|Y)$ . Used with kind permission of Katie Everett.

for  $\lambda \geq 0$ . However, we see from the information diagram in Figure 5.12b that  $I(Z; X)$  contains some information that is relevant to  $Y$ . A sensible alternative objective is to minimizes the residual mutual information,  $I(X; Z|Y)$ . This gives rise to the following objective:

$$\min I(\mathbf{x}; \mathbf{z}|\mathbf{y}) - \lambda' I(\mathbf{y}; \mathbf{z}) \quad (5.182)$$

for  $\lambda' \geq 0$ . This is known as the **conditional entropy bottleneck** or **CEB** [Fis20]. See Figure 5.14 for an illustration.

Since  $I(\mathbf{x}; \mathbf{z}|\mathbf{y}) = I(\mathbf{x}; \mathbf{z}) - I(\mathbf{y}; \mathbf{z})$ , we see that the CEB is equivalent to standard IB with  $\lambda' = \lambda + 1$ . However, it is easier to upper bound  $I(\mathbf{x}; \mathbf{z}|\mathbf{y})$  than  $I(\mathbf{x}; \mathbf{z})$ , since we are conditioning on  $\mathbf{y}$ , which

1 provides information about  $\mathbf{z}$ . In particular, we have  
2

$$\mathbb{I}(\mathbf{x}; \mathbf{z}|\mathbf{y}) = \mathbb{I}(\mathbf{x}; \mathbf{z}) - \mathbb{I}(\mathbf{y}; \mathbf{z}) \quad (5.183)$$

$$= \mathbb{H}(\mathbf{z}) - \mathbb{H}(\mathbf{z}|\mathbf{x}) - [\mathbb{H}(\mathbf{z}) - \mathbb{H}(\mathbf{z}|\mathbf{y})] \quad (5.184)$$

$$= -\mathbb{H}(\mathbf{z}|\mathbf{x}) - \mathbb{H}(\mathbf{z}|\mathbf{y}) \quad (5.185)$$

$$= \int d\mathbf{z}d\mathbf{x} p(\mathbf{x}, \mathbf{z}) \log p(\mathbf{z}|\mathbf{x}) - \int d\mathbf{z}d\mathbf{y} p(\mathbf{z}, \mathbf{y}) \log p(\mathbf{z}|\mathbf{y}) \quad (5.186)$$

$$\leq \int d\mathbf{z}d\mathbf{x} p(\mathbf{x}, \mathbf{z}) \log e(\mathbf{z}|\mathbf{x}) - \int d\mathbf{z}d\mathbf{y} p(\mathbf{z}, \mathbf{y}) \log b(\mathbf{z}|\mathbf{y}) \quad (5.187)$$

$$= \langle \log e(\mathbf{z}|\mathbf{x}) \rangle - \langle \log b(\mathbf{z}|\mathbf{y}) \rangle \quad (5.188)$$

13 Putting it altogether, we get the final CEB objective:  
14

$$\min \beta (\langle \log e(\mathbf{z}|\mathbf{x}) \rangle - \langle \log b(\mathbf{z}|\mathbf{y}) \rangle) - \langle \log d(\mathbf{y}|\mathbf{z}) \rangle \quad (5.189)$$

15 Note that it is generally easier to learn the conditional backwards encoder  $b(\mathbf{z}|\mathbf{y})$  than the  
16 unconditional marginal  $m(\mathbf{z})$ . Also, we know that the tightest upper bound occurs when  $\mathbb{I}(\mathbf{x}; \mathbf{z}|\mathbf{y}) =$   
17  $\mathbb{I}(\mathbf{x}; \mathbf{z}) - \mathbb{I}(\mathbf{y}; \mathbf{z}) = 0$ . The corresponding value of  $\beta$  corresponds to an optimal representation. By  
18 contrast, it is not clear how to measure distance from optimality when using IB.  
19

20212223242526272829303132333435363738394041424344454647

# 6 Optimization

## 6.1 Introduction

In this chapter, we consider solving **optimization problems** of various forms. Abstractly these can all be written as

$$\boldsymbol{\theta}^* \in \operatorname{argmin}_{\boldsymbol{\theta} \in \Theta} \mathcal{L}(\boldsymbol{\theta}) \quad (6.1)$$

where  $\mathcal{L} : \Theta \rightarrow \mathbb{R}$  is the objective or loss function, and  $\Theta$  is the parameter space we are optimizing over. However, this abstraction hides many details, such as whether the problem is constrained or unconstrained, discrete or continuous, convex or non-convex, etc. In the prequel to this book, [Mur22], we discussed some simple optimization algorithms for some common problems that arise in machine learning. In this chapter, we discuss some more advanced methods. For more details on optimization, please consult some of the many excellent textbooks, such as [KW19b; BV04; NW06; Ber15; Ber16] as well as various review articles, such as [BCN18; Sun+19b; PPS18; Pey20].

## 6.2 Automatic differentiation

*This section was written by Roy Frostig.*

This section is concerned with computing (partial) derivatives of complicated functions in an automatic manner. By “complicated” we mean those expressed as a composition of an arbitrary number of more basic operations, such as in deep neural networks. This task is known as **automatic differentiation (AD)**, or **autodiff**. AD is an essential component in optimization and deep learning, and is also used in several other fields across science and engineering. See e.g. Baydin et al. [Bay+15] for a review focused on machine learning and Griewank and Walther [GW08] for a classical textbook.

### 6.2.1 Differentiation in functional form

Before covering automatic differentiation, it is useful to review the mathematics of differentiation. We will use a particular **functional** notation for partial derivatives, rather than the typical one used throughout much of this book. We will refer to the latter as the **named variable** notation for the moment. Named variable notation relies on associating function arguments with names. For instance, given a function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ , the partial derivative of  $f$  with respect to its first scalar argument, at a

1 point  $\mathbf{a} = (a_1, a_2)$ , might be written:  
2

$$\frac{\partial f}{\partial x_1} \Big|_{\mathbf{x}=\mathbf{a}} \tag{6.2}$$

3  
4 This notation is not entirely self-contained. It refers to a name  $\mathbf{x} = (x_1, x_2)$ , implicit or inferred from  
5 context, suggesting the argument of  $f$ . An alternative expression is:  
6

$$\frac{\partial}{\partial a_1} f(a_1, a_2) \tag{6.3}$$

7  
8 where now  $a_1$  serves both as an argument name (or a symbol in an expression) and as a particular  
9 evaluation point. Tracking names can become an increasingly complicated endeavor as we compose  
10 many functions together, each possibly taking several arguments.  
11

12 A functional notation instead defines derivatives as operators on functions. If a function has  
13 multiple arguments, they are identified by position rather than by name, alleviating the need for  
14 auxiliary variable definitions. Some of the following definitions draw on those in Spivak's *Calculus*  
15 on Manifolds [Spi71], in Sussman and Wisdom's *Functional Differential Geometry* [SW13], and  
16 generally appear more regularly in accounts of differential calculus and geometry. These texts are  
17 recommended for a more formal treatment, and a more mathematically general view, of the material  
18 briefly covered in this section.  
19

20 Beside notation, we will rely on some basic multivariable calculus concepts. This includes the  
21 notion of (partial) derivatives, the differential or Jacobian of a function at a point, its role as a linear  
22 approximation local to the point, and various properties of linear maps, matrices, and transposition.  
23 We will focus on a finite-dimensional setting and write  $\{\mathbf{e}_1, \dots, \mathbf{e}_n\}$  for the standard basis in  $\mathbb{R}^n$ .  
24

25  
26 **Linear and multilinear functions.** We use  $F : \mathbb{R}^n \multimap \mathbb{R}^m$  to denote a function  $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$   
27 that is linear, and by  $F[\mathbf{x}]$  its application to  $\mathbf{x} \in \mathbb{R}^n$ . Recall that such a linear map corresponds  
28 to a matrix in  $\mathbb{R}^{m \times n}$  whose columns are  $F[\mathbf{e}_1], \dots, F[\mathbf{e}_n]$ ; both interpretations will prove useful.  
29 Conveniently, function composition and matrix multiplication expressions look similar: to compose  
30 two linear maps  $F$  and  $G$  we can write  $F \circ G$  or, barely abusing notation, consider the matrix  $FG$ .  
31 Every linear map  $F : \mathbb{R}^n \multimap \mathbb{R}^m$  has a transpose  $F : \mathbb{R}^m \multimap \mathbb{R}^n$ , which is another linear map identified  
32 with transposing the corresponding matrix.  
33

34 Repeatedly using the linear arrow symbol, we can denote by:  
35

$$T : \underbrace{\mathbb{R}^n \multimap \cdots \multimap \mathbb{R}^n}_{k \text{ times}} \multimap \mathbb{R}^m \tag{6.4}$$

36  
37 a multilinear, or more specifically  $k$ -linear, map:  
38

$$T : \underbrace{\mathbb{R}^n \times \cdots \times \mathbb{R}^n}_{k \text{ times}} \rightarrow \mathbb{R}^m \tag{6.5}$$

39  
40 which corresponds to an array (or tensor) in  $\mathbb{R}^{m \times n \times \cdots \times n}$ . We denote by  $T[\mathbf{x}_1, \dots, \mathbf{x}_k] \in \mathbb{R}^m$  the  
41 application of such a  $k$ -linear map to vectors  $\mathbf{x}_1, \dots, \mathbf{x}_k \in \mathbb{R}^n$ .  
42

**The derivative operator.** For an open set  $U \subset \mathbb{R}^n$  and a differentiable function  $f : U \rightarrow \mathbb{R}^m$ , denote its **derivative function**:

$$\partial f : U \rightarrow (\mathbb{R}^n \multimap \mathbb{R}^m) \quad (6.6)$$

or equivalently  $\partial f : U \rightarrow \mathbb{R}^{m \times n}$ . This function maps a point  $\mathbf{x} \in U$  to the Jacobian of all partial derivatives evaluated at  $\mathbf{x}$ . The symbol  $\partial$  itself denotes the **derivative operator**, a function mapping functions to their derivative functions. When  $m = 1$ , the map  $\partial f(\mathbf{x})$  recovers the standard gradient  $\nabla f(\mathbf{x})$  at any  $\mathbf{x} \in U$ , by considering the matrix view of the former. Indeed, the nabla symbol  $\nabla$  is sometimes described as an operator as well, such that  $\nabla f$  is a function. When  $n = m = 1$ , the Jacobian is scalar-valued, and  $\partial f$  is the familiar derivative  $f'$ .

In the expression  $\partial f(\mathbf{x})[\mathbf{v}]$ , we will sometimes refer to the argument  $\mathbf{x}$  as the **linearization point** for the Jacobian, and to  $\mathbf{v}$  as the **perturbation**. We call the map:

$$(\mathbf{x}, \mathbf{v}) \mapsto \partial f(\mathbf{x})[\mathbf{v}] \quad (6.7)$$

over linearization points  $\mathbf{x} \in U$  and *input* perturbations  $\mathbf{v} \in \mathbb{R}^n$  the **Jacobian-vector product (JVP)**. We similarly call its transpose:

$$(\mathbf{x}, \mathbf{u}) \mapsto \partial f(\mathbf{x})^\top[\mathbf{u}] \quad (6.8)$$

over linearization points  $\mathbf{x} \in U$  and *output* perturbations  $\mathbf{u} \in \mathbb{R}^m$  the **vector-Jacobian product (VJP)**.

Thinking about maps instead of matrices can help us define higher-order derivatives recursively, as we proceed to do below. It separately suggests how the action of a Jacobian is commonly written in code. When we consider writing  $\partial f(\mathbf{x})$  in a program for a fixed  $\mathbf{x}$ , we often implement it as a function that carries out multiplication by the Jacobian matrix, i.e.  $\mathbf{v} \mapsto \partial f(\mathbf{x})[\mathbf{v}]$ , instead of explicitly representing it as a matrix of numbers in memory. Going a step further, for that matter, we often implement  $\partial f$  as an entire JVP at once, i.e. over any linearization point  $\mathbf{x}$  and perturbation  $\mathbf{v}$ . As a toy example with scalars, consider the cosine:

$$(x, v) \mapsto \partial \cos(x)v = -v \sin(x) \quad (6.9)$$

If we express this at once in code, we can, say, avoid computing  $\sin(x)$  whenever  $v = 0$ .<sup>1</sup>

**Higher-order derivatives.** Suppose the function  $f$  above remains arbitrarily differentiable over its domain  $U \subset \mathbb{R}^n$ . To take another derivative, we write:

$$\partial^2 f : U \rightarrow (\mathbb{R}^n \multimap \mathbb{R}^n \multimap \mathbb{R}^m) \quad (6.10)$$

where  $\partial^2 f(\mathbf{x})$  is a bilinear map representing all second-order partial derivatives. In named variable notation, one might write  $\frac{\partial f(\mathbf{x})}{\partial x_i \partial x_j}$  to refer to  $\partial^2 f(\mathbf{x})[\mathbf{e}_i, \mathbf{e}_j]$ , for example.

<sup>1</sup> 1. This example ignores that such an optimization might be done (best) by a compiler. Then again, for more complex examples, implementing  $(\mathbf{x}, \mathbf{v}) \mapsto \partial f(\mathbf{x})[\mathbf{v}]$  as a single subroutine can help guide compiler optimizations all the same.

The second derivative function  $\partial^2 f$  can be treated coherently as the outcome of applying the derivative operator twice. That is, it makes sense to say that  $\partial^2 = \partial \circ \partial$ . This observation extends recursively to cover arbitrary higher-order derivatives. For  $k \geq 1$ :

$$\partial^k f : U \rightarrow (\underbrace{\mathbb{R}^n \multimap \dots \multimap \mathbb{R}^n}_{k \text{ times}} \multimap \mathbb{R}^m) \quad (6.11)$$

is such that  $\partial^k f(\mathbf{x})$  is a  $k$ -linear map.

With  $m = 1$ , the map  $\partial^2 f(\mathbf{x})$  corresponds to the Hessian matrix at any  $\mathbf{x} \in U$ . Although Jacobians and Hessians suffice to make sense of many machine learning techniques, arbitrary higher-order derivatives are not hard to come by either (e.g. [Kel+20]). As an example, they appear when writing down something as basic as a function's Taylor series approximation, which we can express with our derivative operator as:

$$f(\mathbf{x} + \mathbf{v}) \approx f(\mathbf{x}) + \partial f(\mathbf{x})[\mathbf{v}] + \frac{1}{2!} \partial^2 f(\mathbf{x})[\mathbf{v}, \mathbf{v}] + \dots + \frac{1}{k!} \partial^k f(\mathbf{x})[\mathbf{v}, \dots, \mathbf{v}] \quad (6.12)$$

**Multiple inputs.** Now consider a function of two arguments:

$$g : U \times V \rightarrow \mathbb{R}^m. \quad (6.13)$$

where  $U \subset \mathbb{R}^{n_1}$  and  $V \subset \mathbb{R}^{n_2}$ . For our purposes, a product domain like  $U \times V$  mainly serves to suggest a convenient partitioning of a function's input components. It is isomorphic to a subset of  $\mathbb{R}^{n_1+n_2}$ , corresponding to a single-input function. The latter tells us how the derivative functions of  $g$  ought to look, based on previous definitions, and we will swap between the two views with little warning. Multiple inputs tend to arise in the context of computational circuits and programs: many functions in code are written to accept multiple arguments, and many basic operations (such as  $+$ ) do the same.

With multiple inputs, we can denote by  $\partial_i g$  the derivative function with respect to the  $i$ 'th argument:

$$\partial_1 g : \mathbb{R}^{n_1} \times \mathbb{R}^{n_2} \rightarrow (\mathbb{R}^{n_1} \multimap \mathbb{R}^m), \text{ and} \quad (6.14)$$

$$\partial_2 g : \mathbb{R}^{n_1} \times \mathbb{R}^{n_2} \rightarrow (\mathbb{R}^{n_2} \multimap \mathbb{R}^m). \quad (6.15)$$

Under the matrix view, the function  $\partial_1 g$  maps a pair of points  $\mathbf{x} \in \mathbb{R}^{n_1}$  and  $\mathbf{y} \in \mathbb{R}^{n_2}$  to the matrix of all partial derivatives of  $g$  with respect to its first argument, evaluated at  $(\mathbf{x}, \mathbf{y})$ . We take  $\partial g$  with no subscript to simply mean the concatenation of  $\partial_1 g$  and  $\partial_2 g$ :

$$\partial g : \mathbb{R}^{n_1} \times \mathbb{R}^{n_2} \rightarrow (\mathbb{R}^{n_1} \times \mathbb{R}^{n_2} \multimap \mathbb{R}^m) \quad (6.16)$$

where, for every linearization point  $(\mathbf{x}, \mathbf{y}) \in U \times V$  and perturbations  $\dot{\mathbf{x}} \in \mathbb{R}^{n_1}$ ,  $\dot{\mathbf{y}} \in \mathbb{R}^{n_2}$ :

$$\partial g(\mathbf{x}, \mathbf{y})[\dot{\mathbf{x}}, \dot{\mathbf{y}}] = \partial_1 g(\mathbf{x}, \mathbf{y})[\dot{\mathbf{x}}] + \partial_2 g(\mathbf{x}, \mathbf{y})[\dot{\mathbf{y}}]. \quad (6.17)$$

Alternatively, taking the matrix view:

$$\partial g(\mathbf{x}, \mathbf{y}) = (\partial_1 g(\mathbf{x}, \mathbf{y}) \quad \partial_2 g(\mathbf{x}, \mathbf{y})). \quad (6.18)$$

This convention will simplify our chain rule statement below. When  $n_1 = n_2 = m = 1$ , both sub-matrices are scalar, and  $\partial g_1(x, y)$  recovers the partial derivative that might otherwise be written in named variable notation as:

$$\frac{\partial}{\partial x} g(x, y). \quad (6.19)$$

However, the expression  $\partial g_1$  bears a meaning on its own (as a function) whereas the expression  $\frac{\partial g}{\partial x}$  may be ambiguous without further context. Again composing operators lets us write higher-order derivatives. For instance,  $\partial_2 \partial_1 g(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{m \times n_1 \times n_2}$ , and if  $m = 1$ , the Hessian of  $g$  at  $(\mathbf{x}, \mathbf{y})$  is:

$$\begin{pmatrix} \partial_1 \partial_1 g(\mathbf{x}, \mathbf{y}) & \partial_1 \partial_2 g(\mathbf{x}, \mathbf{y}) \\ \partial_2 \partial_1 g(\mathbf{x}, \mathbf{y}) & \partial_2 \partial_2 g(\mathbf{x}, \mathbf{y}) \end{pmatrix}. \quad (6.20)$$

**Composition and fan-out.** If  $f = g \circ h$  for some  $h : \mathbb{R}^n \rightarrow \mathbb{R}^p$  and  $g : \mathbb{R}^p \rightarrow \mathbb{R}^m$ , then the **chain rule** of calculus observes that:

$$\partial f(\mathbf{x}) = \partial g(h(\mathbf{x})) \circ \partial h(\mathbf{x}) \text{ for all } \mathbf{x} \in \mathbb{R}^n \quad (6.21)$$

How does this interact with our notation for multi-argument functions? For one, it can lead us to consider expressions with **fan-out**, where several sub-expressions are functions of the same input. For instance, assume two functions  $a : \mathbb{R}^n \rightarrow \mathbb{R}^{m_1}$  and  $b : \mathbb{R}^n \rightarrow \mathbb{R}^{m_2}$ , and that:

$$f(\mathbf{x}) = g(a(\mathbf{x}), b(\mathbf{x})) \quad (6.22)$$

for some function  $g$ . Abbreviating  $h(\mathbf{x}) = (a(\mathbf{x}), b(\mathbf{x}))$  so that  $f(\mathbf{x}) = g(h(\mathbf{x}))$ , Equations (6.16) and (6.21) tell us that:

$$\partial f(\mathbf{x}) = \partial g(h(\mathbf{x})) \circ \partial h(\mathbf{x}) \quad (6.23)$$

$$= \partial_1 g(a(\mathbf{x}), b(\mathbf{x})) \circ \partial a(\mathbf{x}) + \partial_2 g(a(\mathbf{x}), b(\mathbf{x})) \circ \partial b(\mathbf{x}) \quad (6.24)$$

Note that  $+$  is meant pointwise here. It also follows from the above that if instead:

$$f(\mathbf{x}, \mathbf{y}) = g(a(\mathbf{x}), b(\mathbf{y})) \quad (6.25)$$

in other words, if we write multiple arguments but exhibit no fan-out, then:

$$\partial_1 f(\mathbf{x}, \mathbf{y}) = \partial_1 g(a(\mathbf{x}), b(\mathbf{y})) \circ \partial a(\mathbf{x}), \text{ and} \quad (6.26)$$

$$\partial_2 f(\mathbf{x}, \mathbf{y}) = \partial_2 g(a(\mathbf{x}), b(\mathbf{y})) \circ \partial b(\mathbf{y}) \quad (6.27)$$

Composition and fan-out rules for derivatives are what let us break down a complex derivative calculation into simpler ones. This is what automatic differentiation techniques rely on when processing the sort of elaborate numerical computations that turn up in modern machine learning and numerical programming.

1 2 **6.2.2 Differentiating chains, circuits, and programs**

3 The purpose of automatic differentiation is to compute derivatives of arbitrary functions provided as  
4 input. Given a function  $f : U \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$  and a linearization point  $\mathbf{x} \in U$ , AD computes either:  
5

- 6 • the JVP  $\partial f(\mathbf{x})[\mathbf{v}]$  for an input perturbation  $\mathbf{v} \in \mathbb{R}^n$ , or  
7
- 8 • the VJP  $\partial f(\mathbf{x})^\top[\mathbf{u}]$  for an output perturbation  $\mathbf{u} \in \mathbb{R}^m$ .

9 In other words, JVPs and VJPs capture the two essential tasks of AD.<sup>2</sup>

10 Deciding what functions  $f$  to handle as input, and how to represent them, is perhaps the most  
11 load-bearing aspect of this setup. Over what *language* of functions should we operate? By a  
12 language, we mean some formal way of describing functions by composing a set of basic primitive  
13 operations. For primitives, we can think of various differentiable array operations (elementwise  
14 arithmetic, reductions, contractions, indexing and slicing, concatenation, etc.), but we will largely  
15 consider primitives and their derivatives as a given, and focus on how elaborately we can compose  
16 them. AD becomes increasingly challenging with increasingly expressive languages. Considering this,  
17 we introduce it in stages.

18

19 **6.2.2.1 Chain compositions and the chain rule**

20 To start, take only functions that are **chain compositions** of basic operations. Chains are a  
21 convenient class of function representations because derivatives *decompose* along the same structure  
22 according to the aptly-named chain rule.  
23

24 As a toy example, consider  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  composed of three operations in sequence:

$$\underline{25} \quad f = c \circ b \circ a \quad (6.28) \quad \underline{26}$$

27 By the chain rule, its derivatives are given by

$$\underline{28} \quad \partial f(\mathbf{x}) = \partial c(b(a(\mathbf{x}))) \circ \partial b(a(\mathbf{x})) \circ \partial a(\mathbf{x}) \quad (6.29) \quad \underline{29}$$

30 Now consider the JVP against an input perturbation  $\mathbf{v} \in \mathbb{R}^n$ :

$$\underline{31} \quad \partial f(\mathbf{x})[\mathbf{v}] = \partial c(b(a(\mathbf{x}))) [\partial b(a(\mathbf{x})) [\partial a(\mathbf{x})[\mathbf{v}]]] \quad (6.30) \quad \underline{32}$$

33 This expression's bracketing highlights a right-to-left evaluation order that corresponds to **forward-**  
34 **mode automatic differentiation**. Namely, to carry out this JVP, it makes sense to compute  
35 prefixes of the original chain:  
36

$$\underline{37} \quad \mathbf{x}, a(\mathbf{x}), b(a(\mathbf{x})) \quad (6.31) \quad \underline{38}$$

39 alongside the partial JVPs, because each is then immediately used as a subsequent linearization  
40 point, respectively:

$$\underline{41} \quad \partial a(\mathbf{x}), \partial b(a(\mathbf{x})), \partial c(b(a(\mathbf{x}))) \quad (6.32) \quad \underline{42}$$

43 Extending this idea to arbitrary chain compositions gives Algorithm 1.

44 

---

45 2. Materializing the Jacobian as a numerical array, as is commonly required in an optimization context, is a special  
46 case of computing a JVP or VJP against the standard basis vectors in  $\mathbb{R}^n$  or  $\mathbb{R}^m$  respectively.

47

---

**Algorithm 1:** Forward-mode automatic differentiation (JVP) on chains

---

```

1 input:  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  as a chain composition  $f = f_T \circ \dots \circ f_1$ 
2 input: linearization point  $\mathbf{x} \in \mathbb{R}^n$  and input perturbation  $\mathbf{v} \in \mathbb{R}^n$ 
3  $\mathbf{x}_0, \mathbf{v}_0 := \mathbf{x}, \mathbf{v}$ 
4 for  $t := 1, \dots, T$  do
5    $\mathbf{x}_t := f_t(\mathbf{x}_{t-1})$ 
6    $\mathbf{v}_t := \partial f_t(\mathbf{x}_{t-1})[\mathbf{v}_{t-1}]$ 
7 output:  $\mathbf{x}_T$ , equal to  $f(\mathbf{x})$ 
8 output:  $\mathbf{v}_T$ , equal to  $\partial f(\mathbf{x})[\mathbf{v}]$ 

```

---

By contrast, we can transpose Equation (6.29) to consider a VJP against an output perturbation  $\mathbf{u} \in \mathbb{R}^m$ :

$$\partial f(\mathbf{x})^\top[\mathbf{u}] = \partial a(\mathbf{x})^\top [\partial b(a(\mathbf{x}))^\top [\partial c(b(a(\mathbf{x})))^\top[\mathbf{u}]]] \quad (6.33)$$

Transposition reverses the Jacobian maps relative to their order in Equation (6.29), and now the bracketed evaluation corresponds to **reverse-mode automatic differentiation**. To carry out this VJP, we can compute the original chain prefixes  $\mathbf{x}$ ,  $a(\mathbf{x})$ , and  $b(a(\mathbf{x}))$  first, and then read them *in reverse* as successive linearization points:

$$\partial c(b(a(\mathbf{x})))^\top, \partial b(a(\mathbf{x}))^\top, \partial a(\mathbf{x})^\top \quad (6.34)$$

Extending this idea to arbitrary chain compositions gives Algorithm 2.

---

**Algorithm 2:** Reverse-mode automatic differentiation (VJP) on chains

---

```

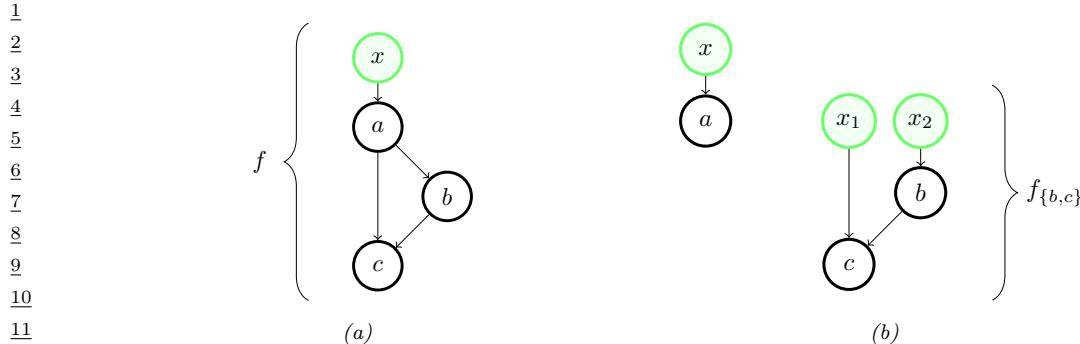
1 input:  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  as a chain composition  $f = f_T \circ \dots \circ f_1$ 
2 input: linearization point  $\mathbf{x} \in \mathbb{R}^n$  and output perturbation  $\mathbf{u} \in \mathbb{R}^n$ 
3  $\mathbf{x}_0 := \mathbf{x}$ 
4 for  $t := 1, \dots, T$  do
5    $\mathbf{x}_t := f_t(\mathbf{x}_{t-1})$ 
6    $\mathbf{u}_{T+1} := \mathbf{u}$ 
7   for  $t := T, \dots, 1$  do
8      $\mathbf{u}_t := \partial f_t(\mathbf{x}_{t-1})^\top[\mathbf{u}_{t+1}]$ 
9   output:  $\mathbf{x}_T$ , equal to  $f(\mathbf{x})$ 
10  output:  $\mathbf{u}_1$ , equal to  $\partial f(\mathbf{x})^\top[\mathbf{u}]$ 

```

---

Although chain compositions impose a very specific structure, they already capture some deep neural network models, such as multi-layer perceptrons (provided matrix multiplication is a primitive operation), as covered in this book's prequel [Mur22, Ch.13].

Reverse-mode AD is faster than forward-mode when the output is scalar valued (as often arises in deep learning, where the output is a loss function). However, reverse-mode AD stores all chain



17 prefixes before its backward traversal, so it consumes more memory than forward-mode. There  
18 are ways to combat this memory requirement in special-case scenarios, such as when the chained  
19 operations are each reversible [MDA15; Gom+17; KKL20]. One can also trade off memory for  
20 computation by discarding some prefixes and re-computing them as needed.

21

### 6.2.2.2 From chains to circuits

24 When primitives can accept multiple inputs, we can naturally extend chains to **circuits**—directed  
25 acyclic graphs over primitive operations, sometimes also called computation graphs. To set up for  
26 this section, we will distinguish between (i) **input nodes** of a circuit, which symbolize a function’s  
27 arguments, and (ii) **primitive nodes**, each of which is labeled by a primitive operation. We assume  
28 that input nodes have no incoming edges and (without loss of generality) exactly one outgoing edge  
29 each, and that the graph has exactly one sink node. The overall function of the circuit is composition  
30 of operations from the input nodes to the sink, where the output of each operation is input to others  
31 according to its outgoing edges.

32 What made AD work in Section 6.2.2.1 is the fact that derivatives decompose along chains thanks  
33 to the aptly-named chain rule. When moving from chains to directed acyclic graphs, do we need  
34 some sort of “graph rule” in order to decompose our calculation along the circuit’s structure? Circuits  
35 introduce two new features: **fan-in** and **fan-out**. In graphical terms, fan-in simply refers to multiple  
36 edges incoming to a node, and fan-out refers to multiple edges outgoing.

37 What do these mean in functional terms? Fan-in happens when a primitive operation accepts  
38 multiple arguments. We observed in Section 6.2.1 that multiple arguments can be treated as one, and  
39 how the chain rule then applies. Fan-out requires slightly more care, specifically for reverse-mode  
40 differentiation.

41 The gist of an answer can be illustrated with a small example. Consider the circuit in Figure 6.1a.  
42 The operation  $a$  precedes  $b$  and  $c$  topologically, with an outgoing edge to each of both. We can cut  $a$   
43 away from  $\{b, c\}$  to produce two new circuits, shown in Figure 6.1b. The first corresponds to  $a$  and  
44 the second corresponds to the remaining computation, given by:

$$45 \\ 46 \quad f_{\{b,c\}}(\mathbf{x}_1, \mathbf{x}_2) = c(\mathbf{x}_1, b(\mathbf{x}_2)). \quad (6.35)$$

47

We can recover the complete function  $f$  from  $a$  and  $f_{\{b,c\}}$  with the help of a function  $\text{dup}$  given by:

$$\text{dup}(\mathbf{x}) = (\mathbf{x}, \mathbf{x}) \equiv \begin{pmatrix} I \\ I \end{pmatrix} \mathbf{x} \quad (6.36)$$

so that  $f$  can be written as a chain composition:

$$f = f_{\{b,c\}} \circ \text{dup} \circ a. \quad (6.37)$$

The circuit for  $f_{\{b,c\}}$  contains no fan-out, and composition rules such as Equation (6.25) tell us its derivatives in terms of  $b$ ,  $c$ , and their derivatives, all via the chain rule. Meanwhile, the chain rule applied to Equation (6.37) says that:

$$\partial f(\mathbf{x}) = \partial f_{\{b,c\}}(\text{dup}(a(\mathbf{x}))) \circ \partial \text{dup}(a(\mathbf{x})) \circ \partial a(\mathbf{x}) \quad (6.38)$$

$$= \partial f_{\{b,c\}}(a(\mathbf{x}), a(\mathbf{x})) \circ \begin{pmatrix} I \\ I \end{pmatrix} \circ \partial a(\mathbf{x}). \quad (6.39)$$

The above expression suggests calculating a JVP of  $f$  by right-to-left evaluation. It is similar to the JVP calculation suggested by Equation (6.30), but with a *duplication* operation  $(I \ I)^\top$  in the middle that arises from the Jacobian of  $\text{dup}$ .

Transposing the derivative of  $f$  at  $\mathbf{x}$ :

$$\partial f(\mathbf{x})^\top = \partial a(\mathbf{x})^\top \circ (I \ I) \circ \partial f_{\{b,c\}}(a(\mathbf{x}), a(\mathbf{x}))^\top. \quad (6.40)$$

Considering right-to-left evaluation, this too is similar to the VJP calculation suggested by Equation (6.33), but with a *summation* operation  $(I \ I)$  in the middle that arises from the *transposed* Jacobian of  $\text{dup}$ . The lesson of using  $\text{dup}$  in this small example is that, more generally, in order to handle fan-out in reverse mode AD, we can process operations in topological order—first forward and then in reverse—and then *sum* partial VJPs along multiple outgoing edges.

### Algorithm 3: Foward-mode circuit differentiation (JVP)

```

1 input:  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  composing  $f_1, \dots, f_T$  in topological order, where  $f_1$  is identity
2 input: linearization point  $\mathbf{x} \in \mathbb{R}^n$  and perturbation  $\mathbf{v} \in \mathbb{R}^n$ 
3  $\mathbf{x}_1, \mathbf{v}_1 := \mathbf{x}, \mathbf{v}$ 
4 for  $t := 2, \dots, T$  do
5   let  $[q_1, \dots, q_r] = \text{Pa}(t)$ 
6    $\mathbf{x}_t := f_t(\mathbf{x}_{q_1}, \dots, \mathbf{x}_{q_r})$ 
7    $\mathbf{v}_t := \sum_{i=1}^r \partial_i f_t(\mathbf{x}_{q_1}, \dots, \mathbf{x}_{q_r})[\mathbf{v}_{q_i}]$ 
8 output:  $\mathbf{x}_T$ , equal to  $f(\mathbf{x})$ 
9 output:  $\mathbf{v}_T$ , equal to  $\partial f(\mathbf{x})[\mathbf{v}]$ 
```

Algorithms 3 and 4 give a complete description of forward- and reverse-mode differentiation on circuits. For brevity they assume a single argument to the entire circuit function. Nodes are indexed  $1, \dots, T$ . The first is the input node associated and the remaining  $T - 1$  are labeled by their operation  $f_2, \dots, f_T$ . We take  $f_1$  to be the identity. For each  $t$ , if  $f_t$  takes  $k$  arguments, let  $\text{Pa}(t)$  be the ordered

---

1  
2   **Algorithm 4:** Reverse-mode circuit differentiation (VJP)  
3   **1 input:**  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  composing  $f_1, \dots, f_T$  in topological order, where  $f_1, f_T$  are identity  
4   **2 input:** linearization point  $\mathbf{x} \in \mathbb{R}^n$  and perturbation  $\mathbf{u} \in \mathbb{R}^m$   
5   **3**  $\mathbf{x}_1 := \mathbf{x}$   
6   **4 for**  $t := 2, \dots, T$  **do**  
7    **5**   let  $[q_1, \dots, q_r] = \text{Pa}(t)$   
8    **6**    $\mathbf{x}_t := f_t(\mathbf{x}_{q_1}, \dots, \mathbf{x}_{q_r})$   
9  
10   **7**  $\mathbf{u}_{(T-1) \rightarrow T} := \mathbf{u}$   
11   **8 for**  $t := T - 1, \dots, 2$  **do**  
12    **9**   let  $[q_1, \dots, q_r] = \text{Pa}(t)$   
13    **10**    $\mathbf{u}'_t := \sum_{c \in \text{Ch}(t)} \mathbf{u}_{t \rightarrow c}$   
14    **11**    $\mathbf{u}_{q_i \rightarrow t} := \partial_i f_t(\mathbf{x}_{q_1}, \dots, \mathbf{x}_{q_r})^\top \mathbf{u}'_t$  for  $i = 1, \dots, r$   
15   **12 output:**  $\mathbf{x}_T$ , equal to  $f(\mathbf{x})$   
16   **13 output:**  $\mathbf{u}_{1 \rightarrow 2}$ , equal to  $\partial f(\mathbf{x})^\top \mathbf{u}$   
17

---

18  
19

20 list of  $k$  indices of its parent nodes (possibly containing duplicates, due to fan-out), and let  $\text{Ch}(t)$  be  
21 the indices of its children (again possibly duplicate). Algorithm 4 takes a few more conventions: that  
22  $f_T$  is the identity, that node  $T$  has  $T - 1$  as its only parent, and that the child of node 1 is node 2.  
23 Fan-out is a feature of *graphs*, but arguably not an essential feature of *functions*. One can always  
24 remove all fan-out from a circuit representation by duplicating nodes. Our interest in fan-out is  
25 precisely to avoid this, allowing for an efficient representation and, in turn, efficient memory use in  
26 Algorithms 3 and 4.

27 Reverse-mode AD on circuits has appeared under various names and formulations over the years.  
28 The algorithm is precisely the **backpropagation** algorithm in neural networks, a term introduced  
29 in the 1980s [RHW86b; RHW86a], and has separately come up in the context of control theory  
30 and sensitivity, as summarized in historical notes by Goodfellow, Bengio, and Courville [GBC16,  
31 Section 6.6].

32

### 33 6.2.2.3 From circuits to programs

34 Graphs are useful for introducing AD algorithms, and they might align well enough with neural  
35 network applications. But computer scientists have spent decades formalizing and studying various  
36 “languages for expressing functions compositionally.” Simply put, this is what programming languages  
37 are for! Can we automatically differentiate numerical functions expressed in, say, Python, Haskell,  
38 or some variant of the lambda calculus? These offer a far more widespread—and intuitively more  
39 expressive—way to describe an input function.<sup>3</sup>

40 In the previous sections, our approach to AD became more complex as we allowed for more  
41 complex graph structure. Something similar happens when we introduce grammatical constructs in a  
42 programming language. How do we adapt AD to handle a language with loops, conditionals, and

43  
44 3. In Python, what the language calls a “function” does not always describe a pure function of the arguments listed  
45 in its syntactic definition; its behavior may rely on side effects or global state, as allowed by the language. Here, we  
46 specifically mean a Python function that is pure and functional. JAX’s documentation details this restriction [Bra+18].

47

recursive calls? What about parallel programming constructs? We have partial answers to questions like these today, although they invite a deeper dive into language details such as type systems and implementation concerns [Yu+18; Inn20; Pas+21b].

One example language construct that we already know how to handle, due to Section 6.2.2.2, is a standard `let` expression. In languages with a means of name or variable binding, multiple appearances of the same variable are analogous to fan-out in a circuit. Figure 6.1a corresponds to a function  $f$  that we could write in a functional language as:

```
9   f(x) =
10  let ax = a(x)
11    in c(ax, b(ax))
```

in which `ax` indeed appears twice after it is bound.

Understanding the interaction between language capacity and automatic differentiability is an ongoing topic of computer science research [PS08a; AP19; Vyt+19; BMP19; MP21]. In the meantime, functional languages have proven quite effective in recent AD systems, both widely-used and experimental. Systems such as JAX, Dex, and others are designed around pure functional programming models, and internally rely on functional program representations for differentiation [Mac+15; BPS16; Sha+19; FJL18; Bra+18; Mac+19; Dex; Fro+21; Pas+21a].

## 6.3 Stochastic gradient descent

In this section, we consider optimizers for unconstrained differentiable objectives. We consider gradient-based solvers which perform iterative updates of the following form

$$\theta_{t+1} = \theta_t - \eta_t \mathbf{C}_t \mathbf{g}_t \quad (6.41)$$

where  $\mathbf{g}_t = \nabla \mathcal{L}(\theta_t)$  is the gradient of the loss, and  $\mathbf{C}_t$  is an optional **conditioning matrix**.

If we set  $\mathbf{C}_t = \mathbf{I}$ , the method is known as **steepest descent** or **gradient descent**. If we set  $\mathbf{C}_t = \mathbf{H}_t^{-1}$ , where  $\mathbf{H}_t = \nabla^2 \mathcal{L}(\theta_t)$  is the Hessian, we get **Newton's method**. There are many variants of Newton's method that are either more numerically stable, or more computationally efficient, or both.

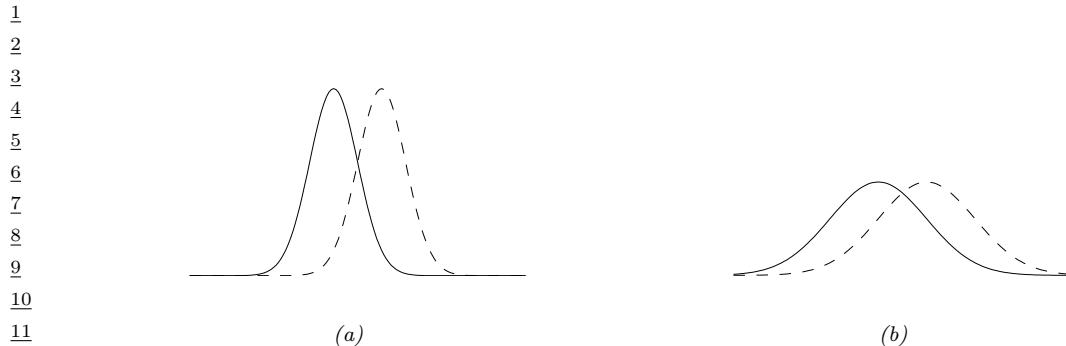
In many problems, we cannot compute the exact gradient, either because the loss is stochastic (e.g., due to random factors in the environment), or because we approximate the loss by randomly subsampling the data. In such cases, we can modify the update in Equation (6.41) to use an unbiased approximation of the gradient. For example, suppose the loss is a finite-sum objective from a supervised learning problem:

$$\mathcal{L}(\theta_t) = \frac{1}{N} \sum_{n=1}^N \ell(y_n, f(x_n; \theta_t)) = \frac{1}{N} \sum_{n=1}^N \mathcal{L}_n(\theta_t) \quad (6.42)$$

We can approximate the gradient using a minibatch  $\mathcal{B}_t$  of size  $B = |\mathcal{B}_t|$  as follows:

$$\hat{\mathbf{g}}_t = \hat{\nabla} \mathcal{L}(\theta_t) = \frac{1}{B} \sum_{n \in \mathcal{B}_t} \nabla \mathcal{L}_n(\theta_t) \quad (6.43)$$

Since the minibatches are randomly sampled, this is a stochastic, but unbiased, estimate of the gradient. If we insert  $\hat{\mathbf{g}}_t$  into Equation (6.41), the method is called **stochastic gradient descent** or **SGD**.



*Figure 6.2: Changing the mean of a Gaussian by a fixed amount (from solid to dotted curve) can have more impact when the (shared) variance is small (as in a) compared to when the variance is large (as in b). Hence the impact (in terms of prediction accuracy) of a change to  $\mu$  depends on where the optimizer is in  $(\mu, \sigma)$  space. From Figure 3 of [Hon+10], reproduced from [Val00]. Used with kind permission of Antti Honkela.*

## 6.4 Natural gradient descent

In this section, we discuss **natural gradient descent (NGD)** [Ama98], which is a second order method for optimizing the parameters of (conditional) probability distributions  $p_{\theta}(y|x)$ . The key idea is to compute parameter updates by measuring distances between the induced distributions, rather than comparing parameter values directly.

For example, consider comparing two Gaussians,  $p_{\theta} = p(y|\mu, \sigma)$  and  $p_{\theta'} = p(y|\mu', \sigma')$ . The (squared) Euclidean distance between the parameter vectors decomposes as  $\|\theta - \theta'\|^2 = (\mu - \mu')^2 + (\sigma - \sigma')^2$ . However, the predictive distribution has the form  $\exp(-\frac{1}{2\sigma^2}(y - \mu)^2)$ , so changes in  $\mu$  need to be measured relative to  $\sigma$ . This is illustrated in Figure 6.2(a-b), which shows two univariate Gaussian distributions (dotted and solid lines) whose means differ by  $\delta$ . In Figure 6.2(a), they share the same small variance  $\sigma^2$ , whereas in Figure 6.2(b), they share the same large variance. It is clear that the value of  $\delta$  matters much more (in terms of the effect on the distribution) when the variance is small. Thus we see that the two parameters interact with each other, which the Euclidean distance cannot capture. This problem gets much worse when we consider more complex models, such as deep neural networks. By modeling such correlations, NGD can converge much faster than other gradient methods.

### 6.4.1 Defining the natural gradient

The key to NGD is to measure the notion of distance between two probability distributions in terms of the KL divergence. As we show in Section 2.4.4, this can be approximated in terms of the Fisher information matrix (FIM). In particular, for any given input  $x$ , we have

$$D_{\text{KL}}(p_{\theta}(y|x) \parallel p_{\theta+\delta}(y|x)) \approx \frac{1}{2} \delta^T \mathbf{F}_x \delta \quad (6.44)$$

where  $\mathbf{F}_x$  is the FIM

$$\mathbf{F}_x(\theta) = -\mathbb{E}_{p_{\theta}(y|x)} [\nabla^2 \log p_{\theta}(y|x)] = \mathbb{E}_{p_{\theta}(y|x)} [(\nabla \log p_{\theta}(y|x))(\nabla \log p_{\theta}(y|x))^T] \quad (6.45)$$

We can compute the average KL between the current and updated distributions using  $\frac{1}{2}\boldsymbol{\delta}^T \mathbf{F} \boldsymbol{\delta}$ , where  $\mathbf{F}$  is the averaged FIM:

$$\mathbf{F}(\boldsymbol{\theta}) = \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [\mathbf{F}_{\mathbf{x}}(\boldsymbol{\theta})] \quad (6.46)$$

NGD uses the inverse FIM as a preconditioning matrix, i.e., we perform updates of the following form:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_t \mathbf{F}(\boldsymbol{\theta}_t)^{-1} \mathbf{g}_t \quad (6.47)$$

The term

$$\mathbf{F}^{-1} \mathbf{g}_t = \mathbf{F}^{-1} \nabla \mathcal{L}(\boldsymbol{\theta}_t) \triangleq \tilde{\nabla} \mathcal{L}(\boldsymbol{\theta}_t) \quad (6.48)$$

is called the **natural gradient**.

## 6.4.2 Interpretations of NGD

### 6.4.2.1 NGD as a trust region method

In Supplementary Section 6.1.3.1 we show that we can interpret standard gradient descent as optimizing a linear approximation to the objective subject to a penalty on the  $\ell_2$  norm of the change in parameters, i.e., if  $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \boldsymbol{\delta}$ , then we optimize

$$M_t(\boldsymbol{\delta}) = \mathcal{L}(\boldsymbol{\theta}_t) + \mathbf{g}_t^T \boldsymbol{\delta} + \eta \|\boldsymbol{\delta}\|_2^2 \quad (6.49)$$

Now let us replace the squared distance with the squared FIM-based distance,  $\|\boldsymbol{\delta}\|_F^2 = \boldsymbol{\delta}^T \mathbf{F} \boldsymbol{\delta}$ . This is equivalent to squared Euclidean distance in the **whitened coordinate system**  $\boldsymbol{\phi} = \mathbf{F}^{\frac{1}{2}} \boldsymbol{\theta}$ , since

$$\|\boldsymbol{\phi}_{t+1} - \boldsymbol{\phi}_t\|_2^2 = \|\mathbf{F}^{\frac{1}{2}}(\boldsymbol{\theta}_t + \boldsymbol{\delta}) - \mathbf{F}^{\frac{1}{2}} \boldsymbol{\theta}_t\|_2^2 = \|\mathbf{F}^{\frac{1}{2}} \boldsymbol{\delta}\|_2^2 = \|\boldsymbol{\delta}\|_F^2 \quad (6.50)$$

The new objective becomes

$$M_t(\boldsymbol{\delta}) = \mathcal{L}(\boldsymbol{\theta}_t) + \mathbf{g}_t^T \boldsymbol{\delta} + \eta \boldsymbol{\delta}^T \mathbf{F} \boldsymbol{\delta} \quad (6.51)$$

Solving  $\nabla_{\boldsymbol{\delta}} M_t(\boldsymbol{\delta}) = \mathbf{0}$  gives the update

$$\boldsymbol{\delta}_t = -\eta \mathbf{F}^{-1} \mathbf{g}_t \quad (6.52)$$

This is the same as the natural gradient direction. Thus we can view NGD as a trust region method, where we use a first-order approximation to the objective, and use FIM-distance in the constraint.

In the above derivation, we assumed  $\mathbf{F}$  was a constant matrix. In most problems, it will change at each point in space, since we are optimizing in a curved space known as a **Riemannian manifold**. For certain models, we can compute the FIM efficiently, allowing us to capture curvature information, even though we use a first-order approximation to the objective.

### 6.4.2.2 NGD as a Gauss-Newton method

If  $p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})$  is an exponential family distribution with natural parameters computed by  $\boldsymbol{\eta} = f(\mathbf{x}, \boldsymbol{\theta})$ , then one can show [Hes00; PB14] that NGD is identical to the generalized Gauss-Newton (GGN) method (Section 17.3.2). Furthermore, in the online setting, these methods are equivalent to performing sequential Bayesian inference using the extended Kalman filter, as shown in [Oll18].

---

1 **6.4.3 Benefits of NGD**

3 The use of the FIM as a preconditioning matrix, rather than the Hessian, has two advantages. First,  
4  $\mathbf{F}$  is always positive definite, whereas  $\mathbf{H}$  can have negative eigenvalues at saddle points, which are  
5 prevalent in high dimensional spaces. Second, it is easy to approximate  $\mathbf{F}$  online from minibatches,  
6 since it is an expectation (wrt the empirical distribution) of outer products of gradient vectors. This  
7 is in contrast to Hessian-based methods [Byr+16; Liu+18a], which are much more sensitive to noise  
8 introduced by the minibatch approximation.

9 In addition, the connection with trust region optimization makes it clear that NGD updates  
10 parameters in a way that matter most for prediction, which allows the method to take larger steps in  
11 uninformative regions of parameter space, which can help avoid getting stuck on plateaus. This can  
12 also help with issues that arise when the parameters are highly correlated.

13 For example, consider a 2d Gaussian with an unusual, highly coupled parameterization, proposed  
14 in [SD12]:

$$\frac{16}{17} p(\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{2\pi} \exp \left[ -\frac{1}{2} \left( (x_1 - \left[ 3\theta_1 + \frac{1}{3}\theta_2 \right])^2 - \frac{1}{2} \left( x_2 - \left[ \frac{1}{3}\theta_1 \right] \right)^2 \right) \right] \quad (6.53)$$

19 The objective is the cross entropy loss:  
20

$$\frac{21}{22} \mathcal{L}(\boldsymbol{\theta}) = -\mathbb{E}_{p^*(\mathbf{x})} [\log p(\mathbf{x}; \boldsymbol{\theta})] \quad (6.54)$$

23 The gradient of this objective is given by

$$\frac{25}{26} \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) \begin{pmatrix} = \mathbb{E}_{p^*(\mathbf{x})} [3(x_1 - [3\theta_1 + \frac{1}{3}\theta_2]) + \frac{1}{3}(x_2 - [\frac{1}{3}\theta_1])] \\ \mathbb{E}_{p^*(\mathbf{x})} [\frac{1}{3}(x_1 - [3\theta_1 + \frac{1}{3}\theta_2])] \end{pmatrix} \quad (6.55)$$

27 Suppose that  $p^*(\mathbf{x}) = p(\mathbf{x}; [0, 0])$ . Then the Fisher matrix is a constant matrix, given by

$$\frac{29}{30} \mathbf{F} = \begin{pmatrix} 3^2 + \frac{1}{3^2} & 1 \\ 1 & \frac{1}{3^2} \end{pmatrix} \quad (6.56)$$

32 Figure 6.3 compares steepest descent in  $\boldsymbol{\theta}$  space with the natural gradient method, which is  
33 equivalent to steepest descent in  $\phi$  space. Both methods start at  $\boldsymbol{\theta} = (1, -1)$ . The global optimum is  
34 at  $\boldsymbol{\theta} = (0, 0)$ . We see that the NG method (blue dots) converges much faster to this optimum and  
35 takes the shortest path, whereas steepest descent takes a very circuitous route. We also see that  
36 the gradient field in the whitened parameter space is more “spherical”, which makes descent much  
37 simpler and faster.

38 Finally, note that since NGD is invariant to how we parameterize the distribution, we will get the  
39 same results even for a standard parameterization of the Gaussian. This is particularly useful if our  
40 probability model is more complex, such as a DNN (see e.g., [SSE18]).

41

42 **6.4.4 Approximating the natural gradient**

44 The main drawback of NGD is the computational cost of computing (the inverse of) the Fisher  
45 Information Matrix (FIM). To speed this up, several methods make assumptions about the form  
46 of  $\mathbf{F}$ , so it can be inverted efficiently. For example, [LeC+98] uses a diagonal approximation for  
47

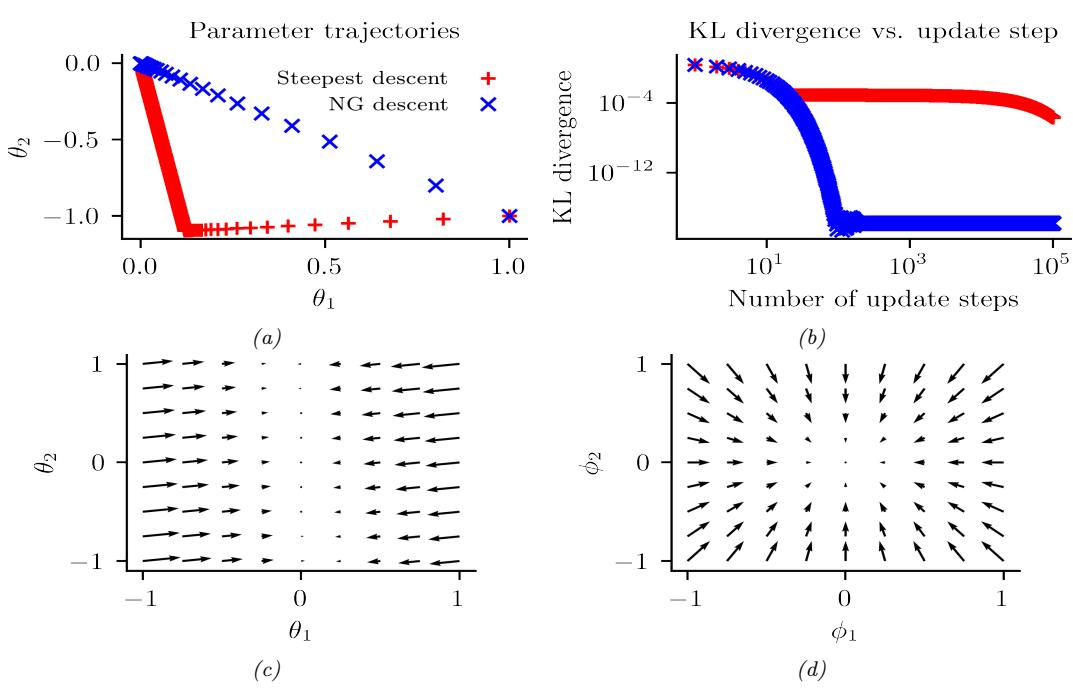


Figure 6.3: Illustration of the benefits of natural gradient vs steepest descent on a 2d problem. (a) Trajectories of the two methods in parameter space (red = steepest descent, blue = NG). They both start in the bottom right, at  $(1, -1)$ . (b) Objective vs number of iterations. (c) Gradient field in the  $\theta$  parameter space. (d) Gradient field in the whitened  $\phi = \mathbf{F}^{\frac{1}{2}}\theta$  parameter space used by NG. Generated by [nat\\_grad\\_demo.ipynb](#).

neural net training; [RMB08] uses a low-rank plus block diagonal approximation; and [GS15] assumes the covariance of the gradients can be modeled by a directed Gaussian graphical model with low treewidth (i.e., the Cholesky factorization of  $\mathbf{F}$  is sparse).

[MG15] propose the **KFAC** method, which stands for ‘‘Kronecker-Factored approximate curvature’’; this approximates the FIM of a DNN as a block diagonal matrix, where each block is a Kronecker product of two small matrices. This method has shown good results on supervised learning of neural nets [GM16; BGM17; Geo+18; Osa+19b] as well as reinforcement learning of neural policy networks [Wu+17]. The KFAC approximation can be justified using the mean field analysis of [AKO18]. In addition, [ZMG19] prove that KFAC will converge to the global optimum of a DNN if it is overparameterized (i.e., acts like an interpolator).

A simpler approach is to approximate the FIM by replacing the model’s distribution with the empirical distribution. In particular, define  $p_{\mathcal{D}}(\mathbf{x}, \mathbf{y}) = \frac{1}{N} \sum_{n=1}^N \delta_{\mathbf{x}_n}(\mathbf{x}) \delta_{\mathbf{y}_n}(\mathbf{y})$ ,  $p_{\mathcal{D}}(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \delta_{\mathbf{x}_n}(\mathbf{x})$  and  $p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{y}) = p_{\mathcal{D}}(\mathbf{x}) p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})$ . Then we can compute the **empirical Fisher Martens** [Mar16] as

1 follows:

$$\underline{3} \quad \mathbf{F} = \mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{y})} [\nabla \log p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) \nabla \log p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})^T] \quad (6.57)$$

$$\underline{4} \quad \approx \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x}, \mathbf{y})} [\nabla \log p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) \nabla \log p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})^T] \quad (6.58)$$

$$\underline{5} \quad = \frac{1}{|\mathcal{D}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \nabla \log p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) \nabla \log p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})^T \quad (6.59)$$

9 This approximation is widely used, since it is simple to compute. In particular, we can compute a  
10 diagonal approximation using the squared gradient vector. (This is similar to ADAGRAD, but only  
11 uses the current gradient instead of a moving average of gradients; the latter is a better approach  
12 when performing stochastic optimization.)

13 Unfortunately, the empirical Fisher does not work as well as the true Fisher [KBH19; Tho+19].  
14 To see why, note that when we reach a flat part of parameter space where the gradient vector goes  
15 to zero, the empirical Fisher will become singular, and hence the algorithm will get stuck on this  
16 plateau. However, the true Fisher takes expectations over the outputs, i.e., it marginalizes out  $\mathbf{y}$ .  
17 This will allow it to detect small changes in the output if we change the parameters. This is why the  
18 natural gradient method can “escape” plateaus better than standard gradient methods.

19 An alternative strategy is to use exact computation of  $\mathbf{F}$ , but solve for  $\mathbf{F}^{-1}\mathbf{g}$  approximately  
20 using truncated conjugate gradient (CG) methods, where each CG step uses efficient methods for  
21 Hessian-vector products [Pea94]. This is called **Hessian free optimization** [Mar10a]. However,  
22 this approach can be slow, since it may take many CG iterations to compute a single parameter  
23 update.

24

#### 25 6.4.5 Natural gradients for the exponential family

26 In this section, we assume  $\mathcal{L}$  is an expected loss of the following form:

$$\underline{28} \quad \mathcal{L}(\boldsymbol{\mu}) = \mathbb{E}_{q_{\boldsymbol{\mu}}(\boldsymbol{\theta})} [\ell(\boldsymbol{\theta})] \quad (6.60)$$

29 where  $q_{\boldsymbol{\mu}}(\boldsymbol{\theta})$  is an exponential family distribution with moment parameters  $\boldsymbol{\mu}$ . This is the basis of  
30 variational optimization (discussed in Section 6.7.3) and natural evolutionary strategies (discussed in  
31 Section 6.9.6).

32 It turns out the gradient wrt the moment parameters is the same as the natural gradient wrt the  
33 natural parameters  $\boldsymbol{\lambda}$ . This follows from the chain rule:

$$\underline{35} \quad \frac{d}{d\boldsymbol{\lambda}} \mathcal{L}(\boldsymbol{\lambda}) = \frac{d\boldsymbol{\mu}}{d\boldsymbol{\lambda}} \frac{d}{d\boldsymbol{\mu}} \mathcal{L}(\boldsymbol{\mu}) = \mathbf{F}(\boldsymbol{\lambda}) \nabla_{\boldsymbol{\mu}} \mathcal{L}(\boldsymbol{\mu}) \quad (6.61)$$

36 where  $\mathcal{L}(\boldsymbol{\mu}) = \mathcal{L}(\boldsymbol{\lambda}(\boldsymbol{\mu}))$ , and where we used Equation (2.215) to write

$$\underline{39} \quad \mathbf{F}(\boldsymbol{\lambda}) = \nabla_{\boldsymbol{\lambda}} \boldsymbol{\mu}(\boldsymbol{\lambda}) = \nabla_{\boldsymbol{\lambda}}^2 A(\boldsymbol{\lambda}) \quad (6.62)$$

40 Hence

$$\underline{42} \quad \tilde{\nabla}_{\boldsymbol{\lambda}} \mathcal{L}(\boldsymbol{\lambda}) = \mathbf{F}(\boldsymbol{\lambda})^{-1} \nabla_{\boldsymbol{\lambda}} \mathcal{L}(\boldsymbol{\lambda}) = \nabla_{\boldsymbol{\mu}} \mathcal{L}(\boldsymbol{\mu}) \quad (6.63)$$

44 It remains to compute the (regular) gradient wrt the moment parameters. The details on how to  
45 do this will depend on the form of the  $q$  and the form of  $\mathcal{L}(\boldsymbol{\lambda})$ . We discuss some approaches to this  
46 problem below.

47

---

### 6.4.5.1 Analytic computation for the Gaussian case

In this section, we assume that  $q(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta} | \mathbf{m}, \mathbf{V})$ . We now show how to compute the relevant gradients analytically.

Following Section 2.3.2.5, the natural parameters of  $q$  are

$$\boldsymbol{\lambda}^{(1)} = \mathbf{V}^{-1}\mathbf{m}, \quad \boldsymbol{\lambda}^{(2)} = -\frac{1}{2}\mathbf{V}^{-1} \quad (6.64)$$

and the moment parameters are

$$\boldsymbol{\mu}^{(1)} = \mathbf{m}, \quad \boldsymbol{\mu}^{(2)} = \mathbf{V} + \mathbf{m}\mathbf{m}^\top \quad (6.65)$$

For simplicity, we derive the result for the scalar case. Let  $m = \mu^{(1)}$  and  $v = \mu^{(2)} - (\mu^{(1)})^2$ . By using the chain rule, the gradient wrt the moment parameters are

$$\frac{\partial \mathcal{L}}{\partial \mu^{(1)}} = \frac{\partial \mathcal{L}}{\partial m} \frac{\partial m}{\partial \mu^{(1)}} + \frac{\partial \mathcal{L}}{\partial v} \frac{\partial v}{\partial \mu^{(1)}} = \frac{\partial \mathcal{L}}{\partial m} - 2 \frac{\partial \mathcal{L}}{\partial v} m \quad (6.66)$$

$$\frac{\partial \mathcal{L}}{\partial \mu^{(2)}} = \frac{\partial \mathcal{L}}{\partial m} \frac{\partial m}{\partial \mu^{(2)}} + \frac{\partial \mathcal{L}}{\partial v} \frac{\partial v}{\partial \mu^{(2)}} = \frac{\partial \mathcal{L}}{\partial v} \quad (6.67)$$

It remains to compute the derivatives wrt  $m$  and  $v$ . If  $\boldsymbol{\theta} \sim \mathcal{N}(\mathbf{m}, \mathbf{V})$ , then from **Bonnet's theorem** [Bon64] we have

$$\frac{\partial}{\partial m_i} \mathbb{E} [\ell(\boldsymbol{\theta})] = \mathbb{E} \left[ \frac{\partial}{\partial \theta_i} \ell(\boldsymbol{\theta}) \right] \quad (6.68)$$

And from **Price's theorem** [Pri58] we have

$$\frac{\partial}{\partial V_{ij}} \mathbb{E} [\ell(\boldsymbol{\theta})] = c_{ij} \mathbb{E} \left[ \frac{\partial^2}{\partial \theta_i \partial \theta_j} \ell(\boldsymbol{\theta}) \right] \quad (6.69)$$

where  $c_{ij} = \frac{1}{2}$  is  $i = j$  and  $c_{ij} = 1$  otherwise. (See [gradient\\_expected\\_value\\_gaussian.ipynb](#) for a “proof by example” of these claims.)

In the multivariate case, the result is as follows [OA09; KR21a]:

$$\nabla_{\boldsymbol{\mu}^{(1)}} \mathbb{E}_{q(\boldsymbol{\theta})} [\ell(\boldsymbol{\theta})] = \nabla_{\mathbf{m}} \mathbb{E}_{q(\boldsymbol{\theta})} [\ell(\boldsymbol{\theta})] - 2 \nabla_{\mathbf{V}} \mathbb{E}_{q(\boldsymbol{\theta})} [\ell(\boldsymbol{\theta})] \mathbf{m} \quad (6.70)$$

$$= \mathbb{E}_{q(\boldsymbol{\theta})} [\nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta})] - \mathbb{E}_{q(\boldsymbol{\theta})} [\nabla_{\boldsymbol{\theta}}^2 \ell(\boldsymbol{\theta})] \mathbf{m} \quad (6.71)$$

$$\nabla_{\boldsymbol{\mu}^{(2)}} \mathbb{E}_{q(\boldsymbol{\theta})} [\ell(\boldsymbol{\theta})] = \nabla_{\mathbf{V}} \mathbb{E}_{q(\boldsymbol{\theta})} [\ell(\boldsymbol{\theta})] \quad (6.72)$$

$$= \frac{1}{2} \mathbb{E}_{q(\boldsymbol{\theta})} [\nabla_{\boldsymbol{\theta}}^2 \ell(\boldsymbol{\theta})] \quad (6.73)$$

Thus we see that the natural gradients rely on both the gradient and Hessian of the loss function  $\ell(\boldsymbol{\theta})$ . We will see applications of this result in Section 6.7.2.2.

### 6.4.5.2 Stochastic approximation for the general case

In general, it can be hard to analytically compute the natural gradient. However, we can compute a Monte Carlo approximation. To see this, let us assume  $\mathcal{L}$  is an expected loss of the following form:

$$\mathcal{L}(\boldsymbol{\mu}) = \mathbb{E}_{q_{\boldsymbol{\mu}}(\boldsymbol{\theta})} [\ell(\boldsymbol{\theta})] \quad (6.74)$$

1 From Equation (6.63) the natural gradient is given by  
2

$$\underline{3} \quad \nabla_{\boldsymbol{\mu}} \mathcal{L}(\boldsymbol{\mu}) = \mathbf{F}(\boldsymbol{\lambda})^{-1} \nabla_{\boldsymbol{\lambda}} \mathcal{L}(\boldsymbol{\lambda}) \quad (6.75)$$

5 For exponential family distributions, both of these terms on the RHS can be written as expectations,  
6 and hence can be approximated by Monte Carlo, as noted by [KL17a]. To see this, note that  
7

$$\underline{8} \quad \mathbf{F}(\boldsymbol{\lambda}) = \nabla_{\boldsymbol{\lambda}} \boldsymbol{\mu}(\boldsymbol{\lambda}) = \nabla_{\boldsymbol{\lambda}} \mathbb{E}_{q_{\boldsymbol{\lambda}}(\boldsymbol{\theta})} [\mathcal{T}(\boldsymbol{\theta})] \quad (6.76)$$

$$\underline{9} \quad \nabla_{\boldsymbol{\lambda}} \mathcal{L}(\boldsymbol{\lambda}) = \nabla_{\boldsymbol{\lambda}} \mathbb{E}_{q_{\boldsymbol{\lambda}}(\boldsymbol{\theta})} [\ell(\boldsymbol{\theta})] \quad (6.77)$$

10 If  $q$  is reparameterizable, we can apply the reparameterization trick (Section 6.5.4) to push the  
11 gradient inside the expectation operator. This lets us sample  $\boldsymbol{\theta}$  from  $q$ , compute the gradients, and  
12 average; we can then pass the resulting stochastic gradients to SGD.  
13

#### 14 6.4.5.3 Natural gradient of the entropy function

16 In this section, we discuss how to compute the natural gradient of the entropy of an exponential  
17 family distribution, which is useful when performing variational inference (Chapter 10). The natural  
18 gradient is given by  
19

$$\underline{20} \quad \tilde{\nabla}_{\boldsymbol{\lambda}} \mathbb{H}(\boldsymbol{\lambda}) = -\nabla_{\boldsymbol{\mu}} \mathbb{E}_{q_{\boldsymbol{\mu}}(\boldsymbol{\theta})} [\log q(\boldsymbol{\theta})] \quad (6.78)$$

21 where, from Equation (2.143), we have  
22

$$\underline{23} \quad \log q(\boldsymbol{\theta}) = \log h(\boldsymbol{\theta}) + \mathcal{T}(\boldsymbol{\theta})^T \boldsymbol{\lambda} - A(\boldsymbol{\lambda}) \quad (6.79)$$

25 Since  $\mathbb{E}[\mathcal{T}(\boldsymbol{\theta})] = \boldsymbol{\mu}$ , we have  
26

$$\underline{27} \quad \nabla_{\boldsymbol{\mu}} \mathbb{E}_{q_{\boldsymbol{\mu}}(\boldsymbol{\theta})} [\log q(\boldsymbol{\theta})] = \nabla_{\boldsymbol{\mu}} \mathbb{E}_{q(\boldsymbol{\theta})} [\log h(\boldsymbol{\theta})] + \nabla_{\boldsymbol{\mu}} \boldsymbol{\mu}^T \boldsymbol{\lambda}(\boldsymbol{\mu}) - \nabla_{\boldsymbol{\mu}} A(\boldsymbol{\lambda}) \quad (6.80)$$

28 where  $h(\boldsymbol{\theta})$  is the base measure. Since  $\boldsymbol{\lambda}$  is a function of  $\boldsymbol{\mu}$ , we have  
29

$$\underline{30} \quad \nabla_{\boldsymbol{\mu}} \boldsymbol{\mu}^T \boldsymbol{\lambda} = \boldsymbol{\lambda} + (\nabla_{\boldsymbol{\mu}} \boldsymbol{\lambda})^T \boldsymbol{\mu} = \boldsymbol{\lambda} + (\mathbf{F}_{\boldsymbol{\lambda}}^{-1} \nabla_{\boldsymbol{\lambda}} \boldsymbol{\lambda})^T \boldsymbol{\mu} = \boldsymbol{\lambda} + \mathbf{F}_{\boldsymbol{\lambda}}^{-1} \boldsymbol{\mu} \quad (6.81)$$

31 and since  $\boldsymbol{\mu} = \nabla_{\boldsymbol{\lambda}} A(\boldsymbol{\lambda})$  we have  
32

$$\underline{33} \quad \nabla_{\boldsymbol{\mu}} A(\boldsymbol{\lambda}) = \mathbf{F}_{\boldsymbol{\lambda}}^{-1} \nabla_{\boldsymbol{\lambda}} A(\boldsymbol{\lambda}) = \mathbf{F}_{\boldsymbol{\lambda}}^{-1} \boldsymbol{\mu} \quad (6.82)$$

35 Hence  
36

$$\underline{37} \quad -\nabla_{\boldsymbol{\mu}} \mathbb{E}_{q_{\boldsymbol{\mu}}(\boldsymbol{\theta})} [\log q(\boldsymbol{\theta})] = -\nabla_{\boldsymbol{\mu}} \mathbb{E}_{q(\boldsymbol{\theta})} [\log h(\boldsymbol{\theta})] - \boldsymbol{\lambda} \quad (6.83)$$

38 If we assume that  $h(\boldsymbol{\theta}) = \text{const}$ , as is often the case, we get  
39

$$\underline{40} \quad \tilde{\nabla}_{\boldsymbol{\lambda}} \mathbb{H}(q) = -\boldsymbol{\lambda} \quad (6.84)$$

#### 42 6.5 Gradients of stochastic functions

44 In this section, we discuss how to compute the gradient of stochastic functions of the form  
45

$$\underline{46} \quad \mathcal{L}(\boldsymbol{\psi}) = \mathbb{E}_{q_{\boldsymbol{\psi}}(\boldsymbol{z})} [\ell(\boldsymbol{\psi}, \boldsymbol{z})] \quad (6.85)$$

1 **6.5.1 Minibatch approximation to finite-sum objectives**

2 In the simplest case,  $q_\psi(\mathbf{z})$  does not depend on  $\psi$ . In this case, we can push gradients inside the  
3 expectation operator,  $\nabla \mathcal{L}(\psi) = \mathbb{E}[\nabla \ell(\psi, \mathbf{z})]$  and then use Monte Carlo sampling for  $\mathbf{z}$  to approximate  
4 the gradient.  
5

6 For example, consider the empirical risk minimization (ERM) problem of minimizing  
7

$$\mathcal{L}(\psi) = \frac{1}{N_D} \sum_{n=1}^{N_D} \ell(\psi, \mathbf{z}_n) \quad (6.86)$$

8 where  $\mathbf{z}_n = (\mathbf{x}_n, \mathbf{y}_n)$  and  
9

$$\ell(\psi, \mathbf{z}_n) = \ell(h(\mathbf{x}_n; \psi), \mathbf{y}_n) \quad (6.87)$$

10 is the per-example loss, where  $h$  is a prediction function. This kind of objective is called a **finite  
11 sum objective**.

12 Now consider trying to minimize this objective. If, at each iteration, we evaluate the objective (and  
13 its gradient) using all  $N_D$  datapoints, the method is called **batch optimization**. However, this can  
14 be very slow if the dataset is large. Fortunately, we can reformulate it as a stochastic optimization  
15 problem, which will be faster to solve. To do this, note that Equation (6.86) can be written as an  
16 expectation wrt the empirical distribution:  
17

$$\mathcal{L}(\psi) = \mathbb{E}_{\mathbf{z} \sim p_D} [\ell(\psi, \mathbf{z})] \quad (6.88)$$

18 Since the distribution is independent of the parameters, we can easily use Monte Carlo sampling to  
19 approximate the objective and its gradient. In particular, we will sample a **minibatch** of  $B = |\mathcal{B}|$   
20 datapoints from the full set  $\mathcal{D}$  at each iteration. More precisely, we have  
21

$$\mathcal{L}(\psi) \approx \frac{1}{B} \sum_{n \in \mathcal{B}} \ell(h(\mathbf{x}_n; \psi), \mathbf{y}_n) \quad (6.89)$$

$$\nabla \mathcal{L}(\psi) \approx \frac{1}{B} \sum_{n \in \mathcal{B}} \nabla \ell(h(\mathbf{x}_n; \psi), \mathbf{y}_n) \quad (6.90)$$

22 These noisy gradients can then be passed to SGD, which is robust to noisy gradients (see Section 6.3).  
23

24 **6.5.2 Optimizing parameters of a distribution**

25 Now suppose the stochasticity depends on the parameters we are optimizing. For example,  $\mathbf{z}$  could  
26 be an action sampled from a stochastic policy  $q_\psi$ , as in RL (Section 35.3.2). In this case, the gradient  
27 is given by  
28

$$\nabla_\psi \mathbb{E}_{q_\psi(\mathbf{z})} [\ell(\psi, \mathbf{z})] = \nabla_\psi \int \ell(\psi, \mathbf{z}) q_\psi(\mathbf{z}) d\mathbf{z} \quad (6.91)$$

$$= \int [\nabla_\psi \ell(\psi, \mathbf{z})] q_\psi(\mathbf{z}) d\mathbf{z} + \int \ell(\psi, \mathbf{z}) [\nabla_\psi q_\psi(\mathbf{z})] d\mathbf{z} \quad (6.92)$$

The first term can be approximated by Monte Carlo sampling:

$$\int [\nabla_{\psi} \ell(\psi, z)] q_{\psi}(z) dz \approx \frac{1}{S} \sum_{s=1}^S \nabla_{\psi} \ell(\psi, z_s) \quad (6.93)$$

where  $z_s \sim q_{\psi}$ . Note that if  $\ell()$  is independent of  $\psi$ , this term vanishes.

Now consider the second term, that takes the gradients of the distribution itself:

$$I \triangleq \int \ell(\psi, z) [\nabla_{\psi} q_{\psi}(z)] dz \quad (6.94)$$

We can no longer use vanilla Monte Carlo sampling to approximate this integral. However, there are various other ways to approximate this (see [Moh+19] for an extensive review). We briefly describe the two main methods in Section 6.5.3 and Section 6.5.4.

### 6.5.3 Score function estimator (likelihood ratio trick)

The simplest way to approximate Equation (6.94) is to exploit the **log derivative trick**, which is the following identity:

$$\nabla_{\psi} q_{\psi}(z) = q_{\psi}(z) \nabla_{\psi} \log q_{\psi}(z) \quad (6.95)$$

With this, we can rewrite Equation (6.94) as follows:

$$I = \int \ell(\psi, z) [q_{\psi}(z) \nabla_{\psi} \log q_{\psi}(z)] dz = \mathbb{E}_{q_{\psi}(z)} [\ell(\psi, z) \nabla_{\psi} \log q_{\psi}(z)] \quad (6.96)$$

This is called the **score function estimator** or **SFE** [Fu15]. (The term “score function” refers to the gradient of a log probability distribution, as explained in Section 2.4.1.) It is also called the **likelihood ratio gradient estimator**. We can now easily approximate this with Monte Carlo:

$$I \approx \frac{1}{S} \sum_{s=1}^S \ell(\psi, z_s) \nabla_{\psi} \log q_{\psi}(z_s) \quad (6.97)$$

We only require that the sampling distribution is differentiable, not the objective  $\ell(\psi, z)$  itself. This allows the method to be used for blackbox stochastic optimization problems, such as variational optimization (Section 6.7.3), black-box variational inference (Section 10.3.2), reinforcement learning (Section 35.3.2), etc.

#### 6.5.3.1 Control variates

The score function estimate can have high variance. One way to reduce this is to use **control variates**, in which we replace  $\ell(\psi, z)$  with

$$\hat{\ell}(\psi, z) = \ell(\psi, z) - c(b(\psi, z) - \mathbb{E}[b(\psi, z)]) \quad (6.98)$$

where  $b(\psi, z)$  is a **baseline function** that is correlated with  $\ell(\psi, z)$ , and  $c > 0$  is a coefficient. Since  $\mathbb{E}[\hat{\ell}(\psi, z)] = \mathbb{E}[\ell(\psi, z)]$ , we can use  $\hat{\ell}$  to compute unbiased gradient estimates of  $\ell$ . The advantage is that this new estimate can result in lower variance, as we show in Section 11.6.3.

---

### 6.5.3.2 Rao-Blackwellisation

Suppose  $q_\psi(\mathbf{z})$  is a discrete distribution. In this case, our objective becomes  $\mathcal{L}(\psi) = \sum_{\mathbf{z}} \ell(\psi, \mathbf{z}) q_\psi(\mathbf{z})$ . For simplicity, let us assume  $\ell(\psi, \mathbf{z}) = \ell(\mathbf{z})$ .

We can now easily compute gradients using  $\nabla_\psi \mathcal{L}(\psi) = \sum_{\mathbf{z}} \ell(\mathbf{z}) \nabla_\psi q_\psi(\mathbf{z})$ . Of course, if  $\mathbf{z}$  can take on exponentially many values (e.g., we are optimizing over the space of strings), this expression is intractable. However, suppose we can partition this sum into two sets, a small set  $S_1$  of high probability values and a large set  $S_2$  of all other values. Then we can enumerate over  $S_1$  and use the score function estimator for  $S_2$ :

$$\nabla_\psi \mathcal{L}(\psi) = \sum_{\mathbf{z} \in S_1} \ell(\mathbf{z}) \nabla_\psi q_\psi(\mathbf{z}) + \mathbb{E}_{q_\psi(\mathbf{z}|\mathbf{z} \in S_2)} [\ell(\mathbf{z}) \nabla_\psi \log q_\psi(\mathbf{z})] \quad (6.99)$$

To compute the second expectation, we can use rejection sampling applied to samples from  $q_\psi(\mathbf{z})$ . This procedure is a form of Rao-Blackwellisation as shown in [Liu+19b], and reduces the variance compared to standard SFE (see Section 11.6.2 for details on Rao-Blackwellisation).

### 6.5.4 Reparameterization trick

The score function estimator can have high variance, even when using a control variate. In this section, we derive a lower variance estimator, which can be applied if  $\ell(\psi, \mathbf{z})$  is differentiable wrt  $\mathbf{z}$ . We additionally require that we can compute a sample from  $q_\psi(\mathbf{z})$  by first sampling  $\epsilon$  from some noise distribution  $q_0$  which is independent of  $\psi$ , and then transforming to  $\mathbf{z}$  using a deterministic and differentiable function  $\mathbf{z} = r(\psi, \epsilon)$ . For example, instead of sampling  $\mathbf{z} \sim \mathcal{N}(\mu, \sigma^2)$ , we can sample  $\epsilon \sim \mathcal{N}(0, 1)$  and compute

$$\mathbf{z} = r(\psi, \epsilon) = \mu + \sigma \epsilon \quad (6.100)$$

where  $\psi = (\mu, \sigma)$ . This allows us to rewrite our stochastic objective as follows:

$$\mathcal{L}(\psi) = \mathbb{E}_{q_\psi(\mathbf{z})} [\ell(\psi, \mathbf{z})] = \mathbb{E}_{q_0(\epsilon)} [\ell(\psi, r(\psi, \epsilon))] \quad (6.101)$$

Since  $q_0(\epsilon)$  is independent of  $\psi$ , we can push the gradient operator inside the expectation, which we can approximate with Monte Carlo:

$$\nabla_\psi \mathcal{L}(\psi) = \mathbb{E}_{q_0(\epsilon)} [\nabla_\psi \ell(\psi, r(\psi, \epsilon))] \approx \frac{1}{S} \sum_{s=1}^S \nabla_\psi \ell(\psi, r(\psi, \epsilon_s)) \quad (6.102)$$

where  $\epsilon_s \sim q_0$ . This is called the **reparameterization gradient** or the **pathwise derivative** [Gla03; Fu15; KW14; RMW14a; TLG14; JO18; FMM18], and is widely used in variational inference (Section 10.3.3), and when fitting VAE models (Section 21.2).

#### 6.5.4.1 Example

As a simple example, suppose we define some arbitrary function, such as  $\ell(z) = z^2 - 3z$ , and then define its expected value as  $\mathcal{L}(\psi) = \mathbb{E}_{\mathcal{N}(z|\mu, v)} [\ell(z)]$ , where  $\psi = (\mu, v)$  and  $v = \sigma^2$ . Suppose we want to compute

$$\nabla_\psi \mathcal{L}(\psi) = \left[ \frac{\partial}{\partial \mu} \mathbb{E} [\ell(z)], \frac{\partial}{\partial v} \mathbb{E} [\ell(z)] \right] \quad (6.103)$$

Since the Gaussian distribution is reparameterizable, we can sample  $z \sim \mathcal{N}(z|\mu, v)$ , and then use automatic differentiation to compute each of these gradient terms, and then average.

However, in the special case of Gaussian distributions, we can also compute the gradient vector directly. In particular, in Section 6.4.5.1 we present Bonnet's theorem, which states that

$$\frac{\partial}{\partial \mu} \mathbb{E} [\ell(z)] = \mathbb{E} \left[ \frac{\partial}{\partial z} \ell(z) \right] \quad (6.104)$$

Similarly, Price's theorem states that

$$\frac{\partial}{\partial v} \mathbb{E} [\ell(z)] = 0.5 \mathbb{E} \left[ \frac{\partial^2}{\partial z^2} \ell(z) \right] \quad (6.105)$$

In `gradient_expected_value_gaussian.ipynb` we show that these two methods are numerically equivalent, as theory suggests.

#### 6.5.4.2 Total derivative

To compute the gradient term inside the expectation in Equation (6.102) we need to use the **total derivative**, since the function  $\ell$  depends on  $\psi$  directly and via the noise sample. Recall that, for a function of the form  $f(\psi_1, \dots, \psi_{d_\psi}, z_1(\psi), \dots, z_{d_z}(\psi))$ , the total derivative wrt  $\psi_i$  is given by the chain rule as follows:

$$\frac{\partial \ell}{\partial \psi_i}^{\text{TD}} = \frac{\partial \ell}{\partial \psi_i} + \sum_j \frac{\partial \ell}{\partial z_j} \frac{\partial z_j}{\partial \psi_i} \quad (6.106)$$

and hence

$$\nabla_{\psi} \ell(\psi, z)^{\text{TD}} = \nabla_{\psi} \ell(\psi, z) + \mathbf{J}^T \nabla_z \ell(\psi, z) \quad (6.107)$$

where  $\mathbf{J} = \frac{\partial z^T}{\partial \psi}$  is the  $d_z \times d_\psi$  Jacobian matrix of the noise transformation:

$$\mathbf{J} = \begin{pmatrix} \frac{\partial z_1}{\partial \psi_1} & \dots & \frac{\partial z_1}{\partial \psi_{d_\psi}} \\ \vdots & \ddots & \vdots \\ \frac{\partial z_{d_z}}{\partial \psi_1} & \dots & \frac{\partial z_{d_z}}{\partial \psi_{d_\psi}} \end{pmatrix} \quad (6.108)$$

Hence we can compute the gradient using

$$\nabla_{\psi} \mathcal{L}(\psi) = \mathbb{E}_{q_0(\epsilon)} [\nabla_{\psi} \ell(\psi, z) + \mathbf{J}(\psi, \epsilon)^T \nabla_z \ell(\psi, r(\psi, \epsilon))] \quad (6.109)$$

We leverage this decomposition in Section 10.3.3.1, where we derive a lower variance gradient estimator in the special case of variational inference.

---

### 6.5.5 The delta method

The **delta method** [Hoe12] approximates the expectation of a function of a random variable by the expectation of the function's Taylor expansion. For example, suppose  $\boldsymbol{\theta} \sim q$  with mean  $\mathbb{E}_{q(\boldsymbol{\theta})}[\boldsymbol{\theta}] = \mathbf{m}$ , and we use a first order expansion. Then we have

$$\mathbb{E}_{q(\boldsymbol{\theta})}[f(\boldsymbol{\theta})] \approx \mathbb{E}_{q(\boldsymbol{\theta})}[f(\mathbf{m}) + (\boldsymbol{\theta} - \mathbf{m})^\top \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\mathbf{m}}] \approx f(\mathbf{m}) \quad (6.110)$$

Now let  $f(\boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta})$ . Then we have

$$\mathbb{E}_{q(\boldsymbol{\theta})}[\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta})] \approx \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\mathbf{m}} \quad (6.111)$$

This is called the **first-order delta method**.

Now let  $f(\boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}}^2 \mathcal{L}(\boldsymbol{\theta})$ . Then we have

$$\mathbb{E}_{q(\boldsymbol{\theta})}[\nabla_{\boldsymbol{\theta}}^2 \mathcal{L}(\boldsymbol{\theta})] \approx \nabla_{\boldsymbol{\theta}}^2 \mathcal{L}(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\mathbf{m}} \quad (6.112)$$

This is called the **second-order delta method**.

### 6.5.6 Gumbel softmax trick

When working with discrete variables, we cannot use the reparameterization trick. However, we can often relax the discrete variables to continuous ones in a way which allows the trick to be used, as we explain below.

Consider a one-hot vector  $\mathbf{d}$  with  $K$  bits, so  $d_k \in \{0, 1\}$  and  $\sum_{k=1}^K d_k = 1$ . This can be used to represent a  $K$ -ary categorical variable  $d$ . Let  $P(d) = \text{Cat}(d|\boldsymbol{\pi})$ , where  $\pi_k = P(d_k = 1)$ , so  $0 \leq \pi_k \leq 1$ . Alternatively we can parameterize the distribution in terms of  $(\alpha_1, \dots, \alpha_K)$ , where  $\pi_k = \alpha_k / (\sum_{k'=1}^K \alpha_{k'})$ . We will denote this by  $d \sim \text{Cat}(d|\boldsymbol{\alpha})$ .

We can sample a one-hot vector  $\mathbf{d}$  from this distribution by computing

$$\mathbf{d} = \text{onehot}(\underset{k}{\operatorname{argmax}} [\epsilon_k + \log \alpha_k]) \quad (6.113)$$

where  $\epsilon_k \sim \text{Gumbel}(0, 1)$  is sampled from the **Gumbel distribution** [Gum54]. We can draw such samples by first sampling  $u_k \sim \text{Unif}(0, 1)$  and then computing  $\epsilon_k = -\log(-\log(u_k))$ . This is called the **Gumbel-Max trick** [MTM14], and gives us a reparameterizable representation for the categorical distribution.

Unfortunately, the derivative of the argmax is 0 everywhere except at the boundary of transitions from one label to another, where the derivative is undefined. However, suppose we replace the argmax with a softmax, and replace the discrete one-hot vector  $\mathbf{d}$  with a continuous relaxation  $\mathbf{x} \in \Delta^{K-1}$ , where  $\Delta^{K-1} = \{\mathbf{x} \in \mathbb{R}^K : x_k \in [0, 1], \sum_{k=1}^K x_k = 1\}$  is the  $K$ -dimensional simplex. Then we can write

$$x_k = \frac{\exp((\log \alpha_k + \epsilon_k)/\tau)}{\sum_{k'=1}^K \exp((\log \alpha_{k'} + \epsilon_{k'})/\tau)} \quad (6.114)$$

where  $\tau > 0$  is a temperature parameter. This is called the **Gumbel-Softmax distribution** [JGP17] or the **concrete distribution** [MMT17]. This smoothly approaches the discrete distribution as  $\tau \rightarrow 0$ , as illustrated in Figure 6.4.

We can now replace  $f(\mathbf{d})$  with  $f(\mathbf{x})$ , which allows us to take reparameterized gradients wrt  $\mathbf{x}$ .

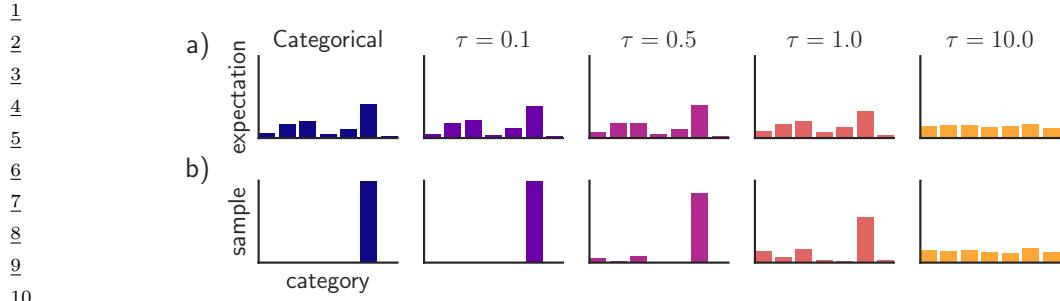


Figure 6.4: Illustration of the Gumbel-Softmax (concrete) distribution with  $K = 7$  states at different temperatures  $\tau$ . The top row shows  $\mathbb{E}[z]$ , and the bottom row shows samples  $z \sim \text{GumbelSoftmax}(\alpha, \tau)$ . The left column shows a discrete (categorical) distribution, which always produces one-hot samples. From Figure 1 of [JGP17]. Used with kind permission of Ben Poole.

### 6.5.7 Stochastic computation graphs

We can represent an arbitrary function containing both deterministic and stochastic components as a **stochastic computation graph**. We can then generalize the AD algorithm (Section 6.2) to leverage score function estimation (Section 6.5.3) and reparameterization (Section 6.5.4) to compute Monte Carlo gradients for complex nested functions. For details, see [Sch+15a; Gaj+19].

### 6.5.8 Straight-through estimator

In this section, we discuss how to approximate the gradient of a quantized version of a signal. For example, suppose we have the following thresholding function, that binarizes its output:

$$f(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases} \quad (6.115)$$

This does not have a well-defined gradient. However, we can use the **straight-through estimator** proposed in [Ben13] as an approximation. The basic idea is to replace  $g(x) = f'(x)$ , where  $f'(x)$  is the derivative of  $f$  wrt input, with  $g(x) = x$  when computing the backwards pass. See Figure 6.5 for a visualization, and [Yin+19b] for an analysis of why this is a valid approximation.

In practice, we sometimes replace  $g(x) = x$  with the **hard tanh** function, defined by

$$\text{HardTanh}(x) = \begin{cases} x & \text{if } -1 \leq x \leq 1 \\ 1 & \text{if } x > 1 \\ -1 & \text{if } x < -1 \end{cases} \quad (6.116)$$

This ensures the gradients that are backpropagated don't get too large. See Section 21.6 for an application of this approach to discrete autoencoders.

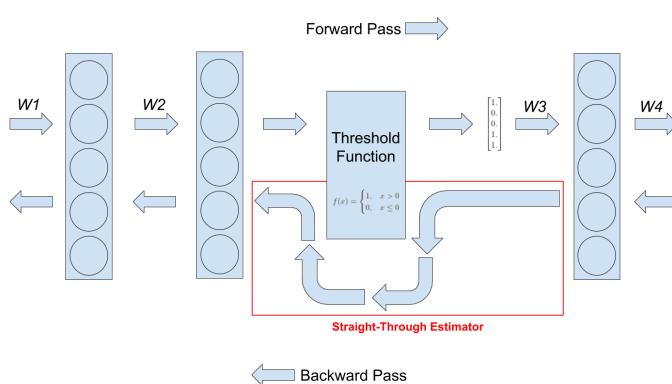


Figure 6.5: Illustration of straight-through estimator when applied to a binary threshold function in the middle of an MLP. From <https://www.hassanaskary.com/python/pytorch/deep%20learning/2020/09/19/intuitive-explanation-of-straight-through-estimators.html>. Used with kind permission of Hassan Askary.

## 6.6 Bound optimization (MM) algorithms

In this section, we consider a class of algorithms known as **bound optimization** or **MM** algorithms. In the context of minimization, MM stands for **majorize-minimize**. In the context of maximization, MM stands for **minorize-maximize**. There are many examples of MM algorithms, such as EM (Section 6.6.3), proximal gradient methods (Section 4.1), the mean shift algorithm for clustering [FH75; Che95; FT05], etc. For more details, see e.g., [HL04; Mai15; SBP17; Nad+19],

### 6.6.1 The general algorithm

In this section, we assume our goal is to *maximize* some function  $\ell(\boldsymbol{\theta})$  wrt its parameters  $\boldsymbol{\theta}$ . The basic approach in MM algorithms is to construct a **surrogate function**  $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^t)$  which is a tight lowerbound to  $\ell(\boldsymbol{\theta})$  such that  $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^t) \leq \ell(\boldsymbol{\theta})$  and  $Q(\boldsymbol{\theta}^t, \boldsymbol{\theta}^t) = \ell(\boldsymbol{\theta}^t)$ . If these conditions are met, we say that  $Q$  minorizes  $\ell$ . We then perform the following update at each step:

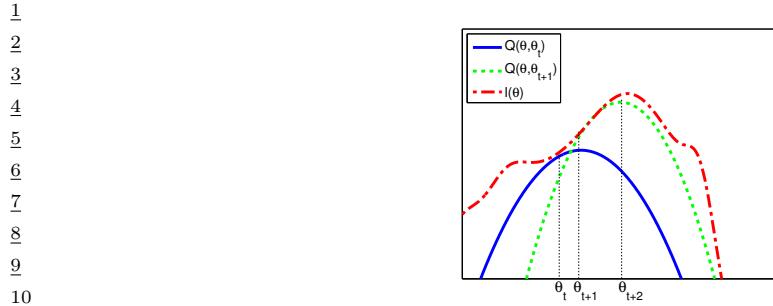
$$\boldsymbol{\theta}^{t+1} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^t) \quad (6.117)$$

This guarantees us monotonic increases in the original objective:

$$\ell(\boldsymbol{\theta}^{t+1}) \geq Q(\boldsymbol{\theta}^{t+1}, \boldsymbol{\theta}^t) \geq Q(\boldsymbol{\theta}^t, \boldsymbol{\theta}^t) = \ell(\boldsymbol{\theta}^t) \quad (6.118)$$

where the first inequality follows since  $Q(\boldsymbol{\theta}^t, \boldsymbol{\theta}')$  is a lower bound on  $\ell(\boldsymbol{\theta}^t)$  for any  $\boldsymbol{\theta}'$ ; the second inequality follows from Equation (6.117); and the final equality follows the tightness property. As a consequence of this result, if you do not observe monotonic increase of the objective, you must have an error in your math and/or code. This is a surprisingly powerful debugging tool.

This process is sketched in Figure 6.6. The dashed red curve is the original function (e.g., the log-likelihood of the observed data). The solid blue curve is the lower bound, evaluated at  $\boldsymbol{\theta}^t$ ; this



11 *Figure 6.6: Illustration of a bound optimization algorithm. Adapted from Figure 9.14 of [Bis06]. Generated*  
12 *by em\_log\_likelihood\_max.ipynb.*

13

14 touches the objective function at  $\theta^t$ . We then set  $\theta^{t+1}$  to the maximum of the lower bound (blue  
15 curve), and fit a new bound at that point (dotted green curve). The maximum of this new bound  
16 becomes  $\theta^{t+2}$ , etc.  
17

18

19 **6.6.2 Example: logistic regression**  
20

21 If  $\ell(\theta)$  is a concave function we want to maximize, then one way to obtain a valid lower bound is to  
22 use a bound on its Hessian, i.e., to find a matrix a negative definite matrix  $\mathbf{B}$  such that  $\mathbf{H}(\theta) \succ \mathbf{B}$ .  
23 In this case, one can show (see [BCN18, App. B]) that

24

$$\ell(\theta) \geq \ell(\theta^t) + (\theta - \theta^t)^T g(\theta^t) + \frac{1}{2}(\theta - \theta^t)^T \mathbf{B}(\theta - \theta^t) \quad (6.119)$$

25

26 where  $g(\theta^t) = \nabla \ell(\theta^t)$ . Therefore the following function is a valid lower bound:

27

$$Q(\theta, \theta^t) = \theta^T (g(\theta^t) - \mathbf{B}\theta^t) + \frac{1}{2}\theta^T \mathbf{B}\theta \quad (6.120)$$

28

29 The corresponding update becomes

30

$$\theta^{t+1} = \theta^t - \mathbf{B}^{-1}g(\theta^t) \quad (6.121)$$

31

32 This is similar to a Newton update, except we use  $\mathbf{B}$ , which is a fixed matrix, rather than  $\mathbf{H}(\theta^t)$ ,  
33 which changes at each iteration. This can give us some of the advantages of second order methods at  
34 lower computational cost.

35 For example, let us fit a multi-class logistic regression model using MM. (We follow the presentation  
36 of [Kri+05], who also consider the more interesting case of *sparse* logistic regression.) The probability  
37 that example  $n$  belongs to class  $c \in \{1, \dots, C\}$  is given by  
38

39

$$p(y_n = c | \mathbf{x}_n, \mathbf{w}) = \frac{\exp(\mathbf{w}_c^T \mathbf{x}_n)}{\sum_{i=1}^C \exp(\mathbf{w}_i^T \mathbf{x}_n)} \quad (6.122)$$

40

41 Because of the normalization condition  $\sum_{c=1}^C p(y_n = c | \mathbf{x}_n, \mathbf{w}) = 1$ , we can set  $\mathbf{w}_C = \mathbf{0}$ . (For example,  
42 in binary logistic regression, where  $C = 2$ , we only learn a single weight vector.) Therefore the  
43 parameters  $\theta$  correspond to a weight matrix  $\mathbf{w}$  of size  $D(C - 1)$ , where  $\mathbf{x}_n \in \mathbb{R}^D$ .

44

If we let  $\mathbf{p}_n(\mathbf{w}) = [p(y_n = 1|\mathbf{x}_n, \mathbf{w}), \dots, p(y_n = C-1|\mathbf{x}_n, \mathbf{w})]$  and  $\mathbf{y}_n = [\mathbb{I}(y_n = 1), \dots, \mathbb{I}(y_n = C-1)]$ , we can write the log-likelihood as follows:

$$\ell(\mathbf{w}) = \sum_{n=1}^N \left[ \sum_{c=1}^{C-1} y_{nc} \mathbf{w}_c^\top \mathbf{x}_n - \log \sum_{c=1}^C \exp(\mathbf{w}_c^\top \mathbf{x}_n) \right] \quad (6.123)$$

The gradient is given by the following:

$$\mathbf{g}(\mathbf{w}) = \sum_{n=1}^N (\mathbf{y}_n - \mathbf{p}_n(\mathbf{w})) \otimes \mathbf{x}_n \quad (6.124)$$

where  $\otimes$  denotes kronecker product (which, in this case, is just outer product of the two vectors).

The Hessian is given by the following:

$$\mathbf{H}(\mathbf{w}) = - \sum_{n=1}^N (\text{diag}(\mathbf{p}_n(\mathbf{w})) - \mathbf{p}_n(\mathbf{w})\mathbf{p}_n(\mathbf{w})^\top) \otimes (\mathbf{x}_n \mathbf{x}_n^\top) \quad (6.125)$$

We can construct a lower bound on the Hessian, as shown in [Boh92]:

$$\mathbf{H}(\mathbf{w}) \succ -\frac{1}{2} [\mathbf{I} - \mathbf{1}\mathbf{1}^\top/C] \otimes (\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top) \triangleq \mathbf{B} \quad (6.126)$$

where  $\mathbf{I}$  is a  $(C-1)$ -dimensional identity matrix, and  $\mathbf{1}$  is a  $(C-1)$ -dimensional vector of all 1s. In the binary case, this becomes

$$\mathbf{H}(\mathbf{w}) \succ -\frac{1}{2} \left(1 - \frac{1}{2}\right) \left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top\right) = -\frac{1}{4} \mathbf{X}^\top \mathbf{X} \quad (6.127)$$

This follows since  $p_n \leq 0.5$  so  $-(p_n - p_n^2) \geq -0.25$ .

We can use this lower bound to construct an MM algorithm to find the MLE. The update becomes

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \mathbf{B}^{-1} \mathbf{g}(\mathbf{w}^t) \quad (6.128)$$

For example, let us consider the binary case, so  $\mathbf{g}^t = \nabla \ell(\mathbf{w}^t) = \mathbf{X}^\top (\mathbf{y} - \boldsymbol{\mu}^t)$ , where  $\boldsymbol{\mu}^t = [\mathbf{p}_n(\mathbf{w}^t), (1 - \mathbf{p}_n(\mathbf{w}^t))]_{n=1}^N$ . The update becomes

$$\mathbf{w}^{t+1} = \mathbf{w}^t - 4(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{g}^t \quad (6.129)$$

The above is faster (per step) than the IRLS (iteratively reweighted least squares) algorithm (i.e., Newton's method), which is the standard method for fitting GLMs. To see this, note that the Newton update has the form

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \mathbf{H}^{-1} \mathbf{g}(\mathbf{w}^t) = \mathbf{w}^t - (\mathbf{X}^\top \mathbf{S}^t \mathbf{X})^{-1} \mathbf{g}^t \quad (6.130)$$

where  $\mathbf{S}^t = \text{diag}(\boldsymbol{\mu}^t \odot (1 - \boldsymbol{\mu}^t))$ . We see that Equation (6.129) is faster to compute, since we can precompute the constant matrix  $(\mathbf{X}^\top \mathbf{X})^{-1}$ .

1 **6.6.3 The EM algorithm**

3 In this section, we discuss the **expectation maximization (EM)** algorithm [DLR77; MK07], which  
4 is an algorithm designed to compute the MLE or MAP parameter estimate for probability models  
5 that have **missing data** and/or **hidden variables**. It is a special case of an MM algorithm.

6 The basic idea behind EM is to alternate between estimating the hidden variables (or missing  
7 values) during the **E step** (expectation step), and then using the fully observed data to compute the  
8 MLE during the **M step** (maximization step). Of course, we need to iterate this process, since the  
9 expected values depend on the parameters, but the parameters depend on the expected values.

10 In Section 6.6.3.1, we show that EM is a **bound optimization** algorithm, which implies that this  
11 iterative procedure will converge to a local maximum of the log likelihood. The speed of convergence  
12 depends on the amount of missing data, which affects the tightness of the bound [XJ96; MD97;  
13 SRG03; KKS20].

14 We now describe the EM algorithm for a generic model. We let  $\mathbf{y}_n$  be the visible data for example  
15  $n$ , and  $\mathbf{z}_n$  be the hidden data.

17 **6.6.3.1 Lower bound**

19 The goal of EM is to maximize the log likelihood of the observed data:

$$\ell(\boldsymbol{\theta}) = \sum_{n=1}^{N_D} \log p(\mathbf{y}_n | \boldsymbol{\theta}) = \sum_{n=1}^{N_D} \log \left[ \sum_{\mathbf{z}_n} p(\mathbf{y}_n, \mathbf{z}_n | \boldsymbol{\theta}) \right] \quad (6.131)$$

24 where  $\mathbf{y}_n$  are the visible variables and  $\mathbf{z}_n$  are the hidden variables. Unfortunately this is hard to  
25 optimize, since the log cannot be pushed inside the sum.

26 EM gets around this problem as follows. First, consider a set of arbitrary distributions  $q_n(\mathbf{z}_n)$  over  
27 each hidden variable  $\mathbf{z}_n$ . The observed data log likelihood can be written as follows:

$$\ell(\boldsymbol{\theta}) = \sum_{n=1}^{N_D} \log \left[ \sum_{\mathbf{z}_n} q_n(\mathbf{z}_n) \frac{p(\mathbf{y}_n, \mathbf{z}_n | \boldsymbol{\theta})}{q_n(\mathbf{z}_n)} \right] \quad (6.132)$$

32 Using Jensen's inequality, we can push the log (which is a concave function) inside the expectation  
33 to get the following lower bound on the log likelihood:

$$\ell(\boldsymbol{\theta}) \geq \sum_n \sum_{\mathbf{z}_n} q_n(\mathbf{z}_n) \log \frac{p(\mathbf{y}_n, \mathbf{z}_n | \boldsymbol{\theta})}{q_n(\mathbf{z}_n)} \quad (6.133)$$

$$= \sum_n \underbrace{\mathbb{E}_{q_n} [\log p(\mathbf{y}_n, \mathbf{z}_n | \boldsymbol{\theta})]}_{\mathcal{L}(\boldsymbol{\theta}, q_n | \mathbf{y}_n)} + \mathbb{H}(q_n) \quad (6.134)$$

$$= \sum_n \mathcal{L}(\boldsymbol{\theta}, q_n | \mathbf{y}_n) \triangleq \mathcal{L}(\boldsymbol{\theta}, \{q_n\} | \mathcal{D}) \quad (6.135)$$

43 where  $\mathbb{H}(q)$  is the entropy of probability distribution  $q$ , and  $\mathcal{L}(\boldsymbol{\theta}, \{q_n\} | \mathcal{D})$  is called the **evidence**  
44 **lower bound** or **ELBO**, since it is a lower bound on the log marginal likelihood,  $\log p(\mathbf{y}_{1:N} | \boldsymbol{\theta})$ , also  
45 called the evidence. Optimizing this bound is the basis of variational inference, as we discuss in  
46 Section 10.1.

47

---

### 6.6.3.2 E step

We see that the lower bound is a sum of  $N$  terms, each of which has the following form:

$$\mathbb{L}(\boldsymbol{\theta}, q_n | \mathbf{y}_n) = \sum_{\mathbf{z}_n} q_n(\mathbf{z}_n) \log \frac{p(\mathbf{y}_n, \mathbf{z}_n | \boldsymbol{\theta})}{q_n(\mathbf{z}_n)} \quad (6.136)$$

$$= \sum_{\mathbf{z}_n} q_n(\mathbf{z}_n) \log \frac{p(\mathbf{z}_n | \mathbf{y}_n, \boldsymbol{\theta}) p(\mathbf{y}_n | \boldsymbol{\theta})}{q_n(\mathbf{z}_n)} \quad (6.137)$$

$$= \sum_{\mathbf{z}_n} q_n(\mathbf{z}_n) \log \frac{p(\mathbf{z}_n | \mathbf{y}_n, \boldsymbol{\theta})}{q_n(\mathbf{z}_n)} + \sum_{\mathbf{z}_n} q_n(\mathbf{z}_n) \log p(\mathbf{y}_n | \boldsymbol{\theta}) \quad (6.138)$$

$$= -D_{\text{KL}}(q_n(\mathbf{z}_n) \| p(\mathbf{z}_n | \mathbf{y}_n, \boldsymbol{\theta})) + \log p(\mathbf{y}_n | \boldsymbol{\theta}) \quad (6.139)$$

where  $D_{\text{KL}}(q \| p) \triangleq \sum_z q(z) \log \frac{q(z)}{p(z)}$  is the Kullback-Leibler divergence (or KL divergence for short) between probability distributions  $q$  and  $p$ . We discuss this in more detail in Section 5.1, but the key property we need here is that  $D_{\text{KL}}(q \| p) \geq 0$  and  $D_{\text{KL}}(q \| p) = 0$  iff  $q = p$ . Hence we can maximize the lower bound  $\mathbb{L}(\boldsymbol{\theta}, \{q_n\} | \mathcal{D})$  wrt  $\{q_n\}$  by setting each one to  $q_n^* = p(\mathbf{z}_n | \mathbf{y}_n, \boldsymbol{\theta})$ . This is called the **E step**. This ensures the ELBO is a tight lower bound:

$$\mathbb{L}(\boldsymbol{\theta}, \{q_n^*\} | \mathcal{D}) = \sum_n \log p(\mathbf{y}_n | \boldsymbol{\theta}) = \ell(\boldsymbol{\theta} | \mathcal{D}) \quad (6.140)$$

To see how this connects to bound optimization, let us define

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^t) = \mathbb{L}(\boldsymbol{\theta}, \{p(\mathbf{z}_n | \mathbf{y}_n; \boldsymbol{\theta}^t)\}) \quad (6.141)$$

Then we have  $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^t) \leq \ell(\boldsymbol{\theta})$  and  $Q(\boldsymbol{\theta}^t, \boldsymbol{\theta}^t) = \ell(\boldsymbol{\theta}^t)$ , as required.

However, if we cannot compute the posteriors  $p(\mathbf{z}_n | \mathbf{y}_n; \boldsymbol{\theta}^t)$  exactly, we can still use an approximate distribution  $q(\mathbf{z}_n | \mathbf{y}_n; \boldsymbol{\theta}^t)$ ; this will yield a non-tight lower-bound on the log-likelihood. This generalized version of EM is known as variational EM [NH98b]. See Section 6.6.6.1 for details.

### 6.6.3.3 M step

In the M step, we need to maximize  $\mathbb{L}(\boldsymbol{\theta}, \{q_n^t\})$  wrt  $\boldsymbol{\theta}$ , where the  $q_n^t$  are the distributions computed in the E step at iteration  $t$ . Since the entropy terms  $\mathbb{H}(q_n)$  are constant wrt  $\boldsymbol{\theta}$ , so we can drop them in the M step. We are left with

$$\ell^t(\boldsymbol{\theta}) = \sum_n \mathbb{E}_{q_n^t(\mathbf{z}_n)} [\log p(\mathbf{y}_n, \mathbf{z}_n | \boldsymbol{\theta})] \quad (6.142)$$

This is called the **expected complete data log likelihood**. If the joint probability is in the exponential family (Section 2.3), we can rewrite this as

$$\ell^t(\boldsymbol{\theta}) = \sum_n \mathbb{E} [\mathcal{T}(\mathbf{y}_n, \mathbf{z}_n)^T \boldsymbol{\theta} - A(\boldsymbol{\theta})] = \sum_n (\mathbb{E} [\mathcal{T}(\mathbf{y}_n, \mathbf{z}_n)]^T \boldsymbol{\theta} - A(\boldsymbol{\theta})) \quad (6.143)$$

where  $\mathbb{E} [\mathcal{T}(\mathbf{y}_n, \mathbf{z}_n)]$  are called the **expected sufficient statistics**.

1 In the M step, we maximize the expected complete data log likelihood to get  
2

$$\underline{3} \quad \underline{4} \quad \boldsymbol{\theta}^{t+1} = \arg \max_{\boldsymbol{\theta}} \sum_n \mathbb{E}_{q_n^t} [\log p(\mathbf{y}_n, \mathbf{z}_n | \boldsymbol{\theta})] \quad (6.144)$$

5

6 In the case of the exponential family, the maximization can be solved in closed-form by matching the  
7 moments of the expected sufficient statistics (Section 2.3.5).

8 We see from the above that the E step does not in fact need to return the full set of posterior  
9 distributions  $\{q(\mathbf{z}_n)\}$ , but can instead just return the sum of the expected sufficient statistics,  
10  $\sum_n \mathbb{E}_{q(\mathbf{z}_n)} [\mathcal{T}(\mathbf{y}_n, \mathbf{z}_n)]$ .

11 A common application of EM is for fitting mixture models; we discuss this in the prequel to this  
12 book, [Mur22]. Below we give a different example.

13

#### 14 6.6.4 Example: EM for an MVN with missing data

15 It is easy to compute the MLE for a multivariate normal when we have a fully observed data matrix:  
16 we just compute the sample mean and covariance. In this section, we consider the case where we have  
17 **missing data or partially observed data**. For example, we can think of the entries of  $\mathbf{Y}$  as being  
18 answers to a survey; some of these answers may be unknown. There are many kinds of missing data,  
19 as we discuss in Section 21.3.4. In this section, we make the missing at random (MAR) assumption,  
20 for simplicity. Under the MAR assumption, the log likelihood of the visible data has the form  
21

$$\underline{22} \quad \underline{23} \quad \log p(\mathbf{X} | \boldsymbol{\theta}) = \sum_n \log p(\mathbf{x}_n | \boldsymbol{\theta}) = \sum_n \log \left[ \int p(\mathbf{x}_n, \mathbf{z}_n | \boldsymbol{\theta}) d\mathbf{z}_n \right] \quad (6.145)$$

24

25 where  $\mathbf{x}_n$  are the visible variables in case  $n$ ,  $\mathbf{z}_n$  are the hidden variables, and  $\mathbf{y}_n = (\mathbf{z}_n, \mathbf{x}_n)$  are all  
26 the variables. Unfortunately, this objective is hard to maximize. since we cannot push the log inside  
27 the expectation. Fortunately, we can easily apply EM, as we explain below.  
28

##### 29 6.6.4.1 E step

30 Suppose we have the parameters  $\boldsymbol{\theta}^{t-1}$  from the previous iteration. Then we can compute the expected  
31 complete data log likelihood at iteration  $t$  as follows:

$$\underline{33} \quad \underline{34} \quad Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{t-1}) = \mathbb{E} \left[ \sum_{n=1}^{N_D} \log \mathcal{N}(\mathbf{y}_n | \boldsymbol{\mu}, \boldsymbol{\Sigma}) | \mathcal{D}, \boldsymbol{\theta}^{t-1} \right] \quad (6.146)$$

35

$$\underline{36} \quad \underline{37} \quad = -\frac{N_D}{2} \log |2\pi\boldsymbol{\Sigma}| - \frac{1}{2} \sum_n \mathbb{E} [(\mathbf{y}_n - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{y}_n - \boldsymbol{\mu})] \quad (6.147)$$

38

$$\underline{39} \quad \underline{40} \quad = -\frac{N_D}{2} \log |2\pi\boldsymbol{\Sigma}| - \frac{1}{2} \text{tr}(\boldsymbol{\Sigma}^{-1} \sum_n \mathbb{E} [(\mathbf{y}_n - \boldsymbol{\mu})(\mathbf{y}_n - \boldsymbol{\mu})^\top]) \quad (6.148)$$

41

$$\underline{42} \quad \underline{43} \quad = -\frac{N_D}{2} \log |\boldsymbol{\Sigma}| - \frac{N_D D}{2} \log(2\pi) - \frac{1}{2} \text{tr}(\boldsymbol{\Sigma}^{-1} \mathbb{E} [\mathbf{S}(\boldsymbol{\mu})]) \quad (6.149)$$

44

where

$$\underline{45} \quad \mathbb{E} [\mathbf{S}(\boldsymbol{\mu})] \triangleq \sum_n \left( \mathbb{E} [\mathbf{y}_n \mathbf{y}_n^\top] + \boldsymbol{\mu} \boldsymbol{\mu}^\top - 2\boldsymbol{\mu} \mathbb{E} [\mathbf{y}_n]^\top \right) \quad (6.150)$$

46

(We drop the conditioning of the expectation on  $\mathcal{D}$  and  $\boldsymbol{\theta}^{t-1}$  for brevity.) We see that we need to compute  $\sum_n \mathbb{E}[\mathbf{y}_n]$  and  $\sum_n \mathbb{E}[\mathbf{y}_n \mathbf{y}_n^\top]$ ; these are the expected sufficient statistics.

To compute these quantities, we use the results from Section 2.2.5.3. We have

$$p(\mathbf{z}_n | \mathbf{x}_n, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{z}_n | \mathbf{m}_n, \mathbf{V}_n) \quad (6.151)$$

$$\mathbf{m}_n \triangleq \boldsymbol{\mu}_h + \boldsymbol{\Sigma}_{hv} \boldsymbol{\Sigma}_{vv}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_v) \quad (6.152)$$

$$\mathbf{V}_n \triangleq \boldsymbol{\Sigma}_{hh} - \boldsymbol{\Sigma}_{hv} \boldsymbol{\Sigma}_{vv}^{-1} \boldsymbol{\Sigma}_{vh} \quad (6.153)$$

where we partition  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  into blocks based on the hidden and visible indices  $h$  and  $v$ . Hence the expected sufficient statistics are

$$\mathbb{E}[\mathbf{y}_n] = (\mathbb{E}[\mathbf{z}_n]; \mathbf{x}_n) = (\mathbf{m}_n; \mathbf{x}_n) \quad (6.154)$$

To compute  $\mathbb{E}[\mathbf{y}_n \mathbf{y}_n^\top]$ , we use the result that  $\text{Cov}[\mathbf{y}] = \mathbb{E}[\mathbf{y}\mathbf{y}^\top] - \mathbb{E}[\mathbf{y}]\mathbb{E}[\mathbf{y}^\top]$ . Hence

$$\mathbb{E}[\mathbf{y}_n \mathbf{y}_n^\top] = \mathbb{E} \left[ \begin{pmatrix} \mathbf{z}_n \\ \mathbf{x}_n \end{pmatrix} \begin{pmatrix} \mathbf{z}_n^\top & \mathbf{x}_n^\top \end{pmatrix} \right] = \begin{pmatrix} \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^\top] & \mathbb{E}[\mathbf{z}_n] \mathbf{x}_n^\top \\ \mathbf{x}_n \mathbb{E}[\mathbf{z}_n]^\top & \mathbf{x}_n \mathbf{x}_n^\top \end{pmatrix} \quad (6.155)$$

$$\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^\top] = \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^\top + \mathbf{V}_n \quad (6.156)$$

#### 6.6.4.2 M step

By solving  $\nabla Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t-1)}) = \mathbf{0}$ , we can show that the M step is equivalent to plugging these ESS into the usual MLE equations to get

$$\boldsymbol{\mu}^t = \frac{1}{N_{\mathcal{D}}} \sum_n \mathbb{E}[\mathbf{y}_n] \quad (6.157)$$

$$\boldsymbol{\Sigma}^t = \frac{1}{N_{\mathcal{D}}} \sum_n \mathbb{E}[\mathbf{y}_n \mathbf{y}_n^\top] - \boldsymbol{\mu}^t (\boldsymbol{\mu}^t)^\top \quad (6.158)$$

Thus we see that EM is *not* equivalent to simply replacing variables by their expectations and applying the standard MLE formula; that would ignore the posterior variance and would result in an incorrect estimate. Instead we must compute the expectation of the sufficient statistics, and plug that into the usual equation for the MLE.

#### 6.6.4.3 Initialization

To get the algorithm started, we can compute the MLE based on those rows of the data matrix that are fully observed. If there are no such rows, we can just estimate the diagonal terms of  $\boldsymbol{\Sigma}$  using the observed marginal statistics. We are then ready to start EM.

#### 6.6.4.4 Example

As an example of this procedure in action, let us consider an imputation problem, where we have  $N_{\mathcal{D}} = 100$  10-dimensional data cases, which we assume to come from a Gaussian. We generate synthetic data where 50% of the observations are missing at random. First we fit the parameters

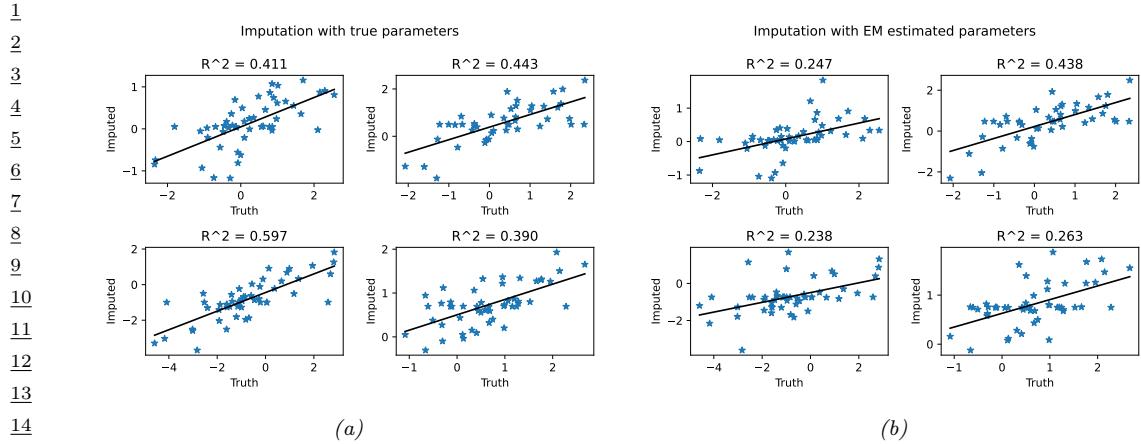


Figure 6.7: Illustration of data imputation using a multivariate Gaussian. (a) Scatter plot of true values vs imputed values using true parameters. (b) Same as (a), but using parameters estimated with EM. We just show the first four variables, for brevity. Generated by [gauss\\_imputation\\_em\\_demo.ipynb](#).

using EM. Call the resulting parameters  $\hat{\theta}$ . We can now use our model for predictions by computing  $\mathbb{E} [z_n | \mathbf{x}_n, \hat{\theta}]$ . Figure 6.7 indicates that the results obtained using the learned parameters are almost as good as with the true parameters. Not surprisingly, performance improves with more data, or as the fraction of missing data is reduced.

### 6.6.5 Example: robust linear regression using Student- $t$ likelihood

In this section, we discuss how to use EM to fit a linear regression model that uses the Student distribution for its likelihood, instead of the more common Gaussian distribution, in order to achieve robustness, as first proposed in [Zel76]. More precisely, the likelihood is given by

$$p(y|\mathbf{x}, \mathbf{w}, \sigma^2, \nu) = \mathcal{T}(y|\mathbf{w}^\top \mathbf{x}, \sigma^2, \nu) \quad (6.159)$$

At first blush it may not be apparent how to do this, since there is no missing data, and there are no hidden variables. However, it turns out that we can introduce “artificial” hidden variables to make the problem easier to solve; this is a common trick. The key insight is that we can represent the Student distribution as a Gaussian scale mixture, as we discussed in Section 28.2.3.1.

We can apply the GSM version of the Student distribution to our problem by associating a latent scale  $z_n \in \mathbb{R}_+$  with each example. The complete data log likelihood is therefore given by

$$\log p(\mathbf{y}, \mathbf{z} | \mathbf{X}, \mathbf{w}, \sigma^2, \nu) = \sum_n -\frac{1}{2} \log(2\pi z_n \sigma^2) - \frac{1}{2z_n \sigma^2} (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 \quad (6.160)$$

$$+ \left(\frac{\nu}{2} - 1\right) \log(z_n) - z_n \frac{\nu}{2} + \text{const} \quad (6.161)$$

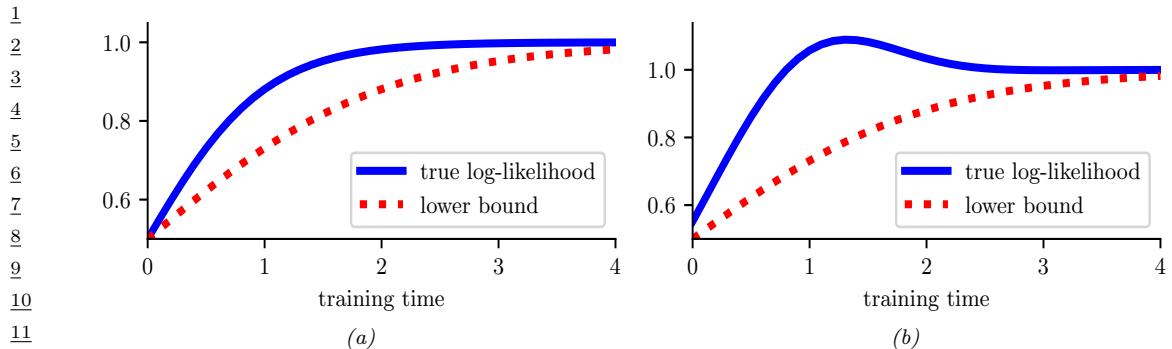


Figure 6.8: Illustration of possible behaviors of variational EM. (a) The lower bound increases at each iteration, and so does the likelihood. (b) The lower bound increases but the likelihood decreases. In this case, the algorithm is closing the gap between the approximate and true posterior. This can have a regularizing effect. Adapted from Figure 6 of [SJJ96]. Generated by `var_em_bound.ipynb`.

Ignoring terms not involving  $\mathbf{w}$ , and taking expectations, we have

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^t) = - \sum_n \frac{\lambda_n^t}{2\sigma^2} (y_n - \mathbf{w}^T \mathbf{x}_n)^2 \quad (6.162)$$

where  $\lambda_n^t \triangleq \mathbb{E}[1/z_n | y_n, \mathbf{x}_n, \mathbf{w}^t]$ . We recognize this as a weighted least squares objective, with weight  $\lambda_n^t$  per data point .

We now discuss how to compute these weights. Using the results from Section 2.2.3.4, one can show that

$$p(z_n | y_n, \mathbf{x}_n, \boldsymbol{\theta}) = \text{IG}\left(\frac{\nu + 1}{2}, \frac{\nu + \delta_n}{2}\right) \quad (6.163)$$

where  $\delta_n = \frac{(y_n - \mathbf{x}^T \mathbf{x}_n)^2}{\sigma^2}$  is the standardized residual. Hence

$$\lambda_n = \mathbb{E}[1/z_n] = \frac{\nu^t + 1}{\nu^t + \delta_n^t} \quad (6.164)$$

So if the residual  $\delta_n^t$  is large, the point will be given low weight  $\lambda_n^t$ , which makes intuitive sense, since it is probably an outlier.

### 6.6.6 Extensions to EM

There are many variations and extensions of the EM algorithm, as discussed in [MK97]. We summarize a few of these below.

#### 6.6.6.1 Variational EM

Suppose in the E step we pick  $q_n^* = \operatorname{argmin}_{q_n \in \mathcal{Q}} D_{\text{KL}}(q_n \| p(\mathbf{z}_n | \mathbf{x}_n, \boldsymbol{\theta}))$ . Because we are optimizing over the space of functions, this is called variational inference (see Section 10.1 for details). If the

1 family of distributions  $\mathcal{Q}$  is rich enough to contain the true posterior,  $q_n = p(\mathbf{z}_n | \mathbf{x}_n, \boldsymbol{\theta})$ , then we can  
2 make the KL be zero. But in general, we might choose a more restrictive class for computational  
3 reasons. For example, we might use  $q_n(\mathbf{z}_n) = \mathcal{N}(\mathbf{z}_n | \boldsymbol{\mu}_n, \text{diag}(\boldsymbol{\sigma}_n))$  even if the true posterior is  
4 correlated.  
5

6 The use of a restricted posterior family  $\mathcal{Q}$  inside the E step of EM is called **variational EM**  
7 [NH98a]. Unlike regular EM, variational EM is not guaranteed to increase the actual log likelihood  
8 itself (see Figure 6.8), but it does monotonically increase the variational lower bound. We can control  
9 the tightness of this lower bound by varying the variational family  $\mathcal{Q}$ ; in the limit in which  $q_n = p_n$ ,  
10 corresponding to exact inference, we recover the same behavior as regular EM. See Section 10.2.5 for  
11 further discussion.

1213

#### 14 6.6.6.2 Hard EM

15

16 Suppose we use a degenerate posterior approximation in the context of variational EM, corresponding  
17 to a point estimate,  $q(\mathbf{z} | \mathbf{x}_n) = \delta_{\hat{\mathbf{z}}_n}(\mathbf{z})$ , where  $\hat{\mathbf{z}}_n = \text{argmax}_{\mathbf{z}} p(\mathbf{z} | \mathbf{x}_n)$ . This is equivalent to **hard**  
18 **EM**, where we ignore uncertainty about  $\mathbf{z}_n$  in the E step.

19 The problem with this degenerate approach is that it is very prone to overfitting, since the number  
20 of latent variables is proportional to the number of datacases [WCS08].

2122

#### 23 6.6.6.3 Monte Carlo EM

24

25 Another approach to handling an intractable E step is to use a Monte Carlo approximation to the  
26 expected sufficient statistics. That is, we draw samples from the posterior,  $\mathbf{z}_n^s \sim p(\mathbf{z}_n | \mathbf{x}_n, \boldsymbol{\theta}^t)$ , and  
27 then compute the sufficient statistics for each completed vector,  $(\mathbf{x}_n, \mathbf{z}_n^s)$ , and then average the  
28 results. This is called **Monte Carlo EM** or **MCEM** [WT90; Nea12].

29 One way to draw samples is to use MCMC (see Chapter 12). However, if we have to wait for  
30 MCMC to converge inside each E step, the method becomes very slow. An alternative is to use  
31 stochastic approximation, and only perform “brief” sampling in the E step, followed by a partial  
32 parameter update. This is called **stochastic approximation EM** [DLM99] and tends to work  
33 better than MCEM.

343536

#### 37 6.6.6.4 Generalized EM

38 Sometimes we can perform the E step exactly, but we cannot perform the M step exactly. However,  
39 we can still monotonically increase the log likelihood by performing a “partial” M step, in which we  
40 merely increase the expected complete data log likelihood, rather than maximizing it. For example,  
41 we might follow a few gradient steps. This is called the **generalized EM** or **GEM** algorithm [MK07].  
42 (This is an unfortunate term, since there are many ways to generalize EM, but it is the standard  
43 terminology.) For example, [Lan95a] proposes to perform one Newton-Raphson step:  
44

45

$$\theta_{t+1} = \theta_t - \eta_t \mathbf{H}_t^{-1} \mathbf{g}_t \quad (6.165)$$

46

where  $0 < \eta_t \leq 1$  is the step size, and

$$\mathbf{g}_t = \frac{\partial}{\partial \boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}_t) |_{\boldsymbol{\theta}=\boldsymbol{\theta}_t} \quad (6.166)$$

$$\mathbf{H}_t = \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} Q(\boldsymbol{\theta}, \boldsymbol{\theta}_t) |_{\boldsymbol{\theta}=\boldsymbol{\theta}_t} \quad (6.167)$$

If  $\eta_t = 1$ , [Lan95a] calls this the **gradient EM algorithm**. However, it is possible to use a larger step size to speed up the algorithm, as in the **quasi-Newton EM algorithm** of [Lan95b]. This method also replaces the Hessian in Equation (6.167), which may not be negative definite (for non exponential family models), with a BFGS approximation. This ensures the overall algorithm is an ascent algorithm. Note, however, when the M step cannot be computed in closed form, EM loses some of its appeal over directly optimizing the marginal likelihood with a gradient based solver.

### 6.6.6.5 ECM algorithm

The **ECM** algorithm stands for “expectation conditional maximization”, and refers to optimizing the parameters in the M step sequentially, if they turn out to be dependent. The **ECME** algorithm, which stands for “ECM either” [LR95], is a variant of ECM in which we maximize the expected complete data log likelihood (the  $Q$  function) as usual, or the observed data log likelihood, during one or more of the conditional maximization steps. The latter can be much faster, since it ignores the results of the E step, and directly optimizes the objective of interest. A standard example of this is when fitting the Student T distribution. For fixed  $\nu$ , we can update  $\boldsymbol{\Sigma}$  as usual, but then to update  $\nu$ , we replace the standard update of the form  $\nu^{t+1} = \arg \max_\nu Q((\boldsymbol{\mu}^{t+1}, \boldsymbol{\Sigma}^{t+1}, \nu), \boldsymbol{\theta}^t)$  with  $\nu^{t+1} = \arg \max_\nu \log p(\mathcal{D} | \boldsymbol{\mu}^{t+1}, \boldsymbol{\Sigma}^{t+1}, \nu)$ . See [MK97] for more information.

### 6.6.6.6 Online EM

When dealing with large or streaming datasets, it is important to be able to learn online, as we discussed in Section 19.7.5. There are two main approaches to **online EM** in the literature. The first approach, known as **incremental EM** [NH98a], optimizes the lower bound  $Q(\boldsymbol{\theta}, q_1, \dots, q_N)$  one  $q_n$  at a time; however, this requires storing the expected sufficient statistics for each data case.

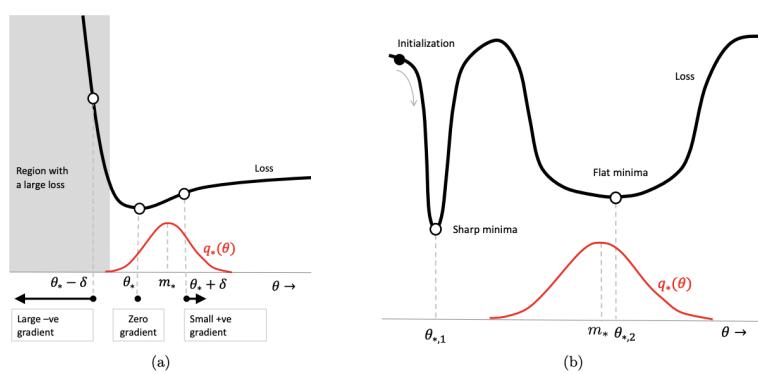
The second approach, known as **stepwise EM** [SI00; LK09; CM09], is based on stochastic gradient descent. This optimizes a local upper bound on  $\ell_n(\boldsymbol{\theta}) = \log p(\mathbf{x}_n | \boldsymbol{\theta})$  at each step. (See [Mai13; Mai15] for a more general discussion of stochastic and incremental bound optimization algorithms.)

## 6.7 The Bayesian learning rule

in this section, we discuss the “**Bayesian learning rule**” [KR21a], which provides a unified framework for deriving many standard (and non-standard) optimization and inference algorithms used in the ML community.

To motivate the BLR, recall the standard **empirical risk minimization** or **ERM** problem, which has the form  $\boldsymbol{\theta}_* = \operatorname{argmin}_{\boldsymbol{\theta}} \bar{\ell}(\boldsymbol{\theta})$ , where

$$\bar{\ell}(\boldsymbol{\theta}) = \sum_{n=1}^N \ell(\mathbf{y}_n, f_{\boldsymbol{\theta}}(\mathbf{x}_n)) + R(\boldsymbol{\theta}) \quad (6.168)$$



*Figure 6.9: Illustration of the robustness obtained by using a Bayesian approach to parameter estimation. (a) When the minimum  $\theta_*$  lies next to a “wall”, the Bayesian solution shifts away from the boundary to avoid large losses due to perturbations of the parameters. (b) The Bayesian solution prefers flat minima over sharp minima, to avoid large losses due to perturbations of the parameters. From Figure 1 of [KR21a]. Used with kind permission of Emtiyaz Khan.*

where  $f_\theta(\mathbf{x})$  is a prediction function,  $\ell(\mathbf{y}, \hat{\mathbf{y}})$  is a loss function, and  $R(\boldsymbol{\theta})$  is some kind of regularizer. Although the regularizer can prevent overfitting, the ERM method can still result in parameter estimates that are not robust. A better approach is to fit a *distribution* over possible parameter values,  $q(\boldsymbol{\theta})$ . If we minimize the expected loss, we will find parameter settings that will work well even if they are slightly perturbed, as illustrated in Figure 6.9, which helps with robustness and generalization. Of course, if the distribution  $q$  collapses to a single delta function, we will end up with the ERM solution. To prevent this, we add a penalty term, that measures the KL divergence from  $q(\boldsymbol{\theta})$  to some prior  $\pi_0(\boldsymbol{\theta}) \propto \exp(R(\boldsymbol{\theta}))$ . This gives rise to the following BLR objective:

$$\mathcal{L}(q) = \mathbb{E}_{q(\boldsymbol{\theta})} \left[ \sum_{n=1}^N \ell(\mathbf{y}_n, f_\theta(\mathbf{x}_n)) \right] + D_{\text{KL}}(q(\boldsymbol{\theta}) \parallel \pi_0(\boldsymbol{\theta})) \quad (6.169)$$

We can rewrite the KL term as

$$D_{\text{KL}}(q(\boldsymbol{\theta}) \parallel \pi_0(\boldsymbol{\theta})) = \mathbb{E}_{q(\boldsymbol{\theta})}[R(\boldsymbol{\theta})] - \mathbb{H}(q(\boldsymbol{\theta})) \quad (6.170)$$

and hence can rewrite the BLR objective as follows:

$$\mathcal{L}(q) = \mathbb{E}_{q(\boldsymbol{\theta})}[\bar{\ell}(\boldsymbol{\theta})] - \mathbb{H}(q(\boldsymbol{\theta})) \quad (6.171)$$

Below we show that different approximations to this objective recover a variety of different methods in the literature.

### 6.7.1 Deriving inference algorithms from BLR

In this section we show how to derive several different inference algorithms from BLR. (We discuss such algorithms in more detail in Chapter 10.)

---

### 6.7.1.1 Bayesian inference as optimization

The BLR objective includes standard exact Bayesian inference as a special case, as first shown in [Opt88]. To see this, let us assume the loss function is derived from a log-likelihood:

$$\ell(y, f_{\theta}(\mathbf{x})) = -\log p(\mathbf{y}|f_{\theta}(\mathbf{x})) \quad (6.172)$$

Let  $\mathcal{D} = \{(\mathbf{x}_n, \mathbf{y}_n) : n = 1 : N\}$  be the data we condition on. The Bayesian posterior can be written as

$$p(\boldsymbol{\theta}|\mathcal{D}) = \frac{1}{Z(\mathcal{D})} \pi_0(\boldsymbol{\theta}) \prod_{n=1}^N p(\mathbf{y}_n|f_{\theta}(\mathbf{x}_n)) \quad (6.173)$$

This can be derived by minimizing the BLR, since

$$\mathcal{L}(q) = -\mathbb{E}_{q(\boldsymbol{\theta})} \left[ \sum_{n=1}^N \log p(\mathbf{y}_n|f_{\theta}(\mathbf{x}_n)) \right] + D_{\text{KL}}(q(\boldsymbol{\theta}) \parallel \pi_0(\boldsymbol{\theta})) \quad (6.174)$$

$$= \mathbb{E}_{q(\boldsymbol{\theta})} \left[ \log \frac{q(\boldsymbol{\theta})}{\frac{\pi_0(\boldsymbol{\theta})}{Z(\mathcal{D})} \prod_{n=1}^N p(\mathbf{y}_n|f_{\theta}(\mathbf{x}_n))} \right] - \log Z(\mathcal{D}) \quad (6.175)$$

$$= D_{\text{KL}}(q(\boldsymbol{\theta}) \parallel p(\boldsymbol{\theta}|\mathcal{D})) - \log Z(\mathcal{D}) \quad (6.176)$$

Since  $Z(\mathcal{D})$  is a constant, we can minimize the loss by setting  $q(\boldsymbol{\theta}) = p(\boldsymbol{\theta}|\mathcal{D})$ .

Of course, we can use other kinds of loss, not just log likelihoods. This results in a framework known as **generalized Bayesian inference** [BHW16; KJD19; KJD21]. See Section 14.1.3 for more discussion.

### 6.7.1.2 Optimization of BLR using natural gradient descent

In general, we cannot compute the exact posterior  $q(\boldsymbol{\theta}) = p(\boldsymbol{\theta}|\mathcal{D})$ , so we seek an approximation. We will assume that  $q(\boldsymbol{\theta})$  is an exponential family distribution, such as a multivariate Gaussian, where the mean represents the standard point estimate of  $\boldsymbol{\theta}$  (as in ERM), and the covariance represents our uncertainty (as in Bayes). Hence  $q$  can be written as follows:

$$q(\boldsymbol{\theta}) = h(\boldsymbol{\theta}) \exp[\boldsymbol{\lambda}^T \mathcal{T}(\boldsymbol{\theta}) - A(\boldsymbol{\lambda})] \quad (6.177)$$

where  $\boldsymbol{\lambda}$  are the natural parameters,  $\mathcal{T}(\boldsymbol{\theta})$  are the sufficient statistics,  $A(\boldsymbol{\lambda})$  is the log partition function, and  $h(\boldsymbol{\theta})$  is the base measure, which is usually a constant. The BLR loss becomes

$$\mathcal{L}(\boldsymbol{\lambda}) = \mathbb{E}_{q_{\boldsymbol{\lambda}}(\boldsymbol{\theta})} [\bar{\ell}(\boldsymbol{\theta})] - \mathbb{H}(q_{\boldsymbol{\lambda}}(\boldsymbol{\theta})) \quad (6.178)$$

We can optimize this using natural gradient descent (Section 6.4). The update becomes

$$\boldsymbol{\lambda}_{t+1} = \boldsymbol{\lambda}_t - \eta_t \tilde{\nabla}_{\boldsymbol{\lambda}} \left[ \mathbb{E}_{q_{\boldsymbol{\lambda}_t}(\boldsymbol{\theta})} [\bar{\ell}(\boldsymbol{\theta})] - \mathbb{H}(q_{\boldsymbol{\lambda}_t}) \right] \quad (6.179)$$

where  $\tilde{\nabla}_{\boldsymbol{\lambda}}$  denotes the natural gradient. We discuss how to compute these natural gradients in Section 6.4.5. In particular, we can convert it to regular gradients wrt the moment parameters  $\boldsymbol{\mu}_t = \boldsymbol{\mu}(\boldsymbol{\lambda}_t)$ . This gives

$$\boldsymbol{\lambda}_{t+1} = \boldsymbol{\lambda}_t - \eta_t \nabla_{\boldsymbol{\mu}} \mathbb{E}_{q_{\boldsymbol{\mu}_t}(\boldsymbol{\theta})} [\bar{\ell}(\boldsymbol{\theta})] + \eta_t \nabla_{\boldsymbol{\mu}} \mathbb{H}(q_{\boldsymbol{\mu}_t}) \quad (6.180)$$

1 From Equation (6.84) we have  
2

$$\nabla_{\boldsymbol{\mu}} \mathbb{H}(q) = -\boldsymbol{\lambda} - \nabla_{\boldsymbol{\mu}} \mathbb{E}_{q_{\boldsymbol{\mu}}(\boldsymbol{\theta})} [\log h(\boldsymbol{\theta})] \quad (6.181)$$

3 Hence the update becomes  
4

$$\boldsymbol{\lambda}_{t+1} = \boldsymbol{\lambda}_t - \eta_t \nabla_{\boldsymbol{\mu}} \mathbb{E}_{q_{\boldsymbol{\mu}_t}(\boldsymbol{\theta})} [\bar{\ell}(\boldsymbol{\theta})] - \eta_t \boldsymbol{\lambda}_t - \eta_t \nabla_{\boldsymbol{\mu}} \mathbb{E}_{q_{\boldsymbol{\mu}}(\boldsymbol{\theta})} [\log h(\boldsymbol{\theta})] \quad (6.182)$$

$$= (1 - \eta_t) \boldsymbol{\lambda}_t - \eta_t \nabla_{\boldsymbol{\mu}} \mathbb{E}_{q_{\boldsymbol{\mu}}(\boldsymbol{\theta})} [\bar{\ell}(\boldsymbol{\theta}) + \log h(\boldsymbol{\theta})] \quad (6.183)$$

5 For distributions  $q$  with constant base measure  $h(\boldsymbol{\theta})$ , this simplifies to  
6

$$\boldsymbol{\lambda}_{t+1} = (1 - \eta_t) \boldsymbol{\lambda}_t - \eta_t \nabla_{\boldsymbol{\mu}} \mathbb{E}_{q_{\boldsymbol{\mu}}(\boldsymbol{\theta})} [\bar{\ell}(\boldsymbol{\theta})] \quad (6.184)$$

7 Hence at the fixed point we have  
8

$$\boldsymbol{\lambda}_* = (1 - \eta) \boldsymbol{\lambda}_* - \eta \nabla_{\boldsymbol{\mu}} \mathbb{E}_{q_{\boldsymbol{\mu}}(\boldsymbol{\theta})} [\bar{\ell}(\boldsymbol{\theta})] \quad (6.185)$$

$$\boldsymbol{\lambda}_* = \nabla_{\boldsymbol{\mu}} \mathbb{E}_{q_{\boldsymbol{\mu}}(\boldsymbol{\theta})} [-\bar{\ell}(\boldsymbol{\theta})] = \tilde{\nabla}_{\boldsymbol{\lambda}} \mathbb{E}_{q_{\boldsymbol{\lambda}}(\boldsymbol{\theta})} [-\bar{\ell}(\boldsymbol{\theta})] \quad (6.186)$$

### 16 6.7.1.3 Conjugate variational inference 17

18 In Section 7.3 we showed how to do exact inference in conjugate models. We can derive Equation (7.11)  
19 from the BLR by using the fixed point condition in Equation (6.186) to write  
20

$$\boldsymbol{\lambda}_* = \nabla_{\boldsymbol{\mu}} \mathbb{E}_{q_*} [-\bar{\ell}(\boldsymbol{\theta})] = \boldsymbol{\lambda}_0 + \sum_{i=1}^N \underbrace{\nabla_{\boldsymbol{\mu}} \mathbb{E}_{q_*} [\log p(\mathbf{y}_i | \boldsymbol{\theta})]}_{\tilde{\lambda}_i(\mathbf{y}_i)} \quad (6.187)$$

21 where  $\tilde{\lambda}_i(\mathbf{y}_i)$  are the sufficient statistics for the  $i$ 'th likelihood term.  
22

23 For models where the joint distribution over the latents factorizes (using a graphical model), we  
24 can further decompose this update into a series of local terms. This gives rise to the variational  
25 message passing scheme discussed in Section 10.2.7.  
26

### 27 6.7.1.4 Partially conjugate variational inference 28

29 In Supplementary Section 10.3.1, we discuss CVI, which performs variational inference for partially  
30 conjugate models, using gradient updates for the non-conjugate parts, and exact Bayesian inference  
31 for the conjugate parts.  
32

33

## 34 6.7.2 Deriving optimization algorithms from BLR 35

36 In this section we show how to derive several different optimization algorithms from BLR. Recall  
37 that in BLR, instead of directly minimizing the loss  
38

$$\bar{\ell}(\boldsymbol{\theta}) = \sum_{n=1}^N \ell(\mathbf{y}_n, f_{\boldsymbol{\theta}}(\mathbf{x}_n)) + R(\boldsymbol{\theta}) \quad (6.188)$$

39 we will instead minimize  
40

$$\mathcal{L}(\boldsymbol{\lambda}) = \mathbb{E}_{q(\boldsymbol{\theta}|\boldsymbol{\lambda})} [\bar{\ell}(\boldsymbol{\theta})] - \mathbb{H}(q(\boldsymbol{\theta}|\boldsymbol{\lambda})) \quad (6.189)$$

41 Below we show that different approximations to this objective recover a variety of different optimization  
42 methods that are used in the literature. (We discuss such algorithms in more detail in Chapter 6.)  
43

1

### 6.7.2.1 Gradient descent

2  
3 In this section, we show how to derive gradient descent as a special case of BLR. We use as our  
4 approximate posterior  $q(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta} | \mathbf{m}, \mathbf{I})$ . In this case the natural and moment parameters are equal,  
5  $\boldsymbol{\mu} = \boldsymbol{\lambda} = \mathbf{m}$ . The base measure satisfies the following (from Equation (2.180)):

6

$$2 \log h(\boldsymbol{\theta}) = -D \log(2\pi) - \boldsymbol{\theta}^\top \boldsymbol{\theta} \quad (6.190)$$

7 Hence

8

$$\nabla_{\boldsymbol{\mu}} \mathbb{E}_q [\log h(\boldsymbol{\theta})] = \nabla_{\boldsymbol{\mu}} (-D \log(2\pi) - \boldsymbol{\mu}^\top \boldsymbol{\mu} - D) = -\boldsymbol{\mu} = -\boldsymbol{\lambda} = -\mathbf{m} \quad (6.191)$$

9 Thus from Equation (6.183) the BLR update becomes

10

$$\mathbf{m}_{t+1} = (1 - \eta_t) \mathbf{m}_t + \eta_t \mathbf{m}_t - \eta_t \nabla_{\mathbf{m}} \mathbb{E}_{q_{\mathbf{m}}(\boldsymbol{\theta})} [\bar{\ell}(\boldsymbol{\theta})] \quad (6.192)$$

11 We can remove the expectation using the first order delta method (Section 6.5.5):

12

$$\nabla_{\mathbf{m}} \mathbb{E}_{q_{\mathbf{m}}(\boldsymbol{\theta})} [\bar{\ell}(\boldsymbol{\theta})] \approx \nabla_{\boldsymbol{\theta}} \bar{\ell}(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\mathbf{m}} \quad (6.193)$$

13 Putting these together gives the gradient descent update:

14

$$\mathbf{m}_{t+1} = \mathbf{m}_t - \eta_t \nabla_{\boldsymbol{\theta}} \bar{\ell}(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\mathbf{m}_t} \quad (6.194)$$

15

### 6.7.2.2 Newton's method

16 In this section, we show how to derive Newton's second order optimization method as a special case  
17 of BLR, as first shown in [Kha+18].

18 Suppose we assume  $q(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta} | \mathbf{m}, \mathbf{S}^{-1})$ . The natural parameters are

19

$$\boldsymbol{\lambda}^{(1)} = \mathbf{S}\mathbf{m}, \quad \boldsymbol{\lambda}^{(2)} = -\frac{1}{2}\mathbf{S} \quad (6.195)$$

20 The mean (moment) parameters are

21

$$\boldsymbol{\mu}^{(1)} = \mathbf{m}, \quad \boldsymbol{\mu}^{(2)} = \mathbf{S}^{-1} + \mathbf{m}\mathbf{m}^\top \quad (6.196)$$

22 Since the base measure is constant (see Equation (2.193)), from Equation (6.184) we have

23

$$\mathbf{S}_{t+1} \mathbf{m}_{t+1} = (1 - \eta_t) \mathbf{S}_t \mathbf{m}_t - \eta_t \nabla_{\boldsymbol{\mu}^{(1)}} \mathbb{E}_{q_{\boldsymbol{\mu}}(\boldsymbol{\theta})} [\bar{\ell}(\boldsymbol{\theta})] \quad (6.197)$$

24

$$\mathbf{S}_{t+1} = (1 - \eta_t) \mathbf{S}_t + 2\eta_t \nabla_{\boldsymbol{\mu}^{(2)}} \mathbb{E}_{q_{\boldsymbol{\mu}}(\boldsymbol{\theta})} [\bar{\ell}(\boldsymbol{\theta})] \quad (6.198)$$

25 In Section 6.4.5.1 we show that

26

$$\nabla_{\boldsymbol{\mu}^{(1)}} \mathbb{E}_{q(\boldsymbol{\theta})} [\bar{\ell}(\boldsymbol{\theta})] = \mathbb{E}_{q(\boldsymbol{\theta})} [\nabla_{\boldsymbol{\theta}} \bar{\ell}(\boldsymbol{\theta})] - \mathbb{E}_{q(\boldsymbol{\theta})} [\nabla_{\boldsymbol{\theta}}^2 \bar{\ell}(\boldsymbol{\theta})] \mathbf{m} \quad (6.199)$$

27

$$\nabla_{\boldsymbol{\mu}^{(2)}} \mathbb{E}_{q(\boldsymbol{\theta})} [\bar{\ell}(\boldsymbol{\theta})] = \frac{1}{2} \mathbb{E}_{q(\boldsymbol{\theta})} [\nabla_{\boldsymbol{\theta}}^2 \bar{\ell}(\boldsymbol{\theta})] \quad (6.200)$$

28 Hence the update for the precision matrix becomes

29

$$\mathbf{S}_{t+1} = (1 - \eta_t) \mathbf{S}_t + \eta_t \mathbb{E}_{q_t} [\nabla_{\boldsymbol{\theta}}^2 \bar{\ell}(\boldsymbol{\theta})] \quad (6.201)$$

1 For the precision weighted mean, we have  
2

$$\underline{3} \quad \mathbf{S}_{t+1} \mathbf{m}_{t+1} = (1 - \eta_t) \mathbf{S}_t \mathbf{m}_t - \eta_t \mathbb{E}_{q_t} [\nabla_{\boldsymbol{\theta}} \bar{\ell}(\boldsymbol{\theta})] + \eta_t \mathbb{E}_{q_t} [\nabla_{\boldsymbol{\theta}}^2 \bar{\ell}(\boldsymbol{\theta})] \mathbf{m}_t \quad (6.202)$$

$$\underline{4} \quad = \mathbf{S}_{t+1} \mathbf{m}_t - \eta_t \mathbb{E}_{q_t} [\nabla_{\boldsymbol{\theta}} \bar{\ell}(\boldsymbol{\theta})] \quad (6.203)$$

5 Hence  
6

$$\underline{7} \quad \mathbf{m}_{t+1} = \mathbf{m}_t - \eta_t \mathbf{S}_{t+1}^{-1} \mathbb{E}_{q_t} [\nabla_{\boldsymbol{\theta}} \bar{\ell}(\boldsymbol{\theta})] \quad (6.204)$$

8 We can recover Newton's method in three steps. First set the learning rate to  $\eta_t = 1$ , based on  
9 an assumption that the objective is convex. Second, treat the iterate as  $\mathbf{m}_t = \boldsymbol{\theta}_t$ . Third, apply the  
10 delta method to get

$$\underline{11} \quad \mathbf{S}_{t+1} = \mathbb{E}_{q_t} [\nabla_{\boldsymbol{\theta}}^2 \bar{\ell}(\boldsymbol{\theta})] \approx \nabla_{\boldsymbol{\theta}}^2 \bar{\ell}(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\mathbf{m}_t} \quad (6.205)$$

12 and  
13

$$\underline{14} \quad \mathbb{E}_q [\nabla_{\boldsymbol{\theta}} \bar{\ell}(\boldsymbol{\theta})] \approx \nabla_{\boldsymbol{\theta}} \bar{\ell}(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\mathbf{m}_t} \quad (6.206)$$

15 This gives Newton's update:  
16

$$\underline{17} \quad \mathbf{m}_{t+1} = \mathbf{m}_t - [\nabla_{\boldsymbol{\theta}}^2 \bar{\ell}(\boldsymbol{\theta})]^{-1} [\nabla_{\boldsymbol{\theta}} \bar{\ell}(\boldsymbol{\theta})] \quad (6.207)$$

### 18 6.7.2.3 Variational online Gauss-Newton

19 In this section, we describe the **Variational Online Gauss-Newton** or **VOGN** method of [Kha+18].  
20 This is an approximate second order optimization method that can be derived from the BLR in  
21 several steps, as we show below.

22 First, we use a diagonal Gaussian approximation to the posterior,  $q_t(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|\boldsymbol{\theta}_t, \mathbf{S}_t^{-1})$ , where  
23  $\mathbf{S}_t = \text{diag}(\mathbf{s}_t)$  is a vector of precisions. Following Section 6.7.2.2, we get the following updates:  
24

$$\underline{25} \quad \boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_t \frac{1}{\mathbf{s}_{t+1}} \odot \mathbb{E}_{q_t} [\nabla_{\boldsymbol{\theta}} \bar{\ell}(\boldsymbol{\theta})] \quad (6.208)$$

$$\underline{26} \quad \mathbf{s}_{t+1} = (1 - \eta_t) \mathbf{s}_t + \eta_t \mathbb{E}_{q_t} [\text{diag}(\nabla_{\boldsymbol{\theta}}^2 \bar{\ell}(\boldsymbol{\theta}))] \quad (6.209)$$

27 where  $\odot$  is elementwise multiplication, and the division by  $\mathbf{s}_{t+1}$  is also elementwise.  
28

29 Second, we use the delta approximation to replace expectations by plugging in the mean. Third  
30 we use a minibatch approximation to the gradient and diagonal Hessian:  
31

$$\underline{32} \quad \hat{\nabla}_{\boldsymbol{\theta}} \bar{\ell}(\boldsymbol{\theta}) = \frac{N}{M} \sum_{i \in \mathcal{M}} \nabla_{\boldsymbol{\theta}} \ell(y_i, f_{\boldsymbol{\theta}}(\mathbf{x}_i)) + \nabla_{\boldsymbol{\theta}} R(\boldsymbol{\theta}) \quad (6.210)$$

$$\underline{33} \quad \hat{\nabla}_{\boldsymbol{\theta}}^2 \bar{\ell}(\boldsymbol{\theta}) = \frac{N}{M} \sum_{i \in \mathcal{M}} \nabla_{\boldsymbol{\theta}}^2 \ell(y_i, f_{\boldsymbol{\theta}}(\mathbf{x}_i)) + \nabla_{\boldsymbol{\theta}}^2 R(\boldsymbol{\theta}) \quad (6.211)$$

34 where  $M$  is the minibatch size.  
35

For some non-convex problems, such as DNNs, the Hessian may be not be positive definite, so we can get better results using a Gauss-Newton approximation, based on the squared gradients instead of the Hessian:

$$\hat{\nabla}_{\theta_j}^2 \bar{\ell}(\boldsymbol{\theta}) \approx \frac{N}{M} \sum_{i \in \mathcal{M}} [\nabla_{\theta_j} \ell(y_i, f_{\boldsymbol{\theta}}(\mathbf{x}_i))]^2 + \nabla_{\theta_j}^2 R(\boldsymbol{\theta}) \quad (6.212)$$

This is also faster to compute.

Putting all this together gives rise to the **Online Gauss-Newton** or **OGN** method of [Osa+19a].

If we drop the delta approximation, and work with expectations. we get the **Variational Online**

**Gauss-Newton** or **VOGN** method of [Kha+18]. We can approximate the expectations by sampling.

In particular, VOGN uses the following **weight perturbation** method

$$\mathbb{E}_{q_t} [\hat{\nabla}_{\boldsymbol{\theta}} \bar{\ell}(\boldsymbol{\theta})] \approx \hat{\nabla}_{\boldsymbol{\theta}} \bar{\ell}(\boldsymbol{\theta}_t + \boldsymbol{\epsilon}_t) \quad (6.213)$$

where  $\boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \text{diag}(\mathbf{s}_t))$ . It also also possible to approximate the Fisher information matrix directly;

this results in the **Variational Online Generalized Gauss-Newton** or **VOGNN** method of

[Osa+19a].

#### 6.7.2.4 Adaptive learning rate SGD

In this section, we show how to derive an update rule which is very similar to the **RMSprop** [Hin14] method, which is widely used in deep learning. The approach we take is similar to that VOGN in Section 6.7.2.3. We use the same diagonal Gaussian approximation,  $q_t(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta} | \boldsymbol{\theta}_t, \mathbf{S}_t^{-1})$ , where  $\mathbf{S}_t = \text{diag}(\mathbf{s}_t)$  is a vector of precisions. We then use the delta method to eliminate expectations:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_t \frac{1}{\mathbf{s}_{t+1}} \odot \nabla_{\boldsymbol{\theta}} \bar{\ell}(\boldsymbol{\theta}_t) \quad (6.214)$$

$$\mathbf{s}_{t+1} = (1 - \eta_t) \mathbf{s}_t + \eta_t \text{diag}(\nabla_{\boldsymbol{\theta}}^2 \bar{\ell}(\boldsymbol{\theta}_t)) \quad (6.215)$$

where  $\odot$  is elementwise multiplication. If we allow for different learning rates we get

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \alpha_t \frac{1}{\mathbf{s}_{t+1}} \odot \nabla_{\boldsymbol{\theta}} \bar{\ell}(\boldsymbol{\theta}_t) \quad (6.216)$$

$$\mathbf{s}_{t+1} = (1 - \beta_t) \mathbf{s}_t + \beta_t \text{diag}(\hat{\nabla}_{\boldsymbol{\theta}}^2 \bar{\ell}(\boldsymbol{\theta}_t)) \quad (6.217)$$

Now suppose we replace the diagonal Hessian approximation with the sum of the squares per-sample gradients:

$$\text{diag}(\nabla_{\boldsymbol{\theta}}^2 \bar{\ell}(\boldsymbol{\theta}_t)) \approx \hat{\nabla} \bar{\ell}(\boldsymbol{\theta}_t) \odot \hat{\nabla} \bar{\ell}(\boldsymbol{\theta}_t) \quad (6.218)$$

If we also change some scaling factors we can get the RMSprop updates:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \alpha \frac{1}{\sqrt{\mathbf{v}_{t+1} + c\mathbf{1}}} \odot \hat{\nabla}_{\boldsymbol{\theta}} \bar{\ell}(\boldsymbol{\theta}_t) \quad (6.219)$$

$$\mathbf{v}_{t+1} = (1 - \beta) \mathbf{v}_t + \beta [\hat{\nabla} \bar{\ell}(\boldsymbol{\theta}_t) \odot \hat{\nabla} \bar{\ell}(\boldsymbol{\theta}_t)] \quad (6.220)$$

This allows us to use standard deep learning optimizers to get a Gaussian approximation to the posterior for the parameters [Osa+19a].

It is also possible to derive the Adam optimizer [KB15] from BLR by adding a momentum term to RMSprop. See [KR21a; Ait18] for details.

1 **6.7.3 Variational optimization**

3 Consider an objective defined in terms of discrete variables. Such objectives are not differentiable and  
4 so are hard to optimize. One advantage of BLR is that it optimizes the parameters of a probability  
5 distribution, and such expected loss objectives are usually differentiable and smooth. This is called  
6 “**variational optimization**” [Bar17], since we are optimizing over a probability distribution.

7 For example, consider the case of a **binary neural network** where  $\theta_d \in \{0, 1\}$  indicates if  
8 weight  $d$  is used or not, we can optimize over the parameters of a Bernoulli distribution,  $q(\boldsymbol{\theta}|\boldsymbol{\lambda}) =$   
9  $\prod_{d=1}^D \text{Ber}(\theta_d|p_d)$ , where  $p_d \in [0, 1]$  and  $\lambda_d = 1/2 \log(p_d/(1-p_d))$  is the log odds. This is the basis of  
10 the **BayesBiNN** approach [MBK20].

11 If we ignore the entropy and regularizer term, we get the following simplified objective:

$$\mathcal{L}(\boldsymbol{\lambda}) = \int \bar{\ell}(\boldsymbol{\theta}) q(\boldsymbol{\theta}|\boldsymbol{\lambda}) d\boldsymbol{\theta} \quad (6.221)$$

16 This method has various names: **stochastic relaxation** [SB12; SB13; MMP13], **stochastic ap-**  
17 **proximation** [HHC12; Hu+12], etc. It is closely related to **evolutionary strategies**, which we  
18 discuss in Section 6.9.6.

19 In the case of functions with continuous domains, we can use a Gaussian for  $q(\boldsymbol{\theta}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . The  
20 resulting integral in Equation (6.221) can then sometimes be solved in closed form, as explained in  
21 [Mob16]. By starting with a broad variance, and gradually reducing it, we hope the method can  
22 avoid poor local optima, similar to simulated annealing (Section 12.9.1). However, we generally get  
23 better results by including the entropy term, because then we can automatically learn to adapt the  
24 variance. In addition, we can often work with natural gradients, which results in faster convergence.

26 **6.8 Bayesian optimization**

29 In this section, we discuss **Bayesian optimization** or **BayesOpt**, which is a model-based approach  
30 to black-box optimization, designed for the case where the objective function  $f : \mathcal{X} \rightarrow \mathbb{R}$  is expensive  
31 to evaluate (e.g., if it requires running a simulation, or training and testing a particular neural net  
32 architecture).

33 Since the true function  $f$  is expensive to evaluate, we want to make as few function calls (i.e., make as  
34 few **queries**  $\mathbf{x}$  to the **oracle**  $f$ ) as possible. This suggests that we should build a **surrogate function**  
35 (also called a **response surface model**) based on the data collected so far,  $\mathcal{D}_n = \{(\mathbf{x}_i, y_i) : i = 1 : n\}$ ,  
36 which we can use to decide which point to query next. There is an inherent tradeoff between picking  
37 the point  $\mathbf{x}$  where we think  $f(\mathbf{x})$  is large (we follow the convention in the literature and assume  
38 we are trying to maximize  $f$ ), and picking points where we are uncertain about  $f(\mathbf{x})$  but where  
39 observing the function value might help us improve the surrogate model. This is another instance of  
40 the exploration-exploitation dilemma.

41 In the special case where the domain we are optimizing over is finite, so  $\mathcal{X} = \{1, \dots, A\}$ , the  
42 BayesOpt problem becomes similar to the **best arm identification** problem in the bandit literature  
43 (Section 34.4). An important difference is that in bandits, we care about the cost of every action we  
44 take, whereas in optimization, we usually only care about the cost of the final solution we find. In  
45 other words, in bandits, we want to minimize cumulative regret, whereas in optimization we want to  
46 minimize simple or final regret.

47

Another related topic is **active learning**. Here the goal is to identify the whole function  $f$  with as few queries as possible, whereas in BayesOpt, the goal is just to identify the maximum of the function.

Bayesian optimization is a large topic, and we only give a brief overview below. For more details, see e.g., [Sha+16; Fra18; Gar22]. (See also <https://distill.pub/2020/bayesian-optimization/> for an interactive tutorial.)

### 6.8.1 Sequential model-based optimization

BayesOpt is an instance of a strategy known as sequential model-based optimization (**SMBO**) [HHLB11]. In this approach, we alternate between querying the function at a point, and updating our estimate of the surrogate based on the new data. More precisely, at each iteration  $n$ , we have a labeled dataset  $\mathcal{D}_n = \{(\mathbf{x}_i, y_i) : i = 1 : n\}$ , which records points  $\mathbf{x}_i$  that we have queried, and the corresponding function values  $y_i = f(\mathbf{x}_i) + \epsilon_i$ , where  $\epsilon_i$  is an optional noise term. We use this data to estimate a probability distribution over the true function  $f$ ; we will denote this by  $p(f|\mathcal{D}_n)$ . We then choose the next point to query  $\mathbf{x}_{n+1}$  using an **acquisition function**  $\alpha(\mathbf{x}; \mathcal{D}_n)$ , which computes the expected utility of querying  $\mathbf{x}$ . (We discuss acquisition functions in Section 6.8.3). After we observe  $y_{n+1} = f(\mathbf{x}_{n+1}) + \epsilon_{n+1}$ , we update our beliefs about the function, and repeat. See Algorithm 5 for some pseudocode.

---

**Algorithm 5:** Bayesian optimization

---

```

1 Collect initial dataset  $\mathcal{D}_0 = \{(\mathbf{x}_i, y_i) : \}$  from random queries  $\mathbf{x}_i$ 
2 Initialize model  $p(f|\mathcal{D}_0)$ 
3 for  $n = 1, 2, \dots$  until convergence do
4   Choose next query point  $\mathbf{x}_{n+1} = \operatorname{argmax}_{\mathbf{x} \in \mathcal{X}} \alpha(\mathbf{x}; \mathcal{D}_n)$ 
5   Measure function value,  $y_{n+1} = f(\mathbf{x}_{n+1}) + \epsilon_n$ 
6   Augment dataset,  $\mathcal{D}_{n+1} = \{\mathcal{D}_n, (\mathbf{x}_{n+1}, y_{n+1})\}$ 
7   Update model,  $p(f|\mathcal{D}_{n+1}) \propto p(f|\mathcal{D}_n)p(y_{n+1}|\mathbf{x}_{n+1}, f)$ 
```

---

This method is illustrated in Figure 6.10. The goal is to find the global optimum of the solid black curve. In the first row, we show the 2 previously queried points,  $x_1$  and  $x_2$ , and their corresponding function values.  $y_1 = f(x_1)$  and  $y_2 = f(x_2)$ . Our uncertainty about the value of  $f$  at those locations is 0 (if we assume no observation noise), as illustrated by the posterior credible interval (shaded blue area) becoming “pinched”. Consequently the acquisition function (shown in green at the bottom) also has value 0 at those previously queried points. The red triangle represents the maximum of the acquisition function, which becomes our next query,  $x_3$ . In the second row, we show the result of observing  $y_3 = f(x_3)$ ; this further reduces our uncertainty about the shape of the function. In the third row, we show the result of observing  $y_4 = f(x_4)$ . This process repeats until we run out of time, or until we are confident there are no better unexplored points to query.

The two main “ingredients” that we need to provide to a BayesOpt algorithm are (1) a way to represent and update the posterior surrogate  $p(f|\mathcal{D}_n)$ , and (2) a way to define and optimize the acquisition function  $\alpha(\mathbf{x}; \mathcal{D}_n)$ . We discuss both of these topics below.

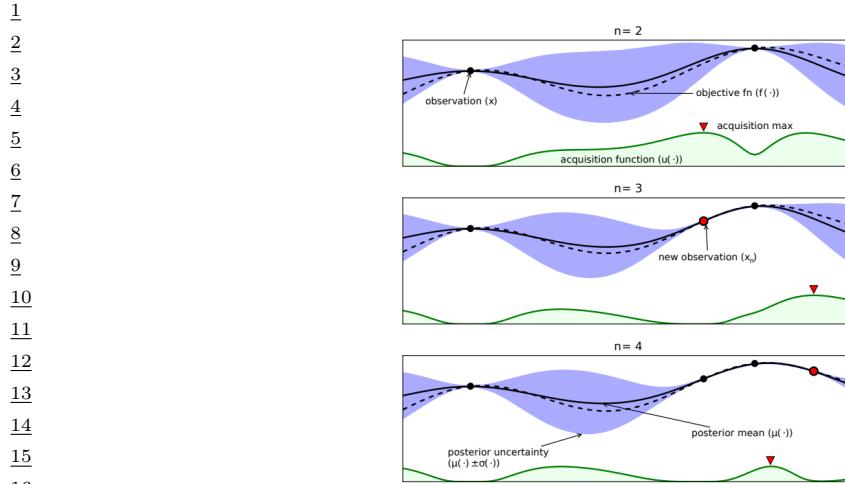


Figure 6.10: Illustration of sequential Bayesian optimization over three iterations. The rows correspond to a training set of size  $t = 2, 3, 4$ . The solid black line is the true, but unknown, function  $f(x)$ . The dotted black line is the posterior mean,  $\mu(x)$ . The shaded blue intervals are the 95% credible interval derived from  $\mu(x)$  and  $\sigma(x)$ . The solid black dots correspond to points whose function value has already been computed, i.e.,  $x_n$  for which  $f(x_n)$  is known. The green curve at the bottom is the acquisition function. The red dot is the proposed next point to query, which is the maximum of the acquisition function. From Figure 1 of [Sha+16]. Used with kind permission of Nando de Freitas.

## 6.8.2 Surrogate functions

In this section, we discuss ways to represent and update the posterior over functions,  $p(f|\mathcal{D}_n)$ .

### 6.8.2.1 Gaussian processes

In BayesOpt, it is very common to use a Gaussian process or GP for our surrogate. GPs are explained in detail in Chapter 18, but the basic idea is that they represent  $p(f(\mathbf{x})|\mathcal{D}_n)$  as a Gaussian,  $p(f(\mathbf{x})|\mathcal{D}_n) = \mathcal{N}(f|\mu_n(\mathbf{x}), \sigma_n^2(\mathbf{x}))$ , where  $\mu_n(\mathbf{x})$  and  $\sigma_n(\mathbf{x})$  are functions that can be derived from the training data  $\mathcal{D}_n = \{(\mathbf{x}_i, y_i) : i = 1 : n\}$  using a simple closed-form equation. The GP requires specifying a kernel function  $\mathcal{K}_{\theta}(\mathbf{x}, \mathbf{x}')$ , which measures similarities between input points  $\mathbf{x}, \mathbf{x}'$ . The intuition is that if two inputs are similar, so  $\mathcal{K}_{\theta}(\mathbf{x}, \mathbf{x}')$  is large, then the corresponding function values are also likely to be similar, so  $f(\mathbf{x})$  and  $f(\mathbf{x}')$  should be positively correlated. This allows us to interpolate the function between the labeled training points; in some cases, it also lets us extrapolate beyond them.

GPs work well when we have little training data, and they support closed form Bayesian updating. However, exact updating takes  $O(N^3)$  for  $N$  samples, which becomes too slow if we perform many function evaluations. There are various methods (Section 18.5.3) for reducing this to  $O(NM^2)$  time, where  $M$  is a parameter we choose, but this sacrifices some of the accuracy.

In addition, the performance of GPs depends heavily on having a good kernel. We can estimate the kernel parameters  $\theta$  by maximizing the marginal likelihood, as discussed in Section 18.6.1. However,

since the sample size is small (by assumption), we can often get better performance by marginalizing out  $\theta$  using approximate Bayesian inference methods, as discussed in Section 18.6.2. See e.g., [WF16] for further details.

### 6.8.2.2 Bayesian neural networks

A natural alternative to GPs is to use a parametric model. If we use linear regression, we can efficiently perform exact Bayesian inference, as shown in Section 15.2. If we use a nonlinear model, such as a DNN, we need to use approximate inference methods. We discuss Bayesian neural networks in detail in Chapter 17. For their application to BayesOpt, see e.g. [Spr+16].

### 6.8.2.3 Other models

We are free to use other forms of regression model. [HHLB11] use an ensemble of random forests; such models can easily handle conditional parameter spaces, as we discuss in Section 6.8.4.2, although bootstrapping (which is needed to get uncertainty estimates) can be slow.

## 6.8.3 Acquisition functions

In BayesOpt, we use an **acquisition function** (also called a **merit function**) to evaluate the expected utility of each possible point we could query:  $\alpha(\mathbf{x}|\mathcal{D}_n) = \mathbb{E}_{p(y|\mathbf{x}, \mathcal{D}_n)}[U(\mathbf{x}, y; \mathcal{D}_n)]$ , where  $y = f(\mathbf{x})$  is the unknown value of the function at point  $\mathbf{x}$ , and  $U()$  is a utility function. Different utility functions give rise to different acquisition functions, as we discuss below. We usually choose functions so that the utility of picking a point that has already been queried is small (or 0, in the case of noise-free observations), in order to encourage exploration.

### 6.8.3.1 Probability of improvement

Let us define  $V_n = \max_{i=1}^n y_i$  to be the best value observed so far (known as the **incumbent**). (If the observations are noisy, using the highest mean value  $\max_i \mathbb{E}[f(\mathbf{x}_i)]$  is a reasonable alternative [WF16].) Then we define the utility of some new point  $\mathbf{x}$  using  $U(\mathbf{x}, y; \mathcal{D}_n) = \mathbb{I}(y > V_n)$ . This gives reward iff the new value is better than the incumbent. The corresponding acquisition function is then given by the expected utility,  $\alpha_{PI}(\mathbf{x}; \mathcal{D}_n) = p(f(\mathbf{x}) > V_n | \mathcal{D}_n)$ . This is known as the **probability of improvement** [Kus64]. If  $p(f | \mathcal{D}_n)$  is a GP, then this quantity can be computed in closed form, as follows:

$$\alpha_{PI}(\mathbf{x}; \mathcal{D}_n) = p(f(\mathbf{x}) > V_n | \mathcal{D}_n) = \Phi(\gamma(\mathbf{x}, V_n)) \quad (6.222)$$

where  $\Phi$  is the cdf of the  $\mathcal{N}(0, 1)$  distribution and

$$\gamma(\mathbf{x}, \tau) = \frac{\mu_n(\mathbf{x}) - \tau}{\sigma_n(\mathbf{x})} \quad (6.223)$$

### 6.8.3.2 Expected improvement

The problem with PI is that all improvements are considered equally good, so the method tends to exploit quite aggressively [Jon01]. A common alternative takes into account the amount of

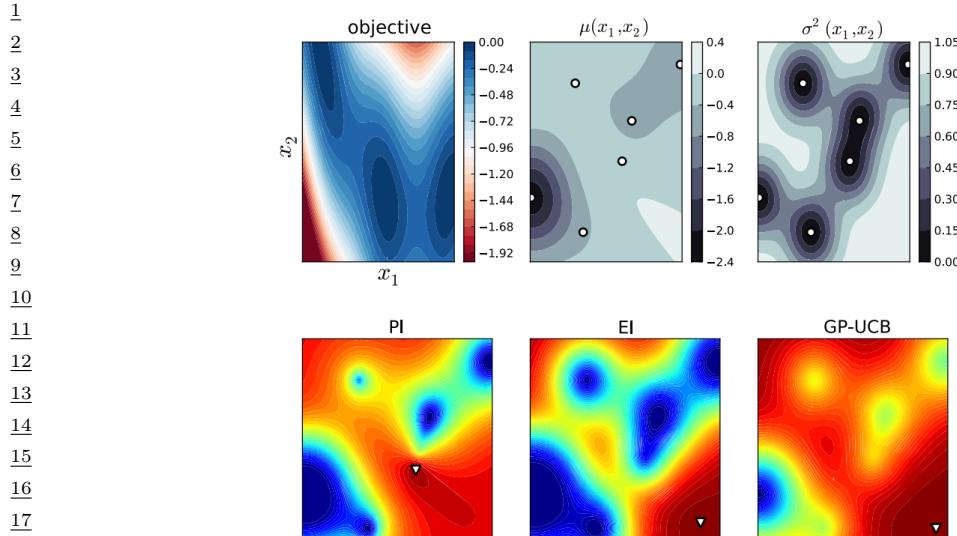


Figure 6.11: The first row shows the objective function, (the Branin function defined on  $\mathbb{R}^2$ ), and its posterior mean and variance using a GP estimate. White dots are the observed data points. The second row shows 3 different acquisition functions (probability of improvement, expected improvement, and upper confidence bound); the white triangles are the maxima of the corresponding acquisition functions. From Figure 6 of [BCF10]. Used with kind permission of Nando de Freitas.

improvement by defining  $U(\mathbf{x}, y; \mathcal{D}_n) = (y - V_n)\mathbb{I}(y > V_n)$  and

$$\alpha_{EI}(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_{\mathcal{D}_n}[U(\mathbf{x}, y)] = \mathbb{E}_{\mathcal{D}_n}[(f(\mathbf{x}) - V_n)\mathbb{I}(f(\mathbf{x}) > V_n)] \quad (6.224)$$

This acquisition function is known as the **expected improvement** (EI) criterion [Moc+96]. In the case of a GP surrogate, this has the following closed form expression:

$$\alpha_{EI}(\mathbf{x}; \mathcal{D}_n) = (\mu_n(\mathbf{x}) - \tau)\Phi(\gamma(\mathbf{x})) + \sigma_n(\mathbf{x})\phi(\gamma(\mathbf{x})) \quad (6.225)$$

where  $\phi()$  is the pdf of the  $\mathcal{N}(0, 1)$  distribution. The first term encourages exploitation (evaluating points with high mean) and the second term encourages exploration (evaluating points with high variance). This is illustrated in Figure 6.10.

#### 6.8.3.3 Upper confidence bound (UCB)

An alternative approach is to compute an **upper confidence bound** or UCB on the function, at some confidence level  $\beta_n$ , and then to define the acquisition function as follows:  $\alpha_{UCB}(\mathbf{x}; \mathcal{D}_n) = \mu_n(\mathbf{x}) + \beta_n\sigma_n(\mathbf{x})$ . This is the same as in the contextual bandit setting, discussed in Section 34.4.5, except we are optimizing over  $\mathbf{x} \in \mathcal{X}$ , rather than a finite set of arms  $a \in \{1, \dots, A\}$ . If we use a GP for our surrogate, the method is known as **GP-UCB** [Sri+10].

---

### 6.8.3.4 Thompson sampling

We introduced **Thompson sampling** in Section 34.4.6 in the context of multi-armed bandits, where the state space is finite,  $\mathcal{X} = \{1, \dots, A\}$ , and the acquisition function  $\alpha(a; \mathcal{D}_n)$  corresponds to the probability that arm  $a$  is the best arm. We can generalize this to real-valued input spaces  $\mathcal{X}$  using

$$\alpha(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_{p(\boldsymbol{\theta}|\mathcal{D}_n)} \left[ \mathbb{I} \left( \mathbf{x} = \operatorname{argmax}_{\mathbf{x}'} f_{\boldsymbol{\theta}}(\mathbf{x}') \right) \right] \quad (6.226)$$

We can compute a single sample approximation to this integral by sampling  $\tilde{\boldsymbol{\theta}} \sim p(\boldsymbol{\theta}|\mathcal{D}_n)$ . We can then pick the optimal action as follows:

$$\mathbf{x}_{n+1} = \operatorname{argmax}_{\mathbf{x}} \alpha(\mathbf{x}; \mathcal{D}_n) = \operatorname{argmax}_{\mathbf{x}} \mathbb{I} \left( \mathbf{x} = \operatorname{argmax}_{\mathbf{x}'} f_{\tilde{\boldsymbol{\theta}}}(\mathbf{x}') \right) = \operatorname{argmax}_{\mathbf{x}} f_{\tilde{\boldsymbol{\theta}}}(\mathbf{x}) \quad (6.227)$$

In other words, we greedily maximize the sampled surrogate.

For continuous spaces, Thompson sampling is harder to apply than in the bandit case, since we can't directly compute the best "arm"  $\mathbf{x}_{n+1}$  from the sampled function. Furthermore, when using GPs, there are some subtle technical difficulties with sampling a function, as opposed to sampling the parameters of a parametric surrogate model (see [HLHG14] for discussion).

### 6.8.3.5 Entropy search

Since our goal in BayesOpt is to find  $\mathbf{x}^* = \operatorname{argmax}_{\mathbf{x}} f(\mathbf{x})$ , it makes sense to try to directly minimize our uncertainty about the location of  $\mathbf{x}^*$ , which we denote by  $p_*(\mathbf{x}|\mathcal{D}_n)$ . We will therefore define the utility as follows:

$$U(\mathbf{x}, y; \mathcal{D}_n) = \mathbb{H}(\mathbf{x}^*|\mathcal{D}_n) - \mathbb{H}(\mathbf{x}^*|\mathcal{D}_n \cup \{(\mathbf{x}, y)\}) \quad (6.228)$$

where  $\mathbb{H}(\mathbf{x}^*|\mathcal{D}_n) = \mathbb{H}(p_*(\mathbf{x}|\mathcal{D}_n))$  is the entropy of the posterior distribution over the location of the optimum. This is known as the information gain criterion; the difference from the objective used in active learning is that here we want to gain information about  $\mathbf{x}^*$  rather than about  $f$  for all  $\mathbf{x}$ . The corresponding acquisition function is given by

$$\alpha_{ES}(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_{p(y|\mathbf{x}, \mathcal{D}_n)} [U(\mathbf{x}, y; \mathcal{D}_n)] = \mathbb{H}(\mathbf{x}^*|\mathcal{D}_n) - \mathbb{E}_{p(y|\mathbf{x}, \mathcal{D}_n)} [\mathbb{H}(\mathbf{x}^*|\mathcal{D}_n \cup \{(\mathbf{x}, y)\})] \quad (6.229)$$

This is known as **entropy search** [HS12].

Unfortunately, computing  $\mathbb{H}(\mathbf{x}^*|\mathcal{D}_n)$  is hard, since it requires a probability model over the input space. Fortunately, we can leverage the symmetry of mutual information to rewrite the acquisition function in Equation (6.229) as follows:

$$\alpha_{PES}(\mathbf{x}; \mathcal{D}_n) = \mathbb{H}(y|\mathcal{D}_n, \mathbf{x}) - \mathbb{E}_{\mathbf{x}^*|\mathcal{D}_n} [\mathbb{H}(y|\mathcal{D}_n, \mathbf{x}, \mathbf{x}^*)] \quad (6.230)$$

where we can approximate the expectation from  $p(\mathbf{x}^*|\mathcal{D}_n)$  using Thompson sampling. Now we just have to model uncertainty about the output space  $y$ . This is known as **predictive entropy search** [HLHG14].

1 **6.8.3.6 Knowledge gradient**

3 So far the acquisition functions we have considered are all greedy, in that they only look one step  
4 ahead. The **knowledge gradient** acquisition function, proposed in [FPD09], looks two steps ahead  
5 by considering the improvement we might expect to get if we query  $\mathbf{x}$ , update our posterior, and  
6 then exploit our knowledge by maximizing wrt our new beliefs. More precisely, let us define the best  
7 value we can find if we query one more point:

$$\underline{9} \quad V_{n+1}(\mathbf{x}, y) = \max_{\mathbf{x}'} \mathbb{E}_{p(f|\mathbf{x}, y, \mathcal{D}_n)} [f(\mathbf{x}')] \quad (6.231)$$

$$\underline{10} \quad V_{n+1}(\mathbf{x}) = \mathbb{E}_{p(y|\mathbf{x}, \mathcal{D}_n)} [V_{n+1}(\mathbf{x}, y)] \quad (6.232)$$

12 We define the KG acquisition function as follows:

$$\underline{14} \quad \alpha_{KG}(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_{\mathcal{D}_n} [(V_{n+1}(\mathbf{x}) - V_n) \mathbb{I}(V_{n+1}(\mathbf{x}) > V_n)] \quad (6.233)$$

16 Compare this to the EI function in Equation (6.224). Thus we pick the point  $\mathbf{x}_{n+1}$  such that  
17 observing  $f(\mathbf{x}_{n+1})$  will give us knowledge which we can then exploit, rather than directly trying to  
18 find a better point with better  $f$  value.

19

20 **6.8.3.7 Optimizing the acquisition function**

22 The acquisition function  $\alpha(\mathbf{x})$  is often multimodal (see e.g., Figure 6.11), since it will be 0 at all the  
23 previously queried points (assuming noise-free observations). Consequently maximizing this function  
24 can be a hard subproblem in itself [WHD18; Rub+20].

25 In the continuous setting, it is common to use multi-restart BFGS or grid search. We can also use  
26 the cross-entropy method (Section 6.9.5), using mixtures of Gaussians [BK10] or VAEs [Fau+18] as  
27 the generative model over  $\mathbf{x}$ . In the discrete, combinatorial setting (e.g., when optimizing biological  
28 sequences), [Bel+19] use regularized evolution, (Section 6.9.3), and [Ang+20] use proximal policy  
29 optimization (Section 35.3.4). Many other combinations are possible.

30

31 **6.8.4 Other issues**

32 There are many other issues that need to be tackled when using Bayesian optimization, a few of  
33 which we briefly mention below.

35

36 **6.8.4.1 Parallel (batch) queries**

37 In some cases, we want to query the objective function at multiple points in parallel; this is known as  
38 **batched Bayesian optimization**. Now we need to optimize over a set of possible queries, which is  
39 computationally even more difficult than the regular case. See [WHD18; DBB20] for some recent  
40 papers on this topic.

41

42 **6.8.4.2 Conditional parameters**

44 BayesOpt is often applied to hyper-parameter optimization. In many applications, some hyperparameters  
45 are only well-defined if other ones take on specific values. For example, suppose we are trying  
46 to automatically tune a classifier, as in the **Auto-Sklearn** system [Feu+15], or the **Auto-Weka**

47

system [Kot+17]. If the method chooses to use a neural network, it also needs to specify the number of layers, and number of hidden units per layer; but if it chooses to use a decision tree, it instead should specify different hyperparameters, such as the maximum tree depth.

We can formalize such problems by defining the search space in terms of a tree or DAG (directed acyclic graph), where different subsets of the parameters are defined at each leaf. Applying GPs to this setting requires non-standard kernels, such as those discussed in [Swe+13; Jen+17]. Alternatively, we can use other forms of Bayesian regression, such as ensembles of random forests [HHLB11], which can easily handle conditional parameter spaces.

#### 6.8.4.3 Multi-fidelity surrogates

In some cases, we can construct surrogate functions with different levels of accuracy, each of which may take variable amounts of time to compute. In particular, let  $f(\mathbf{x}, s)$  be an approximation to the true function at  $\mathbf{x}$  with fidelity  $s$ . The goal is to solve  $\max_{\mathbf{x}} f(\mathbf{x}, 0)$  by observing  $f(\mathbf{x}, s)$  at a sequence of  $(\mathbf{x}_i, s_i)$  values, such that the total cost  $\sum_{i=1}^n c(s_i)$  is below some budget. For example, in the context of hyperparameter selection,  $s$  may control how long we run the parameter optimizer for, or how large the validation set is.

In addition to choosing what fidelity to use for an experiment, we may choose to terminate expensive trials (queries) early, if the results of their cheaper proxies suggest they will not be worth running to completion (see e.g., [Str19; Li+17c; FKH17]). Alternatively, we may choose to resume an earlier aborted run, to collect more data on it, as in the **freeze-thaw algorithm** [SSA14].

#### 6.8.4.4 Constraints

If we want to maximize a function subject to known constraints, we can simply build the constraints into the acquisition function. But if the constraints are unknown, we need to estimate the support of the feasible set in addition to estimating the function. In [GSA14], they propose the weighted EI criterion, given by  $\alpha_{wEI}(\mathbf{x}; \mathcal{D}_n) = \alpha_{EI}(\mathbf{x}; \mathcal{D}_n)h(\mathbf{x}; \mathcal{D}_n)$ , where  $h(\mathbf{x}; \mathcal{D}_n)$  is a GP with a Bernoulli observation model that specifies if  $\mathbf{x}$  is feasible or not. Of course, other methods are possible. For example, [HL+16b] propose a method based on predictive entropy search.

### 6.9 Derivative free optimization

**Derivative free optimization** or **DFO** refers to a class of techniques for optimizing functions without using derivatives. This is useful for blackbox function optimization as well as discrete optimization. If the function is expensive to evaluate, we can use Bayesian optimization (Section 6.8). If the function is cheap to evaluate, we can use stochastic local search methods or evolutionary search methods, as we discuss below.

#### 6.9.1 Local search

In this section, we discuss heuristic optimization algorithms that try to find the global maximum in a discrete, unstructured search space. These algorithms replace the local gradient based update, which has the form  $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \eta_t \mathbf{d}_t$ , with the following discrete analog:

$$\mathbf{x}_{t+1} = \underset{\mathbf{x} \in \text{nbr}(\mathbf{x}_t)}{\operatorname{argmax}} \mathcal{L}(\mathbf{x}) \quad (6.234)$$

1 where  $\text{nbr}(\mathbf{x}_t) \subseteq \mathcal{X}$  is the set of **neighbors** of  $\mathbf{x}_t$ . This is called **hill climbing**, **steepest ascent**,  
 2 or **greedy search**.

3 If the “neighborhood” of a point contains the entire space, Equation (6.234) will return the global  
 4 optimum in one step, but usually such a global neighborhood is too large to search exhaustively.  
 5 Consequently we usually define local neighborhoods. For example, consider the **8-queens problem**.  
 6 Here the goal is to place queens on an  $8 \times 8$  chessboard so that they don’t attack each other (see  
 7 Figure 6.14). The state space has the form  $\mathcal{X} = 64^8$ , since we have to specify the location of each  
 8 queen on the grid. However, due to the constraints, there are only  $8^8 \approx 17M$  feasible states. We  
 9 define the neighbors of a state to be all possible states generated by moving a single queen to another  
 10 square in the same column, so each node has  $8 \times 7 = 56$  neighbors. According to [RN10, p.123], if we  
 11 start at a randomly generated 8-queens state, steepest ascent gets stuck at a local maximum 86% of  
 12 the time, so it only solves 14% of problem instances. However, it is fast, taking an average of 4 steps  
 13 when it succeeds and 3 when it gets stuck.

14 In the sections below, we discuss slightly smarter algorithms that are less likely to get stuck in  
 15 local maxima.

16

#### 17 6.9.1.1 Stochastic local search

18

19 Hill climbing is greedy, since it picks the best point in its local neighborhood, by solving Equa-  
 20 tion (6.234) exactly. One way to reduce the chance of getting stuck in local maxima is to approximately  
 21 maximize this objective at each step. For example, we can define a probability distribution over the  
 22 uphill neighbors, proportional to how much they improve, and then sample one at random. This is  
 23 called **stochastic hill climbing**. If we gradually decrease the entropy of this probability distribution  
 24 (so we become greedier over time), we get a method called simulated annealing, which we discuss in  
 25 Section 12.9.1.

26 Another simple technique is to use greedy hill climbing, but then whenever we reach a local  
 27 maximum, we start again from a different random starting point. This is called **random restart**  
 28 **hill climbing**. To see the benefit of this, consider again the 8-queens problem. If each hill-climbing  
 29 search has a probability of  $p \approx 0.14$  of success, then we expect to need  $R = 1/p \approx 7$  restarts until we  
 30 find a valid solution. The expected number of total steps can be computed as follows. Let  $N_1 = 4$   
 31 be the average number of steps for successful trials, and  $N_0 = 3$  be the average number of steps for  
 32 failures. Then the total number of steps on average is  $N_1 + (R - 1)N_0 = 4 + 6 \times 3 = 22$ . Since each  
 33 step is quick, the overall method is very fast. For example, it can solve an  $n$ -queens problem with  
 34  $n = 1M$  in under a minute.

35 Of course, solving the  $n$ -queens problem is not the most useful task in practice. However, it is  
 36 typical of several real-world **boolean satisfiability problems**, which arise in problems ranging  
 37 from AI planning to model checking (see e.g., [SLM92]). In such problems, simple **stochastic local**  
 38 **search (SLS)** algorithms of the kind we have discussed work surprisingly well (see e.g., [HS05]).

39

#### 40 6.9.1.2 Tabu search

41

42 Hill climbing will stop as soon as it reaches a local maximum or a plateau. Obviously one can perform  
 43 a random restart, but this would ignore all the information that had been gained up to this point. A  
 44 more intelligent alternative is called **tabu search** [GL97]. This is like hill climbing, except it allows  
 45 moves that decrease (or at least do not increase) the scoring function, provided the move is to a new  
 46 47

---

1           

2 **Algorithm 6:** Tabu search.

---

3   1  $t := 0$  // counts iterations

4   2  $c := 0$  // counts number of steps with no progress

5   3 Initialize  $\mathbf{x}_0$

6   4  $\mathbf{x}^* := \mathbf{x}_0$  // current best incumbent

7   5 **while**  $c < c_{\max}$  **do**

8   6     $\mathbf{x}_{t+1} = \operatorname{argmax}_{\mathbf{x} \in \text{nbr}(\mathbf{x}_t) \setminus \{\mathbf{x}_{t-\tau}, \dots, \mathbf{x}_{t-1}\}} f(\mathbf{x})$

9   7    **if**  $f(\mathbf{x}_t) > f(\mathbf{x}^*)$  **then**

10   8       $\mathbf{x}^* := \mathbf{x}_t$

11   9       $c := 0$

12   10     **else**

13   11       $c := c + 1$

14   12     

15   13     $t := t + 1$

16   14   

17   15   **return**  $\mathbf{x}^*$

---

18   16   

19   17   

20 state that has not been seen before. We can enforce this by keeping a tabu list which tracks the  
 21  $\tau$  most recently visited states. This forces the algorithm to explore new states, and increases the  
 22 chances of escaping from local maxima. We continue to do this for up to  $c_{\max}$  steps (known as the  
 23 “tabu tenure”). The pseudocode can be found in Algorithm 6. (If we set  $c_{\max} = 1$ , we get greedy hill  
 24 climbing.)

25 For example, consider what happens when tabu search reaches a hill top,  $\mathbf{x}_t$ . At the next step, it  
 26 will move to one of the neighbors of the peak,  $\mathbf{x}_{t+1} \in \text{nbr}(\mathbf{x}_t)$ , which will have a lower score. At the  
 27 next step, it will move to the neighbor of the previous step,  $\mathbf{x}_{t+2} \in \text{nbr}(\mathbf{x}_{t+1})$ ; the tabu list prevents  
 28 it cycling back to  $\mathbf{x}_t$  (the peak), so it will be forced to pick a neighboring point at the same height or  
 29 lower. It continues in this way, “circling” the peak, possibly being forced downhill to a lower level-set  
 30 (an inverse **basin flooding** operation), until it finds a ridge that leads to a new peak, or until it  
 31 exceeds a maximum number of non-improving moves.

32 According to [RN10, p.123], tabu search increases the percentage of 8-queens problems that can  
 33 be solved from 14% to 94%, although this variant takes an average of 21 steps for each successful  
 34 instance and 64 steps for each failed instance.

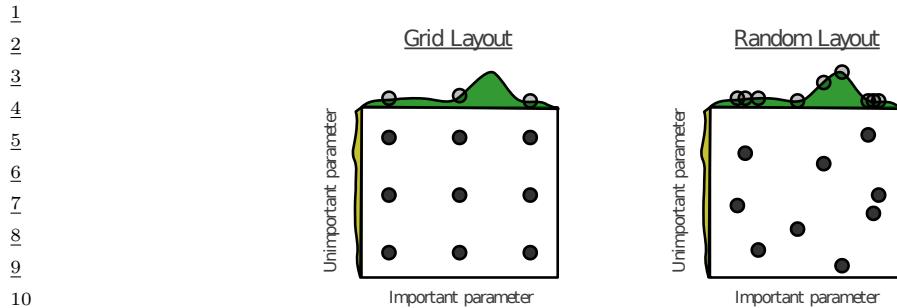
35

### 36 6.9.1.3 Random search

37 A surprisingly effective strategy in problems where we know nothing about the objective is to use  
 38 **random search**. In this approach, each iterate  $\mathbf{x}_{t+1}$  is chosen uniformly at random from  $\mathcal{X}$ . This  
 39 should always be tried as a baseline.

40 In [BB12], they applied this technique to the problem of hyper-parameter optimization for some  
 41 ML models, where the objective is performance on a validation set. In their examples, the search  
 42 space is continuous,  $\Theta = [0, 1]^D$ . It is easy to sample from this at random. The standard alternative  
 43 approach is to quantize the space into a fixed set of values, and then to evaluate them all; this is  
 44 known as **grid search**. (Of course, this is only feasible if the number of dimensions  $D$  is small.)  
 45 They found that random search outperformed grid search. The intuitive reason for this is that many  
 46

47



11 *Figure 6.12: Illustration of grid search (left) vs random search (right). From Figure 1 of [BB12]. Used with*  
 12 *kind permission of James Bergstra.*

13

14

15

16 hyper-parameters do not make much difference to the objective function, as illustrated in Figure 6.12.  
 17 Consequently it is a waste of time to place a fine grid along such unimportant dimensions.

18 RS has also been used to optimize the parameters of MDP policies, where the objective has  
 19 the form  $f(\mathbf{x}) = \mathbb{E}_{\tau \sim \pi_{\mathbf{x}}} [R(\tau)]$  is the expected reward of trajectories generated by using a policy  
 20 with parameters  $\mathbf{x}$ . For policies with few free parameters, RS can outperform more sophisticated  
 21 reinforcement learning methods described in Chapter 35, as shown in [MGR18]. In cases where  
 22 the policy has a large number of parameters, it is sometimes possible to project them to a lower  
 23 dimensional random subspace, and perform optimization (either grid search or random search) in  
 24 this subspace [Li+18a].

25

### 26 6.9.2 Simulated annealing

27

28 **Simulated annealing** [KJV83; LA87] is a **stochastic local search** algorithm (Section 6.9.1.1)  
 29 that attempts to find the global minimum of a black-box function  $\mathcal{E}(\mathbf{x})$ , where  $\mathcal{E}()$  is known as the  
 30 **energy function**. The method works by converting the energy to an (unnormalized) probability  
 31 distribution over states by defining  $p(\mathbf{x}) = \exp(-\mathcal{E}(\mathbf{x}))$ , and then using a variant of the **Metropolis**  
 32 **Hastings** algorithm to sample from a set of probability distributions, designed so that at the final  
 33 step, the method samples from one of the modes of the distribution, i.e., it finds one of the most  
 34 likely states, or lowest energy states. This approach can be used for both discrete and continuous  
 35 optimization. See Section 12.9.1 for details.

36

### 37 6.9.3 Evolutionary algorithms

38

39 Stochastic local search (SLS) maintains a single “best guess” at each step,  $\mathbf{x}_t$ . If we run this for  
 40  $T$  steps, and restart  $K$  times, the total cost is  $TK$ . A natural alternative is to maintain a set or  
 41 **population** of  $K$  good candidates,  $\mathcal{S}_t$ , which we try to improve at each step. This is called an  
 42 **evolutionary algorithm (EA)**. If we run this for  $T$  steps, it also takes  $TK$  time; however, it can  
 43 often get better results than multi-restart SLS, since the search procedure explores more of the space  
 44 in parallel, and information from different members of the population can be shared. Many versions  
 45 of EA are possible, as we discuss below.

46 Since EA algorithms draw inspiration from the biological process of evolution, they also borrow  
 47



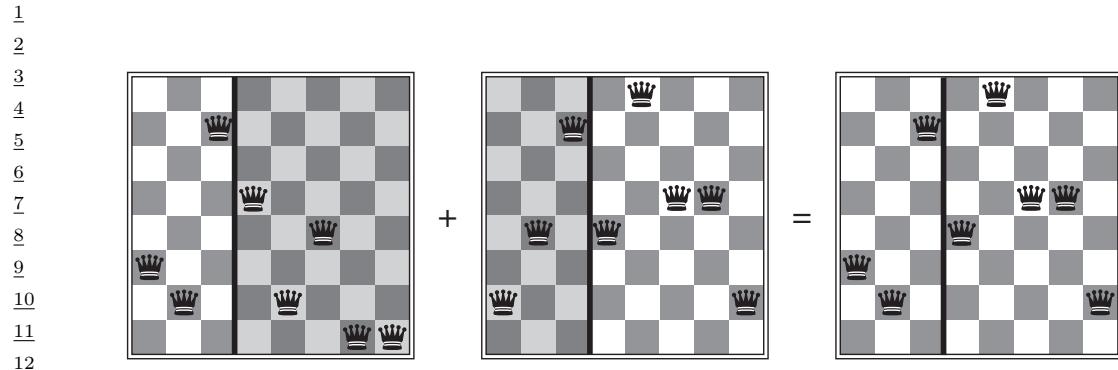
Figure 6.13: Illustration of a genetic algorithm applied to the 8-queens problem. (a) Initial population of 4 strings. (b) We rank the members of the population by fitness, and then compute their probability of mating. Here the integer numbers represent the number of nonattacking pairs of queens, so the global maximum has a value of 28. We pick an individual  $\theta$  with probability  $p(\theta) = \mathcal{L}(\theta)/Z$ , where  $Z = \sum_{\theta \in \mathcal{P}} \mathcal{L}(\theta)$  sums the total fitness of the population. For example, we pick the first individual with probability  $24/78 = 0.31$ , the second with probability  $23/78 = 0.29$ , etc. In this example, we pick the first individual once, the second twice, the third one once, and the last one does not get to breed. (c) A split point on the “chromosome” of each parent is chosen at random. (d) The two parents swap their chromosome halves. (e) We can optionally apply pointwise mutation. From Figure 4.6 of [RN10]. Used with kind permission of Peter Norvig.

a lot of its terminology. The **fitness** of a member of the population is the value of the objective function (possibly normalized across population members). The members of the population at step  $t + 1$  are called the **offspring**. These can be created by randomly choosing a **parent** from  $\mathcal{S}_t$  and applying a random **mutation** to it. This is like asexual reproduction. Alternatively we can create an offspring by choosing two parents from  $\mathcal{S}_t$ , and then combining them in some way to make a child, as in sexual reproduction; combining the parents is called **recombination**. (It is often followed by mutation.)

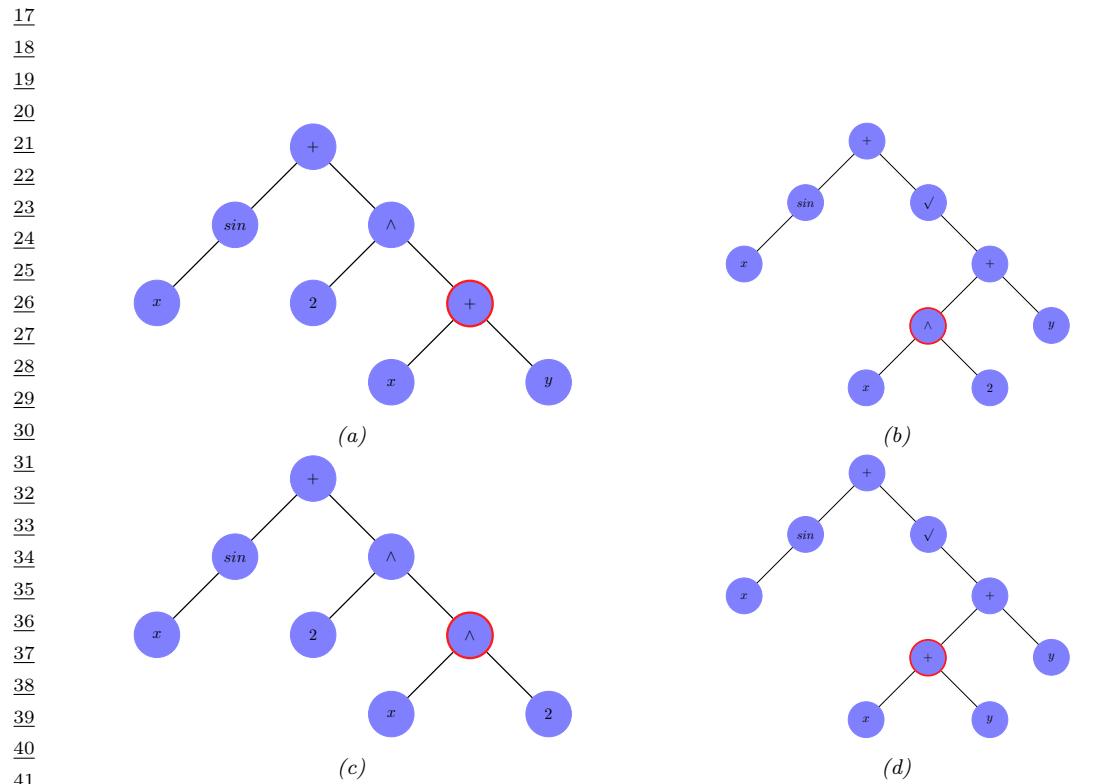
The procedure by which parents are chosen is called the **selection function**. In **truncation selection**, each parent is chosen from the fittest  $K$  members of the population (known as the **elite set**). In **tournament selection**, each parent is the fittest out of  $K$  randomly chosen members. In **fitness proportionate selection**, also called **roulette wheel selection**, each parent is chosen with probability proportional to its fitness relative to the others. We can also “kill off” the oldest members of the population, and then select parents based on their fitness; this is called **regularized evolution** [Rea+19]).

In addition to the selection rule for parents, we need to specify the recombination and mutation rules. There are many possible choices for these heuristics. We briefly mention a few of them below.

- In a **genetic algorithm (GA)** [Gol89; Hol92], we use mutation and a particular recombination method based on **crossover**. To implement crossover, we assume each individual is represented as a vector of integers or binary numbers, by analogy to **chromosomes**. We pick a split point along the chromosome for each of the two chosen parents, and then swap the strings, as illustrated in Figure 6.13.
- In **genetic programming** [Koz92], we use a tree-structured representation of individuals, instead of a bit string. This representation ensures that all crossovers result in valid children,



13 *Figure 6.14:* The 8-queens states corresponding to the first two parents in Figure 6.13(c) and their first child  
 14 in Figure 6.13(d). We see that the encoding 32752411 means that the first queen is in row 3 (counting from  
 15 the bottom left), the second queen is in row 2, etc. The shaded columns are lost in the crossover, but the  
 16 unshaded columns are kept. From Figure 4.7 of [RN10]. Used with kind permission of Peter Norvig.



42 *Figure 6.15:* Illustration of crossover operator in a genetic program. (a-b) the two parents, representing  
 43  $\sin(x) + (x+y)^2$  and  $\sin(x) + \sqrt{x^2+y}$ . The red circles denote the two crossover points. (c-d) the two  
 44 children, representing  $\sin(x) + (x^2)^2$  and  $\sin(x) + \sqrt{x+y+y}$ . Adapted from Figure 9.2 of [Mit97]

as illustrated in Figure 6.15. Genetic programming can be useful for finding good programs as well as other structured objects, such as neural networks. In **evolutionary programming**, the structure of the tree is fixed and only the numerical parameters are evolved.

- In **surrogate assisted EA**, a surrogate function  $\hat{f}(s)$  is used instead of the true objective function  $f(s)$  in order to speed up the evaluation of members of the population (see [Jin11] for a survey). This is similar to the use of response surface models in Bayesian optimization (Section 6.8), except it does not deal with the explore-exploit tradeoff.
- In a **memetic algorithm** [MC03], we combine mutation and recombination with standard local search.

Evolutionary algorithms have been applied to a large number of applications, including training neural networks (this combination is known as **neuroevolution** [Sta+19]). An efficient JAX-based library for (neuro)-evolution can be found at <https://github.com/google/evojax>.

#### 6.9.4 Estimation of distribution (EDA) algorithms

EA methods maintain a population of good candidate solutions, which can be thought of as an implicit (nonparametric) density model over states with high fitness. [BC95] proposed to “remove the genetics from GAs”, by explicitly learning a probabilistic model over the configuration space that puts its mass on high scoring solutions. That is, the population becomes the set of parameters of a generative model,  $\theta_t$ .

One way to learn such a model is as follows. We start by creating a sample of  $K' > K$  candidate solutions from the current model,  $\mathcal{S}_t = \{\mathbf{x}_k \sim p(\mathbf{x}|\theta_t)\}$ . We then rank the samples using the fitness function, and then pick the most promising subset  $\mathcal{S}_t^*$  of size  $K$  using a selection operator (this is known as **truncation selection**). Finally, we fit a new probabilistic model  $p(\mathbf{x}|\theta_{t+1})$  to  $\mathcal{S}_t^*$  using maximum likelihood estimation. This is called the **estimation of distribution** or **EDA** algorithm (see e.g., [LL02; PSCP06; Hau+11; PHL12; Hu+12; San17; Bal17]).

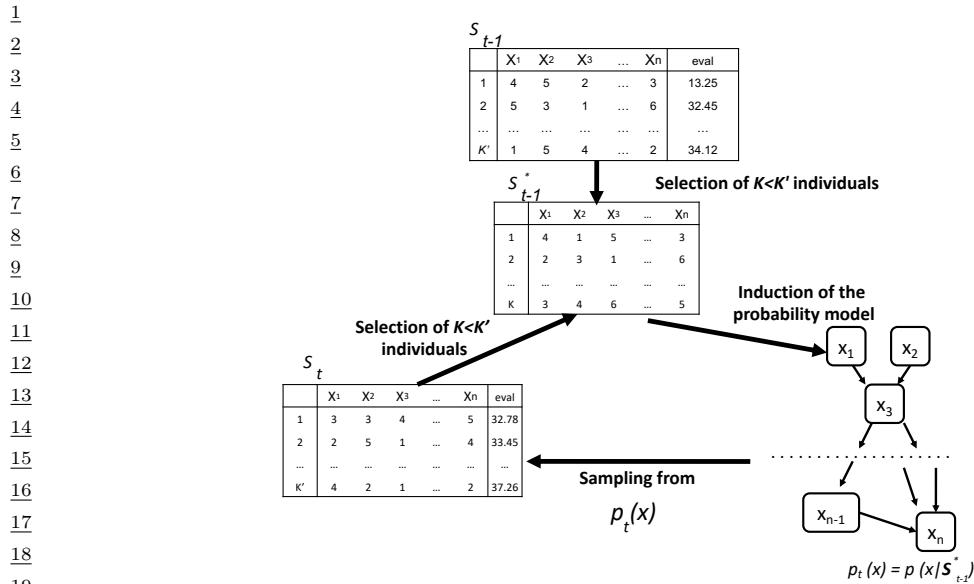
Note that EDA is equivalent to minimizing the cross-entropy between the empirical distribution defined by  $\mathcal{S}_t^*$  and the model distribution  $p(\mathbf{x}|\theta_{t+1})$ . Thus EDA is related to the **cross entropy method**, as described in Section 6.9.5, although CEM usually assumes the special case where  $p(\mathbf{x}|\mathbf{w}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . EDA is also closely related to the EM algorithm, as discussed in [Bro+20a].

As a simple example, suppose the configuration space is bit strings of length  $D$ , and the fitness function is  $f(\mathbf{x}) = \sum_{d=1}^D x_d$ , where  $x_d \in \{0, 1\}$  (this is called the **one-max** function in the EA literature). A simple probabilistic model for this is a fully factored model of the form  $p(\mathbf{x}|\theta) = \prod_{d=1}^D \text{Ber}(x_d|\theta_d)$ . Using this model inside of DBO results in a method called univariate marginal distribution algorithm or **UMDA**.

We can estimate the parameters of the Bernoulli model by setting  $\theta_d$  to the fraction of samples in  $\mathcal{S}_t^*$  that have bit  $d$  turned on. Alternatively, we can incrementally adjust the parameters. The population-based incremental learning (**PBIL**) algorithm [BC95] applies this idea to the factored Bernoulli model, resulting in the following update:

$$\hat{\theta}_{d,t+1} = (1 - \eta_t)\hat{\theta}_{d,t} + \eta_t \bar{\theta}_{d,t} \quad (6.235)$$

where  $\bar{\theta}_{d,t} = \frac{1}{N_t} \sum_{k=1}^{N_t} \mathbb{I}(x_{k,d} = 1)$  is the MLE estimated from the  $K = |\mathcal{S}_t^*|$  samples generated in the current iteration, and  $\eta_t$  is a learning rate.



20     Figure 6.16: Illustration of the BOA algorithm (EDA applied to a generative model structured as a Bayes  
21     net). Adapted from Figure 3 of [PHL12].  
22

23  
24  
25     It is straightforward to use more expressive probability models that capture dependencies between  
26     the parameters (these are known as **building blocks** in the EA literature). For example, in the case  
27     of real-valued parameters, we can use a multivariate Gaussian,  $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . The resulting  
28     method is called the **estimation of multivariate normal algorithm** or **EMNA**, [LL02]. (See  
29     also Section 6.9.5.)

30     For discrete random variables, it is natural to use probabilistic graphical models (Chapter 4) to  
31     capture dependencies between the variables. [BD97] learn a tree-structured graphical model using  
32     the Chow-Liu algorithm (Supplementary Section 30.1.1); [BJV97] is a special case of this where the  
33     graph is a tree. We can also learn more general graphical model structures (see e.g., [LL02]). We  
34     typically use a Bayes net (Section 4.2), since we can use ancestral sampling (Section 4.2.5) to easily  
35     generate samples; the resulting method is therefore called the **Bayesian Optimization Algorithm**  
36     (**BOA**) [PGCP00].<sup>4</sup> The hierarchical BOA (**hBOA**) algorithm [Pel05] extends this by using decision  
37     trees and decision graphs to represent the local CPTs in the Bayes net (as in [CHM97]), rather than  
38     using tables. In general, learning the structure of the probability model for use in EDA is called  
39     linkage learning, by analogy to how genes can be linked together if they can be co-inherited as a  
40     building block.

41     We can also use deep generative models to represent the distribution over good candidates. For  
42     example, [CSF16] use denoising autoencoders and NADE models (Section 22.2), [Bal17] uses a  
43     DNN regressor which is then inverted using gradient descent on the inputs, [PRG17] uses RBMs  
44

45     4. This should not be confused with the Bayesian optimization methods we discuss in Section 6.8, that uses response  
46     surface modeling to model  $p(f(\mathbf{x}))$  rather than  $p(\mathbf{x}^*)$ .

(Section 4.3.3.2), [GSM18] uses VAEs (Section 21.2), etc. Such models might take more data to fit (and therefore more function calls), but can potentially model the probability landscape more faithfully. (Whether that translates to better optimization performance is not clear, however.)

### 6.9.5 Cross-entropy method

The **cross-entropy method** [Rub97; RK04; Boe+05] is a special case of EDA (Section 6.9.4) in which the population is represented by a multivariate Gaussian. In particular, we set  $\boldsymbol{\mu}_{t+1}$  and  $\boldsymbol{\Sigma}_{t+1}$  to the empirical mean and covariance of  $\mathcal{S}_{t+1}^*$ , which are the top  $K$  samples. This is closely related to the SMC algorithm for sampling rare events discussed in Section 13.6.4.

The CEM is sometimes used for model-based RL (Section 35.4), since it is simple and can find reasonably good optima of multi-modal objectives. It is also sometimes used inside of Bayesian optimization (Section 6.8), to optimize the multi-modal acquisition function (see [BK10]).

#### 6.9.5.1 Differentiable CEM

The **differentiable CEM** method of [AY19] replaces the top  $K$  operator with a soft, differentiable approximation, which allows the optimizer to be used as part of an end-to-end differentiable pipeline. For example, we can use this to create a differentiable model predictive control (MPC) algorithm (Section 35.4.1), as described in Section 35.4.5.2.

The basic idea is as follows. Let  $\mathcal{S}_t = \{\mathbf{x}_{t,i} \sim p(\mathbf{x}|\boldsymbol{\theta}_t) : i = 1 : K'\}$  represent the current population, with fitness values  $v_{t,i} = f(\mathbf{x}_{t,i})$ . Let  $v_{t,K}^*$  be the  $K$ 'th smallest value. In CEM, we compute the set of top  $K$  samples,  $\mathcal{S}_t^* = \{i : v_{t,i} \geq v_{t,K}^*\}$ , and then update the model based on these:  $\boldsymbol{\theta}_{t+1} = \operatorname{argmax}_{\boldsymbol{\theta}} \sum_{i \in \mathcal{S}_t} p_t(i) \log p(\mathbf{x}_{t,i}|\boldsymbol{\theta})$ , where  $p_t(i) = \mathbb{I}(i \in \mathcal{S}_t^*) / |\mathcal{S}_t^*|$ . In the differentiable version, we replace the sparse distribution  $\mathbf{p}_t$  with the “soft” dense distribution  $\mathbf{q}_t = \Pi(\mathbf{p}_t; \tau, K)$ , where

$$\Pi(\mathbf{p}; \tau, K) = \underset{\mathbf{0} \leq \mathbf{q} \leq \mathbf{1}}{\operatorname{argmin}} -\mathbf{p}^\top \mathbf{q} - \tau \mathbb{H}(\mathbf{q}) \quad \text{s.t. } \mathbf{1}^\top \mathbf{q} = K \quad (6.236)$$

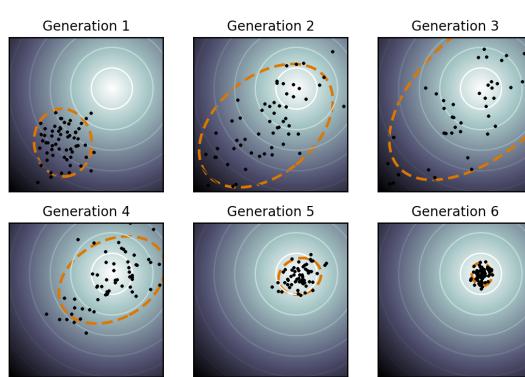
projects the distribution  $\mathbf{p}$  onto the polytope of distributions which sum to  $K$ . (Here  $\mathbb{H}(\mathbf{q}) = -\sum_i q_i \log(q_i) + (1 - q_i) \log(1 - q_i)$  is the entropy, and  $\tau > 0$  is a temperature parameter.) This projection operator (and hence the whole DCEM algorithm) can be backpropagated through using implicit differentiation [AKZK19].

### 6.9.6 Evolutionary strategies

**Evolution strategies** [Wie+14] are a form of distribution-based optimization in which the distribution over the population is represented by a Gaussian,  $p(\mathbf{x}|\boldsymbol{\theta}_t)$  (see e.g., [Sal+17b]). Unlike CEM, the parameters are updated using gradient ascent applied to the expected value of the objective, rather than using MLE on a set of elite samples. More precisely, consider the smoothed objective  $\mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}_{p(\mathbf{x}|\boldsymbol{\theta})} [f(\mathbf{x})]$ . We can use the REINFORCE estimator (Section 6.5.3) to compute the gradient of this objective as follows:

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}_{p(\mathbf{x}|\boldsymbol{\theta})} [f(\mathbf{x}) \nabla_{\boldsymbol{\theta}} \log p(\mathbf{x}|\boldsymbol{\theta})] \quad (6.237)$$

This can be approximated by drawing Monte Carlo samples. We discuss how to compute this gradient below.



*Figure 6.17: Illustration of the CMA-ES method applied to a simple 2d function. The dots represent members of the population, and the dashed orange ellipse represents the multivariate Gaussian. From <https://en.wikipedia.org/wiki/CMA-ES>. Used with kind permission of Wikipedia author Sentewolf.*

#### 6.9.6.1 Natural evolutionary strategies

If the probability model is in the exponential family, we can compute the natural gradient (Section 6.4), rather than the “vanilla” gradient, which can result in faster convergence. Such methods are called **natural evolution strategies** [Wie+14].

#### 6.9.6.2 CMA-ES

The **CMA-ES** method of [Han16], which stands for “covariance matrix adaptation evolution strategy” is a kind of NES. It is very similar to CEM except it updates the parameters in a special way. In particular, instead of computing the new mean and covariance using unweighted MLE on the elite set, we attach weights to the elite samples based on their rank. We then set the new mean to the weighted MLE of the elite set.

The update equations for the covariance are more complex. In particular, “evolutionary paths” are also used to accumulate the search directions across successive generations, and these are used to update the covariance. It can be shown that the resulting updates approximate the natural gradient of  $\mathcal{L}(\theta)$  without explicitly modeling the Fisher information matrix [Oll+17].

Figure 6.17 illustrates the method in action.

## 6.10 Optimal Transport

*This section is written by Marco Cuturi.*

In this section, we focus on **optimal transport** theory, a set of tools that have been proposed, starting with work by [Mon81], to compare two probability distributions. We start from a simple example involving only matchings, and work from there towards various extensions.

### 6.10.1 Warm-up: Matching optimally two families of points

Consider two families  $(\mathbf{x}_1, \dots, \mathbf{x}_n)$  and  $(\mathbf{y}_1, \dots, \mathbf{y}_n)$ , each consisting in  $n > 1$  distinct points taken from a set  $\mathcal{X}$ . A *matching* between these two families is a bijective mapping that assigns to each point  $\mathbf{x}_i$  another point  $\mathbf{y}_j$ . Such an assignment can be encoded by pairing indices  $(i, j) \in \{1, \dots, n\}^2$  such that they define a *permutation*  $\sigma$  in the symmetric group  $\mathcal{S}_n$ . With that convention and given a permutation  $\sigma$ ,  $\mathbf{x}_i$  would be assigned to  $\mathbf{y}_{\sigma_i}$ , the  $\sigma_i$ -th element in the second family.

**Matchings costs.** When matching a family with another, it is natural to consider the cost incurred when pairing any point  $\mathbf{x}_i$  with another point  $\mathbf{y}_j$ , for all possible pairs  $(i, j) \in \{1, \dots, n\}^2$ . For instance,  $\mathbf{x}_i$  might contain information on the current location of a taxi driver  $i$ , and  $\mathbf{y}_j$  that of a user  $j$  who has just requested a taxi; in that case,  $C_{ij} \in \mathbb{R}$  may quantify the cost (in terms of time, fuel or distance) required for taxi driver  $i$  to reach user  $j$ . Alternatively,  $\mathbf{x}_i$  could represent a vector of skills held by a job seeker  $i$  and  $\mathbf{y}_j$  a vector quantifying desirable skills associated with a job posting  $j$ ; in that case  $C_{ij}$  could quantify the numbers of hours required for worker  $i$  to carry out job  $j$ . We will assume without loss of generality that the values  $C_{ij}$  are obtained by evaluating a cost function  $c : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  on the pair  $(\mathbf{x}_i, \mathbf{y}_j)$ , namely  $C_{ij} = c(\mathbf{x}_i, \mathbf{y}_j)$ . In many applications of optimal transport, such cost functions have a geometric interpretation and are typically distance functions on  $\mathcal{X}$  as in Fig. 6.18, in which  $\mathcal{X} = \mathbb{R}^2$ , or as will be later discussed in §6.10.2.4.

**Least-cost Matchings.** Equipped with a cost function  $c$ , the *optimal* matching (or assignment) problem is that of finding a permutation that reaches the smallest total cost, as defined by the function

$$\min_{\sigma} E(\sigma) = \sum_{i=1}^n c(\mathbf{x}_i, \mathbf{y}_{\sigma_i}). \quad (6.238)$$

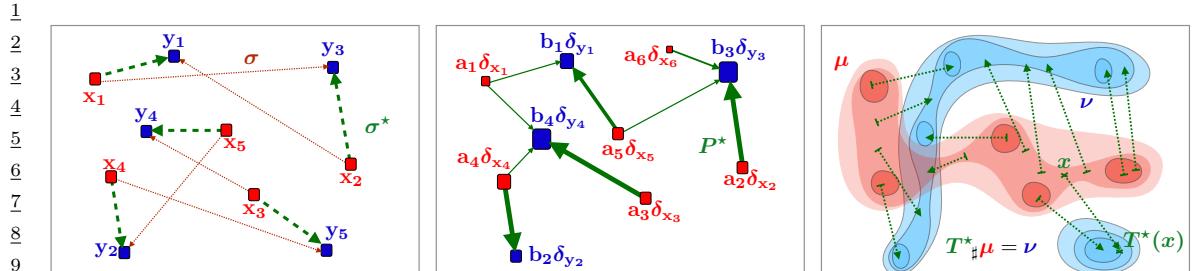
The optimal matching problem is arguably one of the simplest combinatorial optimization problems, tackled as early as the 19th century [JB65]. Although a naive enumeration of all permutations would require evaluating objective  $E$  a total of  $n!$  times, the Hungarian algorithm [Kuh55] was shown to provide the optimal solution in polynomial time [Mun57], and later refined to require in the worst case  $O(n^3)$  operations.

### 6.10.2 From Optimal Matchings to Kantorovich and Monge formulations

The optimal matching problem is relevant to many applications, but it suffers from a few limitations. One could argue that most of the optimal transport literature arises from the necessity to overcome these limitations and extend (6.238) to more general settings. An obvious issue arises when the number of points available in both families is not the same. The second limitation arises when considering a continuous setting, namely when trying to match (or morph) two probability densities, rather than families of atoms (discrete measures).

#### 6.10.2.1 Mass splitting

Suppose again that all points  $\mathbf{x}_i$  and  $\mathbf{y}_j$  describe skills, respectively held by a worker  $i$  and needed for a task  $j$  to be fulfilled in a factory. Since finding a matching is equivalent to finding a permutation in  $\{1, \dots, n\}$ , problem (6.238) cannot handle cases in which the number of workers is larger (or smaller) than the number of tasks. More problematically, the assumption that every single task is indivisible,



11 *Figure 6.18: (left) Matching a family of 5 points to another is equivalent to considering a permutation in*  $\{1, \dots, n\}$ . *When to each pair*  $(\mathbf{x}_i, \mathbf{y}_j) \in \mathbb{R}^2$  *is associated a cost equal to the distance*  $\|\mathbf{x}_i - \mathbf{y}_j\|$ , *the optimal matching problem involves finding a permutation*  $\sigma$  *that minimizes*  $\|\mathbf{x}_i - \mathbf{y}_{\sigma(i)}\|$  *for*  $i$  *in*  $\{1, 2, 3, 4, 5\}$ . *(middle)* The Kantorovich formulation of optimal transport generalizes optimal matchings, and arises when comparing discrete measures, that is families of weighted points that do not necessarily share the same size but do share the same total mass. The relevant variable is a matrix  $P$  of size  $n \times m$ , which must satisfy row-sum and column-sum constraints, and which minimizes its dot product with matrix  $C_{ij}$ . *(right)* another direct extension of the matching problem lies when, intuitively, the number  $n$  of points that is described is such that the considered measures become continuous densities. In that setting, and unlike the Kantorovich setting, the goal is to seek a map  $T : \mathcal{X} \rightarrow \mathcal{Y}$  which, to any point  $x$  in the support of the input measure  $\mu$  is associated a point  $y = T(x)$  in the support of  $\nu$ . The push-forward constraint  $T_*\mu = \nu$  ensures that  $\nu$  is recovered by applying map  $T$  to all points in the support of  $\mu$ ; the optimal map  $T^*$  is that which minimizes the distance between  $x$  and  $T(x)$ , averaged over  $\mu$ .

2324

25 or that workers are only able to dedicate themselves to a single task, is hardly realistic. Indeed,  
26 certain tasks may require more (or less) dedication than that provided by a single worker, whereas  
27 some workers may only be able to work part-time, or, on the contrary, be willing to put extra hours.  
28 The rigid machinery of permutations falls short of handling such cases, since permutations are by  
29 definition one-to-one associations. The Kantorovich formulation allows for *mass-splitting*, the idea  
30 that the effort provided by a worker or needed to complete a given task can be split. In practice,  
31 to each of the  $n$  workers is associated, in addition to  $\mathbf{x}_i$ , a positive number  $\mathbf{a}_i > 0$ . That number  
32 represents the amount of time worker  $i$  is able to provide. Similarly, we introduce numbers  $\mathbf{b}_j > 0$   
33 describing the amount of time needed to carry out each of the  $m$  tasks ( $n$  and  $m$  do not necessarily  
34 coincide). Worker  $i$  is therefore described as a pair  $(\mathbf{a}_i, \mathbf{x}_i)$ , mathematically equivalent to a *weighted*  
35 *Dirac measure*  $\mathbf{a}_i\delta_{\mathbf{x}_i}$ . The overall workforce available to the factory is described as a discrete measure  
36  $\sum_i \mathbf{a}_i\delta_{\mathbf{x}_i}$ , whereas its tasks are described in  $\sum_j \mathbf{b}_j\delta_{\mathbf{y}_j}$ . If one assumes further that the factory has  
37 a balanced workload, namely that  $\sum_i \mathbf{a}_i = \sum_j \mathbf{b}_j$ , then the Kantorovich [Kan42] formulation of  
38 optimal transport is:  
39

$$\text{OT}_C(\mathbf{a}, \mathbf{b}) \triangleq \min_{P \in \mathbb{R}_+^{n \times m}, P\mathbf{1}_m = \mathbf{a}, P^T\mathbf{1}_m = \mathbf{b}} \langle P, C \rangle \triangleq \sum_{i,j} P_{ij} C_{ij}. \quad (\text{K})$$

40 The interpretation behind such matrices is simple: each coefficient  $P_{ij}$  describes an allocation of  
41 time for worker  $i$  to spend on task  $j$ . The  $i$ -th row-sum must be equal to the total  $\mathbf{a}_i$  for the time  
42 constraint of worker  $i$  to be satisfied, whereas the  $j$ -th column-sum must be equal to  $\mathbf{b}_j$ , reflecting  
43 that the time needed to complete task  $j$  has been budgeted.

47

---

### 6.10.2.2 Monge formulation and optimal push-forward maps

By introducing mass-splitting, the Kantorovich formulation of optimal transport allows for a far more general comparison between discrete measures of different sizes and weights (middle plot of Fig. 6.18). Naturally, this flexibility comes with a downside: one can no longer associate to each point  $\mathbf{x}_i$  another point  $\mathbf{y}_j$  to which it is uniquely associated, as was the case with the classical matching problem. Interestingly, this property can be recovered in the limit where the measures become densities. Indeed, the Monge [Mon81] formulation of optimal transport allows to recover precisely that property, on the condition (loosely speaking) that measure  $\mu$  admits a density. In that setting, the analogous mathematical object guaranteeing that  $\mu$  is mapped onto  $\nu$  is that of *push-forward* maps morphing  $\mu$  to  $\nu$ , namely maps  $T$  such that for any measurable set  $A \subset \mathcal{X}$ ,  $\mu(T^{-1}(A)) = \nu(A)$ . When  $T$  is differentiable, and  $\mu, \nu$  have densities  $p$  and  $q$  w.r.t. the Lebesgue measure in  $\mathbb{R}^d$ , this statement is equivalent, thanks to the change of variables formula, to ensuring almost everywhere that:

$$q(T(x)) = p(x)|J_T(x)|, \quad (6.239)$$

where  $|J_T(x)|$  stands for the determinant of the Jacobian matrix of  $T$  evaluated at  $x$ .

Writing  $T_{\sharp}\mu = \nu$  when  $T$  does satisfy these conditions, the Monge [Mon81] problem consists in finding the best map  $T$  that minimizes the average cost between  $\mathbf{x}$  and its displacement  $T(\mathbf{x})$ ,

$$\inf_{T: T_{\sharp}\mu = \nu} \int_{\mathcal{X}} c(\mathbf{x}, T(\mathbf{x})) \mu(d\mathbf{x}). \quad (\text{M})$$

$T$  is therefore a map that pushes forward  $\mu$  to  $\nu$  globally, but which results, on average, in the smallest average cost. While very intuitive, the Monge problem turns out to be extremely difficult to solve in practice, since it is non-convex. Indeed, one can easily check that the constraint  $\{T_{\sharp}\mu = \nu\}$  is not convex, since one can easily find counter-examples for which  $T_{\sharp}\mu = \nu$  and  $T'_{\sharp}\nu$  yet  $(\frac{1}{2}T + \frac{1}{2}T')_{\sharp}\mu \neq \nu$ . Luckily, Kantorovich's approach also works for continuous measures, and yields a comparatively much simpler linear program.

### 6.10.2.3 Kantorovich formulation

The Kantorovich problem (K) can also be extended to a continuous setting: Instead of optimizing over a subset of matrices in  $\mathbb{R}^{n \times m}$ , consider  $\Pi(\mu, \nu)$ , the subset of joint probability distributions  $\mathcal{P}(\mathcal{X} \times \mathcal{X})$  with marginals  $\mu$  and  $\nu$ , namely

$$\Pi(\mu, \nu) \triangleq \{\pi \in \mathcal{P}(\mathcal{X}^2) : \forall A \subset \mathcal{X}, \pi(A \times \mathcal{X}) = \mu(A) \text{ and } \pi(\mathcal{X} \times A) = \nu(A)\}. \quad (6.240)$$

Note that  $\Pi(\mu, \nu)$  is not empty since it always contains the product measure  $\mu \otimes \nu$ . With this definition, the continuous formulation of (K) can be obtained as

$$\text{OT}_c(\mu, \nu) \triangleq \inf_{\pi \in \Pi(\mu, \nu)} \int_{\mathcal{X}^2} c \, d\pi. \quad (\text{K2})$$

Notice that (K2) subsumes directly (K), since one can check that they coincide when  $\mu, \nu$  are discrete measures, with respective probability weights  $\mathbf{a}, \mathbf{b}$  and locations  $(\mathbf{x}_1, \dots, \mathbf{x}_n)$  and  $(\mathbf{y}_1, \dots, \mathbf{y}_m)$ .

1  
2 **6.10.2.4 Wasserstein distances**

3 When  $c$  is equal to a metric  $d$  exponentiated by an integer, the optimal value of the Kantorovich  
4 problem is called the Wasserstein *distance* between  $\mu$  and  $\nu$ :  
5

6 
$$W_p(\mu, \nu) \triangleq \left( \inf_{\pi \in \Pi(\mu, \nu)} \int_{\mathcal{X}^2} d(\mathbf{x}, \mathbf{y})^p \, d\pi(\mathbf{x}, \mathbf{y}) \right)^{1/p}. \quad (6.241)$$
  
7  
8

9 While the symmetry and the fact that  $W_p(\mu, \nu) = 0 \Rightarrow \mu = \nu$  are relatively easy to prove provided  $d$   
10 is a metric, proving the triangle inequality is slightly more challenging, and builds on a result known  
11 as the gluing lemma ([Vil08, p.23]). The  $p$ -th power of  $W_p(\mu, \nu)$  is often abbreviated as  $W_p^p(\mu, \nu)$ .  
12

13  
14 **6.10.3 Solving optimal transport**

15 **6.10.3.1 Duality and cost concavity**

16 Both (K) and (K2) are linear programs: their constraints and objective functions only involve  
17 summations. In that sense they admit a dual formulation (here, again, (DK2) subsumes (DK)):

19  
20 
$$\max_{\mathbf{f} \in \mathbb{R}^n, \mathbf{g} \in \mathbb{R}^m} \mathbf{f}^T \mathbf{a} + \mathbf{g}^T \mathbf{b} \quad (\text{DK})$$
  
21 
$$\mathbf{f} \oplus \mathbf{g} \leq C$$

22  
23 
$$\sup_{f \oplus g \leq c} \int_{\mathcal{X}} f \, d\mu + \int_{\mathcal{X}} g \, d\nu; \quad (\text{DK2})$$
  
24

25 where the sign  $\oplus$  denotes tensor addition for vectors,  $\mathbf{f} \oplus \mathbf{g} = [\mathbf{f}_i + \mathbf{g}_j]_{ij}$ , or functions,  $f \oplus g : \mathbf{x}, \mathbf{y} \mapsto$   
26  $f(\mathbf{x}) + g(\mathbf{y})$ . In other words, the dual problem looks for a pair of vectors (or functions) that attain  
27 the highest possible expectation when summed against  $\mathbf{a}$  and  $\mathbf{b}$  (or integrated against  $\mu, \nu$ ), pending  
28 the constraint that they do not differ too much across points  $\mathbf{x}, \mathbf{y}$ , as measured by  $c$ .

29 The dual problems in (K) and (K2) have two variables. Focusing on the continuous formulation, a  
30 closer inspection shows that it is possible, given a function  $f$  for the first measure, to compute the  
31 best possible candidate for function  $g$ . That function  $g$  should be as large as possible, yet satisfy the  
32 constraint that  $g(\mathbf{y}) \leq c(\mathbf{x}, \mathbf{y}) - f(\mathbf{x})$  for all  $\mathbf{x}, \mathbf{y}$ , making

33  
34 
$$\forall \mathbf{y} \in \mathcal{X}, \bar{f}(\mathbf{y}) \triangleq \inf_{\mathbf{x}} c(\mathbf{x}, \mathbf{y}) - f(\mathbf{x}), \quad (6.242)$$
  
35

36 the optimal choice.  $\bar{f}$  is called the  $c$ -transform of  $f$ . Naturally, one may choose to start instead from  
37  $g$ , to define an alternative  $c$ -transform:  
38

39  
40 
$$\forall \mathbf{x} \in \mathcal{X}, \tilde{g}(\mathbf{x}) \triangleq \inf_{\mathbf{y}} c(\mathbf{x}, \mathbf{y}) - g(\mathbf{y}). \quad (6.243)$$

41 Since these transformations can only improve solutions, one may even think of applying alternatively  
42 these transformations to an arbitrary  $f$ , to define  $\bar{f}$ ,  $\tilde{\bar{f}}$  and so on. One can show, however, that this  
43 has little interest, since  
44

45  
46 
$$\tilde{\bar{f}} = \bar{f}. \quad (6.244)$$
  
47

This remark allows, nonetheless, to narrow down the set of candidate functions to those that have already undergone such transformations. This reasoning yields the so-called set of  $c$ -concave functions,  $\mathcal{F}_c \triangleq \{f \mid \exists g : \mathcal{X} \rightarrow \mathbb{R}, f = \tilde{g}\}$ , which can be shown, equivalently, to be the set of functions  $f$  such that  $f = \bar{f}$ . One can therefore focus our attention to  $c$ -concave functions to solve (DK2) using a so-called semi-dual formulation,

$$\sup_{f \in \mathcal{F}_c} \int_{\mathcal{X}} f \, d\mu + \int_{\mathcal{X}} \bar{f} \, d\nu. \quad (\text{DK2})$$

Going from (DK2) to (DK2), we have removed a dual variable  $g$  and narrowed down the feasible set to  $\mathcal{F}_c$ , at the cost of introducing the highly non-linear transform  $\bar{f}$ . This reformulation is, however, very useful, in the sense that it allows to restrict our attention on  $c$ -concave functions, notably for two important classes of cost functions  $c$ : distances and squared-Euclidean norms.

### 6.10.3.2 Kantorovich-Rubinstein duality and Lipschitz potentials

A striking result illustrating the interest of  $c$ -concavity is provided when  $c$  is a metric  $d$ , namely when  $p = 1$  in (6.241). In that case, one can prove (exploiting notably the triangle inequality of the  $d$ ) that a  $d$ -concave function  $f$  is 1-Lipschitz (one has  $|f(\mathbf{x}) - f(\mathbf{y})| \leq d(\mathbf{x}, \mathbf{y})$  for any  $\mathbf{x}, \mathbf{y}$ ) and such that  $\bar{f} = -f$ . This result translates therefore in the following identity:

$$W_1(\mu, \nu) = \sup_{f \text{1-Lipschitz}} \int_{\mathcal{X}} f \, (d\mu - d\nu). \quad (6.245)$$

This result has numerous practical applications. This supremum over 1-Lipschitz functions can be efficiently approximated using Wavelet coefficients of densities in low-dimensions [SJ08], or heuristically in more general cases by training neural networks parameterized to be 1-Lipschitz [ACB17] using ReLU activation functions, and bounds on the entries of the weight matrices.

### 6.10.3.3 Monge maps as gradients of convex functions: the Brenier theorem

Another application of  $c$ -concavity lies in the case  $c(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|^2$ , which corresponds, up to the factor  $\frac{1}{2}$ , to the squared  $W_2$  distance used between densities in an Euclidean space. The remarkable result, shown first by [Bre91], is that the Monge map solving (M) between two measures for that cost (taken for granted  $\mu$  is regular enough, here assumed to have a density w.r.t. Lebesgue measure) exists and is necessarily the gradient of a convex function. In loose terms, one can show that

$$T^* = \arg \min_{T: T_\# \mu = \nu} \int_{\mathcal{X}} \frac{1}{2} \|\mathbf{x} - T(\mathbf{x})\|_2^2 \mu(d\mathbf{x}). \quad (\text{M})$$

exists, and is the gradient of a convex function  $u : \mathbb{R}^d \rightarrow \mathbb{R}$ , namely  $T^* = \nabla u$ . Conversely, for any convex function  $u$ , the optimal transport map between  $\mu$  and the displacement  $\nabla u_\# \mu$  is necessarily equal to  $\nabla u$ .

We provide a sketch of the proof: one can always exploit, for any reasonable cost function  $c$  (e.g. lower bounded and lower semi continuous), primal-dual relationships: Consider an optimal coupling  $P^*$  for (K2), as well as an optimal  $c$ -concave dual function  $f^*$  for (DK2). This implies in particular that  $(f^*, g^* = \bar{f}^*)$  is optimal for (DK2). Complementary slackness conditions for this pair

1 of linear programs imply that if  $\mathbf{x}_0, \mathbf{y}_0$  is in the support of  $P^*$ , then necessarily (and sufficiently)  
2  $f^*(\mathbf{x}_0) + \bar{f}^*(\mathbf{y}_0) = c(\mathbf{x}_0, \mathbf{y}_0)$ . Suppose therefore that  $\mathbf{x}_0, \mathbf{y}_0$  is indeed in the support of  $P^*$ . From the  
3 equality  $f^*(\mathbf{x}_0) + \bar{f}^*(\mathbf{y}_0) = c(\mathbf{x}_0, \mathbf{y}_0)$  one can trivially obtain that  $\bar{f}^*(\mathbf{y}_0) = c(\mathbf{x}_0, \mathbf{y}_0) - f^*(\mathbf{x}_0)$ . Yet,  
4 recall also that, by definition,  $\bar{f}^*(\mathbf{y}_0) = \inf_{\mathbf{x}} c(\mathbf{x}, \mathbf{y}_0) - f^*(\mathbf{x})$ . Therefore,  $\mathbf{x}_0$  has the special property  
5 that it minimizes  $\mathbf{x} \rightarrow c(\mathbf{x}, \mathbf{y}_0) - f^*(\mathbf{x})$ . If, at this point, one recalls that  $c$  is assumed in this section  
6 to be  $c(\mathbf{x}, \mathbf{y}) = \frac{1}{2}\|\mathbf{x} - \mathbf{y}\|^2$ , one has therefore that  $\mathbf{x}_0$  verifies

$$\mathbf{x}_0 \in \operatorname{argmin}_{\mathbf{x}} \frac{1}{2}\|\mathbf{x} - \mathbf{y}_0\|^2 - f^*(\mathbf{x}). \quad (6.246)$$

10 Assuming  $f^*$  is differentiable, which one can prove by  $c$ -concavity, this yields the identity

$$\mathbf{y}_0 - \mathbf{x}_0 - \nabla f^*(\mathbf{x}_0) = 0 \Rightarrow \mathbf{y}_0 = \mathbf{x}_0 - \nabla f^*(\mathbf{x}_0) = \nabla\left(\frac{1}{2}\|\cdot\|^2 - f^*\right)(\mathbf{x}_0). \quad (6.247)$$

14 Therefore, if  $(\mathbf{x}_0, \mathbf{y}_0)$  is in the support of  $P^*$ ,  $\mathbf{y}_0$  is uniquely determined, which proves  $P^*$  is in fact a  
15 Monge map “disguised” as a coupling, namely

$$\mathbf{P}^* = (\operatorname{Id}, \nabla\left(\frac{1}{2}\|\cdot\|^2 - f^*\right))_{\sharp}\mu. \quad (6.248)$$

19 The end of the proof can be worked out as follows: For any function  $h : \mathcal{X} \rightarrow \mathbf{R}$ , one can show, using  
20 the definitions of  $c$ -transforms and the Legendre transform, that  $\frac{1}{2}\|\cdot\|^2 - h$  is convex if and only if  $h$   
21 is  $c$ -concave. An intermediate step in that proof relies on showing that  $\frac{1}{2}\|\cdot\|^2 - \bar{h}$  is equal to the  
22 Legendre transform of  $\frac{1}{2}\|\cdot\|^2 - h$ . The function  $\frac{1}{2}\|\cdot\|^2 - f^*$  above is therefore convex, by  $c$ -concavity  
23 of  $f^*$ , and the optimal transport map is itself the gradient of a convex function.

24 Knowing that an optimal transport map for the squared-Euclidean cost is necessarily the gradient  
25 of a convex function can prove very useful to solve (DK2). Indeed, this knowledge can be leveraged  
26 to restrict estimation to relevant families of functions, namely gradients of input-convex neural  
27 networks[AJK17], as proposed in [Mak+20] or [Kor+20], as well as arbitrary convex functions with  
28 desirable smoothness and strong-convexity constants [PdC20].

29

### 30 6.10.3.4 Closed forms for univariate and Gaussian distributions

31 Many metrics between probability distributions have closed form expressions for simple cases. The  
32 Wasserstein distance is no exception, and can be computed in close form in two important scenarios.  
33 When distributions are univariate and the cost  $c(\mathbf{x}, \mathbf{y})$  is either a convex function of the difference  
34  $\mathbf{x} - \mathbf{y}$ , or when  $\partial c / \partial \mathbf{x} \partial \mathbf{y} < 0$  a.e., then the Wasserstein distance is essentially a comparison between  
35 the quantile functions of  $\mu$  and  $\nu$ . Recall that for a measure  $\rho$ , its quantile function  $Q_\rho$  is a function  
36 that takes values in  $[0, 1]$  and is valued in the support of  $\rho$ , and corresponds to the (generalized)  
37 inverse map of  $F_\rho$ , the cumulative distribution function (cdf) of  $\rho$ . With these notations, one has  
38 that

$$\text{OT}_c(\mu, \nu) = \int_{[0,1]} c(Q_\mu(u), Q_\nu(u)) \, du \quad (6.249)$$

43 In particular, when  $c$  is  $\mathbf{x}, \mathbf{y} \mapsto |\mathbf{x} - \mathbf{y}|$  then  $\text{OT}_c(\mu, \nu)$  corresponds to the Kolmogorov-Smirnov  
44 statistic, namely the area between the cdf of  $\mu$  and that of  $\nu$ . If  $c$  is  $\mathbf{x}, \mathbf{y} \mapsto (\mathbf{x} - \mathbf{y})^2$ , we recover  
45 simply the squared-Euclidean norm between the quantile functions of  $\mu$  and  $\nu$ . Note finally that the  
46 Monge map is also available in closed form, and is equal to  $Q_\nu \circ F_\mu$ .

47

The second closed form applies to so-called elliptically contoured distributions, chiefly among them Gaussian multivariate distributions [Gel90]. For two Gaussians  $\mathcal{N}(\mathbf{m}_1, \Sigma_1)$  and  $\mathcal{N}(\mathbf{m}_2, \Sigma_2)$  their 2-Wasserstein distance decomposes as

$$W_2^2(\mathcal{N}(\mathbf{m}_1, \Sigma_1), \mathcal{N}(\mathbf{m}_2, \Sigma_2)) = \|\mathbf{m}_1 - \mathbf{m}_2\|^2 + \mathcal{B}^2(\Sigma_1, \Sigma_2) \quad (6.250)$$

where the Bures metric  $\mathcal{B}$  reads:

$$\mathcal{B}^2(\Sigma_1, \Sigma_2) = \text{trace} \left( \Sigma_1 + \Sigma_2 - 2 \left( \Sigma_1^{\frac{1}{2}} \Sigma_2 \Sigma_1^{\frac{1}{2}} \right)^{\frac{1}{2}} \right). \quad (6.251)$$

Notice in particular that these quantities are well-defined even when the covariance matrices are not invertible, and that they collapse to the distance between means as both covariances become 0. When the first covariance matrix is invertible, one has that the optimal Monge map is given by

$$T \triangleq \mathbf{x} \mapsto A(\mathbf{x} - \mathbf{m}_1) + \mathbf{m}_2, \text{ where } A \triangleq \Sigma_1^{-\frac{1}{2}} \left( \Sigma_1^{\frac{1}{2}} \Sigma_2 \Sigma_1^{\frac{1}{2}} \right)^{\frac{1}{2}} \Sigma_1^{-\frac{1}{2}} \quad (6.252)$$

It is easy to show that  $T^*$  is indeed optimal: The fact that  $T_{\sharp}\mathcal{N}(\mathbf{m}_1, \Sigma_1) = \mathcal{N}(\mathbf{m}_2, \Sigma_2)$  follows from the knowledge that the affine push-forward of a Gaussian is another Gaussian. Here  $T$  is designed to push precisely the first Gaussian onto the second (and  $A$  designed to recover random variables with variance  $\Sigma_2$  when starting from random variables with variance  $\Sigma_1$ ). The optimality of  $T$  can be recovered by simply noticing that is the gradient of a convex quadratic form, since  $A$  is positive definite, and closing this proof using the Brenier theorem above.

### 6.10.3.5 Exact evaluation using linear program solvers

We have hinted, using duality and  $c$ -concavity, that methods based on stochastic optimization over 1-Lipschitz or convex neural networks can be employed to estimate Wasserstein distances when  $c$  is the Euclidean distance or its square. These approaches are, however, non-convex and can only reach local optima. Apart from these two cases, and the closed forms provided above, the only reliable approach to compute Wasserstein distances appears when both  $\mu$  and  $\nu$  are discrete measures: In that case, one can instantiate and solve the discrete (K) problem, or its dual (DK) formulation. The primal problem is a canonical example of network flow problems, and can be solved with the network-simplex method in  $O(nm(n+m)\log(n+m))$  complexity [AMO88], or, alternatively, with the comparable auction algorithm [BC89]. These approaches suffer from computational limitations: their cubic cost is intractable for large scale scenarios; their combinatorial flavor makes it harder to solve to parallelize simultaneously the computation of multiple optimal transport problems with a common cost matrix  $C$ .

An altogether different issue, arising from statistics, should discourage further users from using these LP formulations, notably in high-dimensional settings. Indeed, the bottleneck practitioners will most likely encounter when using (K) is that, in most scenarios, their goal will be to approximate the distance between two continuous measures  $\mu, \nu$  using only i.i.d samples contained in empirical measures  $\hat{\mu}_n, \hat{\nu}_n$ . Using (K) to approximate the corresponding (K2) is doomed to fail, as various results [FG15] have shown in relevant settings (notably for measures in  $\mathbb{R}^d$ ) that the *sample complexity*

of the estimator provided by (K) to approximate (K2) is of order  $1/n^{1/q}$ . In other words, the gap between  $W_2(\mu, \nu)$  and  $W_2(\hat{\mu}_n, \hat{\nu}_n)$  is large on expectation, decreases extremely slowly as  $n$  increases in high dimensions, and solving exactly (K2) between these samples is compute power that is mostly wasted on overfitting. To address this curse of dimensionality, it is therefore extremely important in practice to approach (K2) using a more careful strategy, one that involves regularizations that can leverage prior assumptions on  $\mu$  and  $\nu$ . While all approaches outlined above using neural networks can be interpreted under this light, we focus in the following on a specific approach that results in a convex problem that is relatively simple to implement, embarrassingly parallel and with quadratic complexity.

11

### 6.10.3.6 Obtaining smoothness using entropic regularization

A computational approach to speed-up the resolution of (K) was proposed in [Cut13], building on earlier contributions [Wil69; KY94] and a filiation to the Schrödinger bridge problem in the special case where  $c = d^2$  [Léo14]. The idea rests upon regularizing the transportation cost by the Kullback-Leibler divergence of the coupling to the product measure of  $\mu, \nu$ ,

$$W_{c,\gamma}(\mu, \nu) \triangleq \inf_{\pi \in \Pi(\mu, \nu)} \int_{\mathcal{X}^2} d(\mathbf{x}, \mathbf{y})^p d\pi(\mathbf{x}, \mathbf{y}) + \gamma D_{\text{KL}}(\pi \| \mu \otimes \nu). \quad (6.253)$$

20

When instantiated on discrete measures, this problem is equivalent to the following  $\gamma$ -strongly convex problem on the set of transportation matrices (which should be compared to (K))

$$\text{OT}_{C,\gamma}(\mathbf{a}, \mathbf{b}) = \min_{P \in \mathbb{R}_+^{n \times m}, P\mathbf{1}_m = \mathbf{a}, P^T\mathbf{1}_m = \mathbf{b}} \langle P, C \rangle \triangleq \sum_{i,j} P_{ij} C_{ij} - \gamma \mathbb{H}(P) + \gamma (\mathbb{H}(\mathbf{a}) + \mathbb{H}(\mathbf{b})), \quad (6.254)$$

Where and is itself equivalent to the following dual problem (which should be compared to (DK))

$$\text{OT}_{C,\gamma}(\mathbf{a}, \mathbf{b}) = \max_{\mathbf{f} \in \mathbb{R}^n, \mathbf{g} \in \mathbb{R}^m} \mathbf{f}^T \mathbf{a} + \mathbf{g}^T \mathbf{b} - \gamma (e^{\mathbf{f}/\gamma})^T K e^{\mathbf{g}/\gamma} + \gamma (1 + \mathbb{H}(\mathbf{a}) + \mathbb{H}(\mathbf{b})) \quad (6.255)$$

and  $K \triangleq e^{-C/\gamma}$  is the elementwise exponential of  $-C/\gamma$ . This regularization has several benefits. Primal-dual relationships show an explicit link between the (unique) solution  $P_\gamma^*$  and a pair of optimal dual variables  $(\mathbf{f}^*, \mathbf{g}^*)$  as

$$P_\gamma^* = \text{diag}(e^{\mathbf{f}/\gamma}) K \text{diag}(e^{\mathbf{g}/\gamma}) \quad (6.256)$$

Problem(6.255) can be solved using a fairly simple strategy that has proved very sturdy in practice: a simple block-coordinate ascent (optimizing alternatively the objective in  $\mathbf{f}$  and then  $\mathbf{g}$ ), resulting in the famous Sinkhorn algorithm [Sin67], here expressed with log-sum-exp updates, starting from an arbitrary initialization for  $\mathbf{g}$ , to carry out these two updates sequentially, until they converge:

$$\mathbf{f} \leftarrow \gamma \log \mathbf{a} - \gamma \log K e^{\mathbf{g}/\gamma} \quad \mathbf{g} \leftarrow \gamma \log \mathbf{b} - \gamma \log K^T e^{\mathbf{f}/\gamma} \quad (6.257)$$

The convergence of this algorithm has been amply studied (see [CK21] and references therein). Convergence is naturally slower as  $\gamma$  decreases, reflecting the hardness of approaching LP solutions, as studied in [AWR17]. This regularization also has statistical benefits since, as argued in [Gen+19], the sample complexity of the regularized Wasserstein distance improves to a  $O(1/\sqrt{n})$  regime, with, however, a constant in  $1/\gamma^{q/2}$  that deteriorates as dimension grows.

47

---

## 6.11 Submodular optimization

*This section was written by Jeff Bilmes.*

This section provides a brief overview of submodularity in machine learning.<sup>5</sup> Submodularity has an extremely simple definition. However, the “simplest things are often the most complicated to understand fully” [Sam74], and while submodularity has been studied extensively over the years, it continues to yield new and surprising insights and properties, some of which are extremely relevant to data science, machine learning, and artificial intelligence. A submodular function operates on subsets of some finite *ground set*,  $V$ . Finding a guaranteed good subset of  $V$  would ordinarily require an amount of computation exponential in the size of  $V$ . Submodular functions, however, have certain properties that make optimization either tractable or approximable where otherwise neither would be possible. The properties are quite natural, however, so submodular functions are both flexible and widely applicable to real problems. Submodularity involves an intuitive and natural diminishing returns property, stating that adding an element to a smaller set helps more than adding it to a larger set. Like convexity, submodularity allows one to efficiently find provably optimal or near-optimal solutions. In contrast to convexity, however, where little regarding maximization is guaranteed, submodular functions can be both minimized and (approximately) maximized. Submodular maximization and minimization, however, require very different algorithmic solutions and have quite different applications. It is sometimes said that submodular functions are a discrete form of convexity. This is not quite true, as submodular functions are like both convex and concave functions, but also have properties that are similar simultaneously to both convex and concave functions at the same time, but then some properties of submodularity are like neither convexity nor concavity. Convexity and concavity, for example, can be conveyed even as univariate functions. This is impossible for submodularity, as submodular functions are defined based only on the response of the function to changes amongst different variables in a multidimensional discrete space.

### 6.11.1 Intuition, Examples, and Background

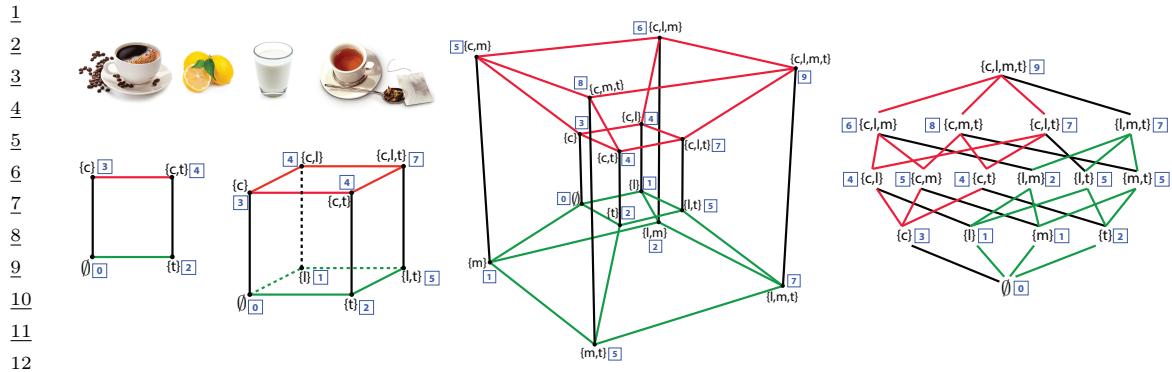
Let us define a *set function*  $f : 2^V \rightarrow \mathbb{R}$  as one that assigns a value to every subset of  $V$ . The notation  $2^V$  is the power set of  $V$ , and has size  $2^{|V|}$  which means that  $f$  lives in space  $\mathbb{R}^{2^n}$  — i.e., since there are  $2^n$  possible subsets of  $V$ ,  $f$  can return  $2^n$  distinct values. We use the notation  $X + v$  as shorthand for  $X \cup \{v\}$ . Also, the value of an element in a given context is so widely used a concept, we have a special notation for it — the incremental value *gain* of  $v$  in the context if  $X$  is defined as  $f(v|X) = f(X + v) - f(X)$ . Thus, while  $f(v)$  is the value of element  $v$ ,  $f(v|X)$  is the value of element  $v$  if you already have  $X$ . We also define the gain of set  $X$  in the context of  $Y$  as  $f(X|Y) = f(X \cup Y) - f(Y)$ .

#### 6.11.1.1 Coffee, Lemon, Milk, Tea

As a simple example, will explore the manner in which the value of everyday items may interact and combine, namely coffee, lemon, milk, and tea. Consider the value relationships amongst the four items coffee ( $c$ ), lemon ( $l$ ), milk ( $m$ ), and tea ( $t$ ) as shown in Figure 6.19. Suppose you just woke up,

---

<sup>5</sup> 5. A greatly extended version of the material in this section may be found at [Bil22].



13 *Figure 6.19: The value relationships between coffee ( $c$ ), lemon ( $l$ ), milk ( $m$ ), and tea ( $t$ ). On the left, we first*  
 14 *see a simple square showing the relationships between coffee and tea, and see that they are substitutive (or*  
 15 *submodular). In this, and all of the shapes, the vertex label set is indicated in curly braces and the value at*  
 16 *that vertex is a blue integer in a box. We next see a three-dimensional cube that adds lemon to coffee and tea*  
 17 *set. We see that tea and lemon are complementary (supermodular) but coffee and lemon are additive (modular,*  
 18 *or independent). We next see a four-dimensional hypercube (tesseract) showing all of the value relationships*  
 19 *described in the text. The four-dimensional hypercube is also shown as a lattice (on the right) showing the*  
 20 *same relationships as well as two (red and green, also shown in the tesseract) of the eight three-dimensional*  
 21 *cubes contained within.*

22  
 23  
 24 and there is a function  $f : 2^V \rightarrow \mathbb{R}$  that provides the average valuation for any subset of the items in  
 25  $V = \{c, l, m, t\}$ . You can think of this function as giving the average price a typical person  
 26 would be willing to pay for any subset of items. Since nothing should cost nothing, we would expect  
 27 that  $f(\emptyset) = 0$ . Clearly, one needs either coffee or tea in the morning, so  $f(c) > 0$  and  $f(t) > 0$ , and  
 28 coffee is usually more expensive than tea, so that  $f(c) > f(t)$  pound for pound. Also more items cost  
 29 more, so that, for example,  $0 < f(c) < f(c, m) < f(c, m, t) < f(c, l, m, t)$ . Thus, the function  $f$  is  
 30 strictly *monotone*, or  $f(X) < f(Y)$  whenever  $X \subset Y$ .

31 The next thing we note is that coffee and tea may substitute for each other – they both have  
 32 the same effect, waking you up. They are mutually redundant, and they decrease each other's  
 33 value since once you have had a cup of coffee, a cup of tea is less necessary and less desirable. Thus,  
 34  $f(c, t) < f(c) + f(t)$ , which is known as a *subadditive* relationship, the whole is less than the sum  
 35 of the parts. On the other hand, some items complement each other. For example, milk and coffee  
 36 are better combined together than when both are considered in isolation, or  $f(m, c) > f(m) + f(c)$ ,  
 37 a *superadditive* relationship, the whole is more than the sum of the parts. A few of the items  
 38 do not affect each others' price. For example, lemon and milk cost the same together as apart, so  
 39  $f(l, m) = f(l) + f(m)$ , an *additive* or *modular* relationship — such a relationship is perhaps midway  
 40 between a subadditive and a superadditive relationship, and can be seen as a form of independence.

41 Things become more interesting when we consider three or more items together. For example,  
 42 once you have tea, lemon becomes less valuable when you acquire milk since there might be those  
 43 that prefer milk to lemon in their tea. Similarly, milk becomes less valuable once you have acquired  
 44 lemon since there are those who prefer lemon in their tea to milk. So, once you have tea, lemon and  
 45 milk are substitutive, you would never use both as the lemon would only curdle the milk. These  
 46 are *submodular* relationships,  $f(l|m, t) < f(l|t)$  and  $f(m|l, t) < f(m|t)$  each of which implies that  
 47

$f(l, t) + f(m, t) > f(l, m, t) + f(t)$ . The value of lemon (respectively milk) with tea decreases in the larger context of having milk (respectively lemon) with tea, typical of submodular relationships.

Not all of the items are in a submodular relationship, as sometimes the presence of an item can increase the value of another item. For example, once you have milk, then tea becomes still more valuable when you also acquire lemon, since tea with the choice of either lemon or milk is more valuable than tea with the option only of milk. Similarly, once you have milk, lemon becomes more valuable when you acquire tea, since lemon with milk alone is not nearly as valuable as lemon with tea, even if milk is at hand. This means that  $f(t|l, m) > f(t|m)$  and  $f(l|t, m) > f(l|m)$  implying  $f(l, m) + f(m, t) < f(l, m, t) + f(m)$ . These are known as *supermodular* relationships, where the value increases as the context increases.

We have asked for a set of relationships amongst various subsets of the four items  $V = \{c, l, m, t\}$ , Is there a function that offers a value to each  $X \subseteq V$  that satisfies all of the above relationships? Figure 6.19 in fact shows such a function. On the left, we see a two-dimensional square whose vertices indicate the values over subsets of  $\{c, t\}$  and we can quickly verify that the sum of the blue boxes on north-west (corresponding to  $f(\{c\})$ ) and south-east corners (corresponding to  $f(\{t\})$ ) is greater than the sum of the north-east and south-west corners, expressing the required submodular relationship. Next on the right is a three-dimensional cube that adds the relationship with lemon. Now we have six squares and we see that the values at each of the vertices all satisfy the above requirements — we verify this by considering the valuations at the four corners of every one of the six faces of the cube. Since  $|V| = 4$  we need a four-dimensional hypercube to show all values and this may be shown in two ways. It is first shown as a tesseract, a well-known three-dimensional projection of a four-dimensional hypercube. In the figure, all vertices are labeled both with subsets of  $V$  as well as the function value  $f(X)$  as the blue number in a box. The figure on the right shows a *lattice* version of the four-dimensional hypercube, where corresponding three-dimensional cubes are shown in green and red.

We thus see that a set function is defined for all subsets of a ground set, and that they correspond to valuations at all vertices of the hypercube. For the particular function over valuations of subsets of coffee, lemon, milk, and tea, we have seen submodular, supermodular, and modular relationships all in one function. Therefore, the overall function  $f$  defined in Figure 6.19 is neither submodular, supermodular, nor modular. For combinatorial auctions, there is often a desire to have a diversity of such manners of relationships [LLN06] — representation of these relationships can be handled by a difference of submodular functions [NB05; IB12] or a sum of a submodular and supermodular function [BB18] (further described below). In machine learning, however, most of the time we are interested in functions that are submodular (or modular, or supermodular) everywhere.

### 6.11.2 Submodular Basic Definitions

For a function to be submodular, it must satisfy the submodular relationship for all subsets. We arrive at the following definition.

**Definition 6.11.1 (Submodular Function).** A given set function  $f : 2^V \rightarrow \mathbb{R}$  is submodular if for all  $X, Y \subseteq V$ , we have the following inequality:

$$f(X) + f(Y) \geq f(X \cup Y) + f(X \cap Y) \quad (6.258)$$

1 There are also many other equivalent definitions of submodularity [Bil22] some of which are more  
 2 intuitive and easier to understand. For example, submodular functions are those set functions that  
 3 satisfy the property of diminishing returns. If we think of a function  $f(X)$  as measuring the value of  
 4 a set  $X$  that is a subset of a larger set of data items  $X \subseteq V$ , then the submodular property means  
 5 that the incremental “value” of adding a data item  $v$  to set  $X$  decreases as the size of  $X$  grows. This  
 6 gives us a second classic definition of submodularity.  
 7

8 **Definition 6.11.2** (Submodular Function via Diminishing Returns). *A given set function  $f : 2^V \rightarrow \mathbb{R}$   
 9 is submodular if for all  $X, Y \subseteq V$ , where  $X \subseteq Y$  and for all  $v \notin Y$ , we have the following inequality:*  
 10

$$11 \quad f(X + v) - f(X) \geq f(Y + v) - f(Y) \quad (6.259)$$

12 The property that the incremental value of lemon with tea is less than the incremental value  
 13 of lemon once milk is already in the tea is equivalent to Equation 6.258 if we set  $X = \{m, t\}$  and  
 14  $Y = \{l, t\}$  (i.e.  $f(m, t) + f(l, t) > f(l, m, t) + f(t)$ ). It is naturally also equivalent to Equation 6.259  
 15 if we set  $X = \{t\}$ ,  $Y = \{m, t\}$ , and with  $v = l$  (i.e.,  $f(l|m, t) < f(l|t)$ ).

16 There are many functions that are submodular, one famous one being Shannon entropy seen  
 17 as a function of subsets of random variables. We first point out that there are non-negative (i.e.,  
 18  $f(A) \geq 0, \forall A$ ), monotone non-decreasing (i.e.,  $f(A) \leq f(B)$  whenever  $A \subseteq B$ ) submodular functions  
 19 that are not entropic [Yeu91b; ZY97; ZY98], so submodularity is not just a trivial restatement  
 20 of the class of entropy functions. When a function is monotone non-decreasing, submodular, and  
 21 *normalized* so that  $f(\emptyset) = 0$ , it is often referred to as a **polymatroid function**. Thus, while the  
 22 entropy function is a polymatroid function, it does not encompass all polymatroid functions even  
 23 though all polymatroid functions satisfy the properties Claude Shannon mentioned as being natural  
 24 for an “information” function (see Section 6.11.7).

25 A function  $f$  is supermodular if and only if  $-f$  is submodular. If a function is both submodular  
 26 and supermodular, it is known as a *modular* function. It is always the case that modular functions  
 27 may take the form of a vector-scalar pair  $(m, c)$  where  $m : 2^V \rightarrow \mathbb{R}$  and where  $c \in \mathbb{R}$  is a constant,  
 28 and where for any  $A \subseteq V$ , we have that  $m(A) = c + \sum_{v \in A} m_v$ . If the modular function is normalized,  
 29 so that  $m(\emptyset) = 0$ , then  $c = 0$  and the modular function can be seen simply as a vector  $m \in \mathbb{R}^V$ .  
 30 Hence, we sometimes say that the modular function  $x \in \mathbb{R}^V$  offers a value for set  $A$  as the partial  
 31 sum  $x(A) = \sum_{v \in A} x(v)$ . Many combinatorial problems use modular functions as objectives. For  
 32 example, the graph cut problem uses modular function defined over the edges, judges a cut in a  
 33 graph as the modular function applied to the edges that comprise the cut.

34 As can be seen from the above, and by considering Figure 6.19, a submodular function, and in  
 35 fact any set function,  $f : 2^V \rightarrow \mathbb{R}$  can be seen as a function defined only on the vertices of the  
 36  $n$ -dimensional unit hypercube  $[0, 1]^n$ . Given any set  $X \subseteq V$ , we define  $\mathbf{1}_X \in \{0, 1\}^V$  to be the  
 37 characteristic vector of set  $X$  defined as  $\mathbf{1}_X(v) = 1$  if  $v \in X$  and  $\mathbf{1}_X(v) = 0$  otherwise. This gives us  
 38 a way to map from any set  $X \subseteq V$  to a binary vector  $\mathbf{1}_X$ . We also see that  $\mathbf{1}_X$  is itself a modular  
 39 function since  $\mathbf{1}_X \in \{0, 1\}^V \subset \mathbb{R}^V$ .

40 Submodular functions share a number of properties in common with both convex and concave  
 41 functions [Lov83], including wide applicability, generality, multiple representations, and closure  
 42 under a number of common operators (including mixtures, truncation, complementation, and certain  
 43 convolutions). There is one important submodular closure property that we state here — that is that  
 44 if we take a non-negative weighted (or conical) combinations of submodular functions, we preserve  
 45 submodularity. In other words, if we have a set of  $k$  submodular functions,  $f_i : 2^V \rightarrow \mathbb{R}$ ,  $i \in [k]$ , and  
 46

47

we form  $f(X) = \sum_{i=1}^k \omega_i f_i(X)$  where  $\omega_i \geq 0$  for all  $i$ , then Definition 6.11.1 immediately implies that  $f$  is also submodular. When we consider Definition 6.11.1, we see that submodular functions live in a cone in  $2^n$ -dimensional space defined by the intersection of an exponential number of half-spaces each one of which is defined by one of the inequalities of the form  $f(X) + f(Y) \geq f(X \cup Y) + f(X \cap Y)$ . Each submodular function is therefore a point in that cone. It is therefore not surprising that taking conical combinations of such points stays within this cone.

### 6.11.3 Example Submodular Functions

As mentioned above, there are many functions that are submodular besides entropy. Perhaps the simplest such function is  $f(A) = \sqrt{|A|}$  which is the composition of the square-root function (which is concave) with the cardinality  $|A|$  of the set  $A$ . The gain function is  $f(A + v) - f(A) = \sqrt{k+1} - \sqrt{k}$  if  $|A| = k$ , which we know to be a decreasing in  $k$ , thus establishing the submodularity of  $f$ . In fact, if  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  is any concave function, then  $f(A) = \phi(|A|)$  will be submodular for the same reason.<sup>6</sup> Generalizing this slightly further, a function defined as  $f(A) = \phi(\sum_{a \in A} m(a))$  is also submodular, whenever  $m(a) \geq 0$  for all  $a \in V$ . This yields a composition of a concave function with a modular function  $f(A) = \phi(m(A))$  since  $\sum_{a \in A} m(a) = m(A)$ . We may take sums of such functions as well as add a final modular function without losing submodularity, leading to  $f(A) = \sum_{u \in U} \phi_u(\sum_{a \in A} m_u(a)) + \sum_{a \in A} m_{\pm}(a)$  where  $\phi_u$  can be a distinct concave function for each  $u$ ,  $m_{u,a}$  is a non-negative real value for all  $u$  and  $a$ , and  $m_{\pm}(a)$  is an arbitrary real number. Therefore,  $f(A) = \sum_{u \in U} \phi_u(m_u(A)) + m_{\pm}(A)$  where  $m_u$  is a  $u$ -specific non-negative modular function and  $m_{\pm}$  is an arbitrary modular function. Such functions are sometimes known as **feature-based** submodular functions [BB17] because  $U$  can be a set of non-negative features (in the machine-learning “bag-of-words” sense) and this function measures a form of dispersion over  $A$  as determined by the set of features  $U$ .

A function such as  $f(A) = \sum_{u \in U} \phi_u(m_u(A))$  tends to award high diversity to a set  $A$  that has a high valuation by a distinct set of the features  $U$ . The reason is that, due to the concave nature of  $\phi_u$ , any addition to the argument  $m_u(A)$  by adding, say,  $v$  to  $A$  would diminish as  $A$  gets larger. In order to produce a set larger than  $A$  that has a much larger valuation, one must use a feature  $u' \neq u$  that has not yet diminished as much.

*Facility location* is another well-known submodular function — perhaps an appropriate nickname would be the “ $k$ -means of submodular functions,” due to its applicability, utility, ease-of-use (it needs only an affinity matrix), and similarity to  $k$ -medoids problems. The facility location function is defined using an affinity matrix as follows:  $f(A) = \sum_{v \in V} \max_{a \in A} \text{sim}(a, v)$  where  $\text{sim}(a, v)$  is a non-negative measure of the affinity (or similarity) between element  $a$  and  $v$ . Here, every element  $v \in V$  must have a representative within the set  $A$  and the representative for each  $v \in V$  is chosen to be the element  $a \in A$  most similar to  $v$ . This function is also a form of dispersion or diversity function because, in order to maximize it, every element  $v \in V$  must have some element similar to it in  $A$ . The overall score is then the sum of the similarity between each element  $v \in V$  and  $v$ 's representative. This function is monotone (since as  $A$  includes more elements to become  $B \supseteq A$ , it is possible only to find an element in  $B$  more similar to a given  $v$  than an element in  $A$ ).

While the facility location looks quite different from a feature based function, it is possible

6. While we will not be extensively discussing supermodular functions in this section,  $f(A) = \phi(|A|)$  is supermodular for any convex function  $\phi$ .

1 to precisely represent any facility location function with a feature based function. Consider  
2 just  $\max_{a \in A} x_a$  and, without loss of generality, assume that  $0 \leq x_1 \leq x_2 \leq \dots \leq x_n$ . Then  
3  $\max_{a \in A} x_a = \sum_{i=1}^n y_i \min(|A \cap \{i, i+1, \dots, n\}|, 1)$  where  $y_i = x_i - x_{i-1}$  and we set  $x_0 = 0$ . We note  
4 that this is a sum of weighted concave composed with modular functions since  $\min(\alpha, 1)$  is concave  
5 in  $\alpha$ , and  $|A \cap \{i, i+1, \dots, n\}|$  is a modular function in  $A$ . Thus, the facility location function, a  
6 sum of these, is merely a feature based function.  
7

8 Feature based functions, in fact, are quite expressive, and can be used to represent many different  
9 submodular functions including set cover and graph-based functions. For example, we can define a *set*  
10 *cover function*, given a set of sets  $\{U_v\}_{v \in V}$ , via  $f(X) = |\bigcup_{v \in X} U_v|$ . If  $f(X) = |U|$  where  $U = \bigcup_{v \in V} U_v$   
11 then  $X$  indexes a set that fully covers  $U$ . This can also be represented as  $f(X) = \sum_{u \in U} \min(1, m_u(X))$   
12 where  $m_u(X)$  is a modular function where  $m_u(v) = 1$  if and only if  $u \in U_v$  and otherwise  $m_u(v) = 0$ .  
13 We see that this is a feature based submodular function since  $\min(1, x)$  is concave in  $x$ , and  $U$  is a  
14 set of features.

15 This construct can be used to produce the vertex cover function if we set  $U = V$  to be the set of  
16 vertices in a graph, and set  $m_u(v) = 1$  if and only if vertices  $u$  and  $v$  are adjacent in the graph and  
17 otherwise set  $m_u(v) = 0$ . Similarly, the edge cover function can be expressed by setting  $V$  to be the  
18 set of edges in a graph,  $U$  to be the set of vertices in the graph, and  $m_u(v) = 1$  if and only edge  $v$  is  
19 incident to vertex  $u$ .

20 A generalization of the set cover function is the *probabilistic coverage* function. Let  $P[B_{u,v} = 1]$  be  
21 the probability of the presence of feature (or concept)  $u$  within element  $v$ . Here, we treat  $B_{u,v}$  as a  
22 Bernoulli random variable for each element  $v$  and feature  $u$  so that  $P[B_{u,v} = 1] = 1 - P[B_{u,v} = 0]$ .  
23 Then we can define the probabilistic coverage function as  $f(X) = \sum_{u \in U} f_u(X)$  where, for feature  
24  $u$ , we have  $f_u(X) = 1 - \prod_{v \in X} (1 - P[B_{u,v} = 1])$  which indicates the degree to which feature  $u$  is  
25 “covered” by  $X$ . If we set  $P[B_{u,v} = 1] = 1$  if and only if  $u \in U_v$  and otherwise  $P[B_{u,v} = 1] = 0$ , then  
26  $f_u(X) = \min(1, m_u(X))$  and the set cover function can be represented as  $\sum_{u \in U} f_u(X)$ . We can  
27 generalize this in two ways. First, to make it softer and more probabilistic we allow  $P[B_{u,v} = 1]$  to  
28 be any number between zero and one. We also allow each feature to have a non-negative weight. This  
29 yields the general form of the probabilistic coverage function, which is defined by taking a weighted  
30 combination over all features:  $f_u(X) = \sum_{v \in X} \omega_v f_u(X)$  where  $\omega_u \geq 0$  is a weight for feature  $u$ .  
31 Observe that  $1 - \prod_{v \in X} (1 - P[B_{u,v} = 1]) = 1 - \exp(-m_u(X)) = \phi(m_u(X))$  where  $m_u$  is a modular  
32 function with evaluation  $m_u(X) = \sum_{v \in X} \log(1/(1 - P[B_{u,v} = 1]))$  and for  $z \in \mathbb{R}$ ,  $\phi(z) = 1 - \exp(-z)$   
33 is a concave function. Thus, the probabilistic coverage function (and its set cover specialization) is  
34 also feature based function.

35 Another common submodular function is the graph cut function. Here, we measure the value of a  
36 subset of  $V$  by the edges that cross between a set of nodes and all but that set of nodes. We are given  
37 an undirected non-negative weighted graph  $\mathcal{G} = (V, E, w)$  where  $V$  is the set of nodes,  $E \subseteq V \times V$  is  
38 the set of edges, and  $w \in \mathbb{R}_+^E$  are non-negative edge weights corresponding to symmetric matrix (so  
39  $w_{i,j} = w_{j,i}$ ). For any  $e \in E$ , we have  $e = \{i, j\}$  for some  $i, j \in V$  with  $i \neq j$ , the graph cut function  
40  $f : 2^V \rightarrow \mathbb{R}$  is defined as  $f(X) = \sum_{i \in X, j \in \bar{X}} w_{i,j}$  where  $w_{i,j} \geq 0$  is the weight of edge  $e = \{i, j\}$   
41 ( $w_{i,j} = 0$  if the edge does not exist), and where  $\bar{X} = V \setminus X$  is the complement of set  $X$ . Notice that  
42

43

44

45

46

47

1 we can write the graph cut function as follows:

$$\begin{aligned} \underline{3} \quad f(X) &= \sum_{i \in X, j \in \bar{X}} w_{i,j} = \sum_{i,j \in V} w_{i,j} \mathbf{1}\{i \in X, j \in \bar{X}\} \end{aligned} \quad (6.260)$$

$$\begin{aligned} \underline{6} \quad &= \frac{1}{2} \sum_{i,j \in V} w_{i,j} \min(|X \cap \{i,j\}|, 1) + \frac{1}{2} \sum_{i,j \in V} w_{i,j} \min(|(V \setminus X) \cap \{i,j\}|, 1) - \frac{1}{2} \sum_{i,j \in V} w_{i,j} \end{aligned} \quad (6.261)$$

$$\begin{aligned} \underline{8} \quad &= \tilde{f}(X) + \tilde{f}(V \setminus X) - \tilde{f}(V) \end{aligned} \quad (6.262)$$

10 where  $\tilde{f}(X) = \frac{1}{2} \sum_{i,j \in V} w_{i,j} \min(|X \cap \{i,j\}|, 1)$ . Therefore, since  $\min(\alpha, 1)$  is concave, and since  
11  $m_{i,j}(X) = |X \cap \{i,j\}|$  is modular,  $\tilde{f}(X)$  is submodular for all  $i, j$ . Also, since  $\tilde{f}(X)$  is submodular, so  
12 is  $\tilde{f}(V \setminus X)$  (in  $X$ ). Therefore, the graph cut function can be expressed as a sum of non-normalized  
13 feature-based functions. Note that here the second modular function is not normalized and is  
14 non-increasing, and also we subtract the constant  $\tilde{f}(V)$  to achieve equality.

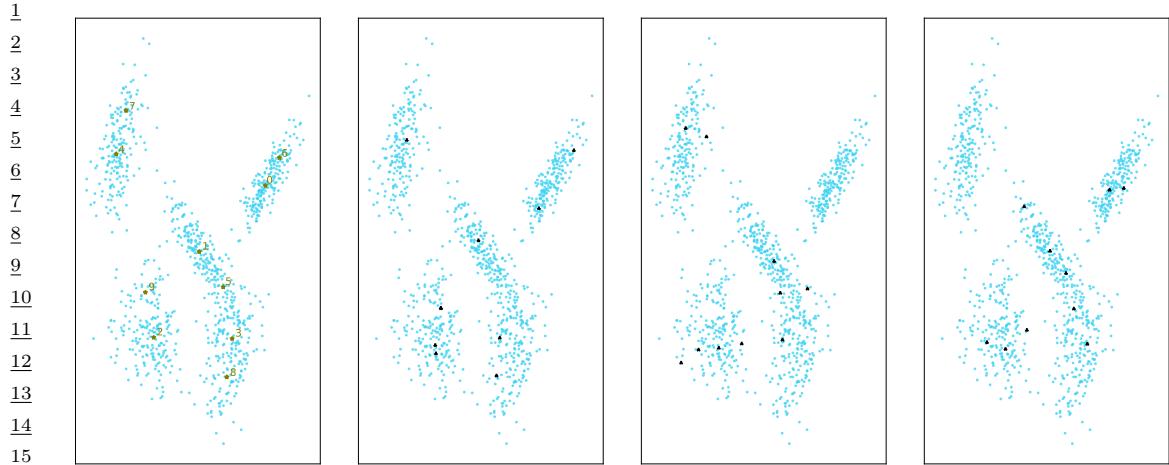
15 Another way to view the graph cut function is to consider the non-negative weights as a modular  
16 function defined over the edges. That is, we view  $w \in \mathbb{R}_+^E$  as a modular function  $w : 2^E \rightarrow \mathbb{R}_+$  where  
17 for every  $A \subseteq E$ ,  $w(A) = \sum_{e \in A} w(e)$  is the weight of the edges  $A$  where  $w(e)$  is the weight of edge  
18  $e$ . Then the graph cut function becomes  $f(X) = w(\{(a, b) \in E : a \in X, b \in X \setminus X\})$ . We view  
19  $\{(a, b) \in E : a \in X, b \in X \setminus X\}$  as a set-to-set mapping function, that maps subsets of nodes to  
20 subsets of edges, and the edge weight modular function  $w$  measures the weight of the resulting edges.  
21 This immediately suggests that other functions can measure the weight of the resulting edges as  
22 well, including non-modular functions. One example is to use a polymatroid function itself leading  
23  $h(X) = g(\{(a, b) \in E : a \in X, b \in X \setminus X\})$  where  $g : 2^E \rightarrow \mathbb{R}_+$  is a submodular function defined  
24 on subsets of edges. The function  $h$  is known as the **cooperative cut** function, and it is neither  
25 submodular nor supermodular in general but there are many useful and practical algorithms that  
26 can be used to optimize it [JB16] thanks to its internal yet exposed and thus available to exploit  
27 submodular structure.

28 While feature based functions are flexible and powerful, there is a strictly broader class of  
29 submodular functions, unable to be expressed by feature-based functions, and that are related to deep  
30 neural networks. Here, we create a recursively nested composition of concave functions with sums  
31 of compositions of concave functions. An example is  $f(A) = \phi(\sum_{u \in U} \omega_u \phi_u(\sum_{a \in A} m_{u,a}))$ , where  $\phi$   
32 is an outer concave function composed with a feature based function, with  $m_{u,a} \geq 0$  and  $\omega_u \geq 0$ .  
33 This is known as a two-layer **deep submodular function** (DSF). A three-layer DSF has the form  
34  $f(A) = \phi(\sum_{c \in C} \omega_c \phi_c(\sum_{u \in U} \omega_{u,c} \phi_u(\sum_{a \in A} m_{u,a})))$ . DSFs strictly expand the class of submodular  
35 functions beyond feature based functions, meaning that there are feature based functions that can  
36 not[BB17] represent deep submodular functions, even simple ones.

### 38 6.11.4 Submodular Optimization

40 Submodular functions, while discrete, would not be very useful if it was not possible to optimize  
41 over them efficiently. There are many natural problems in machine learning that can be cast as  
42 submodular optimization and that can be addressed relatively efficiently.

43 When one wishes to encourage diversity, information, spread, high complexity, independence,  
44 coverage, or dispersion, one usually will maximize a submodular function, in the form of  $\max_{A \in \mathcal{C}} f(A)$   
45 where  $\mathcal{C} \subseteq 2^V$  is a constraint set, a set of subsets we are willing to accept as feasible solutions (more  
46 on this below).



*Figure 6.20: Far Left: cardinality constrained (to ten) submodular maximization of a facility location function over 1000 points in two dimensions. Similarities are based on a Gaussian kernel  $\text{sim}(a, v) = \exp(-d(a, v))$  where  $d(\cdot, \cdot)$  is a distance. Selected points are green stars and the greedy order is also shown next to each selected point. Right three plots: different uniformly-at-random subsets of size ten.*

20

21

22

23 Why is submodularity, in general, a good model for diversity? Submodular functions are such  
24 that once you have some elements, any other elements not in your possession but that are similar  
25 to, explained by, or represented by the elements in your possession become less valuable. Thus, in  
26 order to maximize the function, one must choose other elements that are dissimilar to, or not well  
27 represented by, the ones you already have. That is, the elements similar to the ones you own are  
28 diminished in value relative to their original values, while the elements dissimilar to the ones you  
29 have do not have diminished value relative to their original values. Thus, maximizing a submodular  
30 function successfully involves choosing elements that are jointly dissimilar amongst each other, which  
31 is a definition of diversity. Diversity in general is a critically important aspect in machine learning  
32 and artificial intelligence. For example, bias in data science and machine learning can often be seen  
33 as some lack of diversity somewhere. Submodular functions have the potential to encourage (and  
34 even ensure) diversity, enhance balance, and reduce bias in artificial intelligence.

35 Note that in order for a submodular function to appropriately model diversity, it is important  
36 for it to be instantiated appropriately. Figure 6.20 shows an example in two dimensions. The plot  
37 compares the ten points chosen according to a facility location instantiated with a Gaussian kernel,  
38 along with the random samples of size ten. We see that the facility location selected points are more  
39 diverse and tend to cover the space much better than any of the randomly selected points each of  
40 which miss large regions of the space and/or show cases where points near each other are jointly  
41 selected.

42 When one wishes for homogeneity, conformity, low complexity, coherence, or cooperation, one will  
43 usually minimize a submodular function, in the form of  $\min_{A \in \mathcal{C}} f(A)$ . For example, if  $V$  is a set of  
44 pixels in an image, one might wish to choose a subset of pixels corresponding to a particular object  
45 over which the properties (i.e., color, luminance, texture) are relatively homogeneous. Finding a set  
46  $X$  of size  $k$ , even if  $k$  is large, need not have a large valuation  $f(X)$ , in fact it could even have the  
47

least valuation. Thus, semantic image segmentation could work even if the object being segmented and isolated consists of the majority of image pixels.

#### 6.11.4.1 Submodular Maximization

While the cardinality constrained submodular maximization problem is NP complete [Fei98], it was shown in [NWF78; FNW78] that the very simple and efficient greedy algorithm finds an approximate solution guaranteed to be within  $1 - 1/e \approx 0.63$  of the optimal solution. Moreover, the approximation ratio achieved by the simple greedy algorithm is provably the best achievable in polynomial time, assuming  $P \neq NP$  [Fei98]. The greedy algorithm proceeds as follows: Starting with  $X_0 = \emptyset$ , we repeat the following greedy step for  $i = 0 \dots (k-1)$ :

$$X_{i+1} = X_i \cup (\operatorname{argmax}_{v \in V \setminus X_i} f(X_i \cup \{v\})) \quad (6.263)$$

What the above approximation result means is that if  $X^* \in \operatorname{argmax}\{f(X) : |X| \leq k\}$ , and if  $\tilde{X}$  is the result of the greedy procedure, then  $f(\tilde{X}) \geq (1 - 1/e)f(X^*)$ .

The  $1 - 1/e$  guarantee is a powerful constant factor approximation result since it holds regardless of the size of the initial set  $V$  and regardless of which polymatroid function  $f$  is being optimized. It is possible to make this algorithm run extremely fast using various acceleration tricks [FNW78; NWF78; Min78].

A minor bit of additional information about a polymatroid function, however, can improve the approximation guarantee. Define the total curvature if the polymatroid function  $f$  as  $\kappa = 1 - \min_{v \in V} f(v|V-v)/f(v)$  where we assume  $f(v) > 0$  for all  $v$  (if not, we may prune them from the ground set since such elements can never improve a polymatroid function valuation). We thus have  $0 \leq \kappa \leq 1$ , and [CC84] showed that the greedy algorithm gives a guarantee of  $\frac{1}{\kappa}(1 - e^{-\kappa}) \geq 1 - 1/e$ . In fact, this is an equality (and we get the same bound) when  $\kappa = 1$ , which is the fully curved case. As  $\kappa$  gets smaller, the bound improves, until we reach the  $\kappa = 0$  case and the bound becomes unity. Observe that  $\kappa = 0$  if and only if the function is modular, in which case the greedy algorithm is optimal for the cardinality constrained maximization problem. In some cases, non-submodular functions can be decomposed into components that each might be more amenable to approximation. We see below that any set function can be written as a difference of submodular [NB05; IB12] functions, and sometimes (but not always) a given  $h$  can be composed into a monotone submodular plus a monotone supermodular function, or a BP function [BB18], i.e.,  $h = f + g$  where  $f$  is submodular and  $g$  is supermodular.  $g$  has an easily computed quantity called the supermodular curvature  $\kappa^g = 1 - \min_{v \in V} g(v)/g(v|V-v)$  that, together with the submodular curvature, can be used to produce an approximation ratio having the form  $\frac{1}{\kappa}(1 - e^{-\kappa(1-\kappa^g)})$  for greedily maximizing of  $h$ .

#### 6.11.4.2 Discrete Constraints

There are many other types of constraints one might desire besides a cardinality limitation. The next simplest constraint allows each element  $v$  to have a non-negative cost, say  $m(v) \in \mathbb{R}_+$ . In fact, this means that the costs are modular, i.e., the cost of any set  $X$  is  $m(X) = \sum_{v \in X} m(v)$ . A submodular maximization problem subject to a *knapsack constraint* then takes the form  $\max_{X \subseteq V : m(X) \leq b} f(X)$  where  $b$  is a non-negative budget. While the greedy algorithm does not solve this problem directly, a

<sup>1</sup> slightly modified cost-scaled version of the greedy algorithm [Svi04] does solve this problem for any  
<sup>2</sup> set of knapsack costs. This has been used for various multi-document summarization tasks [LB11;  
<sup>3</sup> LB12].

<sup>4</sup> There is no single direct analogy for a convex set when one is optimizing over subsets of the set  $V$ ,  
<sup>5</sup> but there are a few forms of discrete constraints that are both mathematically interesting and that  
<sup>6</sup> often occur repeatedly in applications.

<sup>7</sup> The first form are the independent subsets of a matroids. The independent sets of a matroid  
<sup>8</sup> are useful to represent a constraint set for submodular maximization [Cal+07; LSV09; Lee+10],  
<sup>9</sup>  $\max_{X \in \mathcal{I}} f(X)$ , and this can be useful in many ways. We can see this by showing a simple example  
<sup>10</sup> of what is known as a *partition matroid*. Consider a partition  $V = \{V_1, V_2, \dots, V_m\}$  of  $V$  into  $m$   
<sup>11</sup> mutually disjoint subsets that we call blocks. Suppose also that for each of the  $m$  blocks, there is a  
<sup>12</sup> positive integer limit  $\ell_i$  for  $i \in [m]$ . Consider next the set of sets formed by taking all subsets of  $V$   
<sup>13</sup> such that each subset has intersection with  $V_i$  no more than  $\ell_i$  for each  $i$ . I.e., consider

$$\mathcal{I}_p = \{X : \forall i \in [m], |V_i \cap X| \leq \ell_i\}. \quad (6.264)$$

<sup>14</sup> Then  $(V, \mathcal{I}_p)$  is a matroid. The corresponding submodular maximization problem is a natural  
<sup>15</sup> generalization of the cardinality constraint in that, rather than having a fixed number of elements  
<sup>16</sup> beyond which we are uninterested, the set of elements  $V$  is organized into groups, and here we have  
<sup>17</sup> a fixed per-group limit beyond which we are uninterested. This is useful for fairness applications  
<sup>18</sup> since the solution must be distributed over the blocks of the matroid. Still, there are many much  
<sup>19</sup> more powerful types of matroids that one can use [Oxl11; GM12].

<sup>20</sup> Regardless of the matroid, the problem  $\max_{X \in \mathcal{I}} f(X)$  can be solved, with a  $1/2$  approximation  
<sup>21</sup> factor, using the same greedy algorithm as above [NWF78; FNW78]. Indeed, the greedy algorithm  
<sup>22</sup> has an intimate relationship with submodularity, a fact that is well studied in some of the seminal  
<sup>23</sup> works on submodularity [Edm70; Lov83; Sch04]. It is also possible to define constraints consisting  
<sup>24</sup> of an *intersection of matroids*, meaning that the solution must be simultaneously independent in  
<sup>25</sup> multiple distinct matroids. Adding on to this, we might wish a set to be independent in multiple  
<sup>26</sup> matroids and also satisfy a knapsack constraint. Knapsack constraints are not matroid constraints,  
<sup>27</sup> since there can be multiple maximal cost solutions that are not the same size (as must be the case in  
<sup>28</sup> a matroid). It is also possible to define discrete constraints using level sets of another completely  
<sup>29</sup> different submodular function [IB13] — given two submodular functions  $f$  and  $g$ , this leads to  
<sup>30</sup> optimization problems of the form  $\max_{X \subseteq V: g(X) \leq \alpha} f(X)$  (the submodular cost submodular knapsack,  
<sup>31</sup> or SCSK, problem) and  $\min_{X \subseteq V: g(X) \geq \alpha} f(X)$  (the submodular cost submodular cover, or SCSC,  
<sup>32</sup> problem). Other examples include covering constraints [IN09], and cut constraints [JB16]. Indeed,  
<sup>33</sup> the type of constraints on submodular maximization for which good and scalable algorithms exist is  
<sup>34</sup> quite vast, and still growing.

<sup>35</sup> One last note on submodular maximization. In the above, the function  $f$  has been assumed to be  
<sup>36</sup> a polymatroid function. There are many submodular functions that are not monotone [Buc+12].  
<sup>37</sup> One example we saw before, namely the graph cut function. Another example is the log of the  
<sup>38</sup> determinant (log-determinant) of a submatrix of a positive-definite matrix (which is the Gaussian  
<sup>39</sup> entropy plus a constant). Suppose that  $\mathbf{M}$  is an  $n \times n$  symmetric positive-definite (SPD) matrix,  
<sup>40</sup> and that  $\mathbf{M}_X$  is a row-column submatrix (i.e., it is an  $|X| \times |X|$  matrix consisting of the rows and  
<sup>41</sup> columns of  $\mathbf{M}$  consisting of the elements in  $X$ ). Then the function defined as  $f(X) = \log \det(\mathbf{M}_X)$   
<sup>42</sup> is submodular but not necessarily monotone non-decreasing. In fact, the submodularity of the  
<sup>43</sup> log-determinant function is one of the reasons that *determinantal point processes* (DPPs), which  
<sup>44</sup>

1 instantiate probability distributions over sets in such a way that high probability is given to those  
 2 subsets that are diverse according to  $\mathbf{M}$ , are useful for certain tasks where we wish to probabilistically  
 3 model diversity [KT11]. Diversity of a set  $X$  here is measured by the volume of the parallelepiped  
 4 which is known to be computed as the determinant of the submatrix  $\mathbf{M}_X$  and taking the log of this  
 5 volume makes the function submodular in  $X$ . A DPP in fact is an example of a log-submodular  
 6 probabilistic model (more in Section 6.11.10).

#### 10 6.11.4.3 Submodular Function Minimization

11 In the case of a polymatroid function, unconstrained minimization is again trivial. However, even in  
 12 the unconstrained case, the minimization of an arbitrary (i.e., not necessarily monotone) submodular  
 13 function  $\min_{X \subseteq V} f(X)$  might seem hopelessly intractable. Unconstrained submodular maximization  
 14 is NP-hard (albeit approximable) and this is not surprising given that there are an exponential  
 15 number of sets needing to be considered. Remarkably, submodular minimization does not require  
 16 exponential computation, is not NP-hard, and in fact, there are polynomial time algorithms for  
 17 doing so, something that is not at all obvious. This is one of the important characteristics that  
 18 submodular functions share with convex functions, their common amenability to minimization.  
 19 Starting in the very late 1960s, and spearheaded by individuals such as Jack Edmonds [Edm70],  
 20 there was a concerted effort in the discrete mathematics community in search of either an algorithm  
 21 that could minimize a submodular function in polynomial time or a proof that such a problem was  
 22 NP-hard. The nut was finally cracked in a classic paper [GLS81] on the ellipsoid algorithm that gave  
 23 a polynomial time algorithm for submodular function minimization (SFM). While the algorithm was  
 24 polynomial, it was a continuous algorithm, and it was not practical, so the search continued for a  
 25 purely combinatorial strongly polynomial time algorithm. Queyranne [Que98] then proved that an  
 26 algorithm [NI92] worked for this problem when the set function also satisfies a symmetry condition  
 27 (i.e.,  $\forall X \subseteq V, f(X) = f(V \setminus X)$ ), which only requires  $O(n^3)$  time. The result finally came in around  
 28 year 2000 using two mostly independent methods [IFF00; Sch00]. These algorithms, however, also  
 29 were impractical in that while they are polynomial time, they have unrealistically high polynomial  
 30 degree (i.e.,  $\tilde{O}(V^7 * \gamma + V^8)$  for [Sch00] and  $\tilde{O}(V^7 * \gamma)$  for [IFF00]). This led to additional work on  
 31 combinatorial algorithms for SFM leading to algorithms that could perform SFM in time  $\tilde{O}(V^5\gamma + V^6)$   
 32 in [IO09]. Two practical algorithms for SFM include the Fujishige-Wolfe procedure [Fuj05; Wol76]<sup>7</sup>  
 33 as well as the Frank-Wolfe procedure, each of which minimize the 2-norm on a polyhedron  $B_f$   
 34 associated with the submodular function  $f$  and which is defined below (it should also be noted  
 35 that the Frank-Wolfe algorithm can also be used to minimize the convex extension of the function,  
 36 something that is relatively easy to compute via the Lovász extension [Lov83]). More recent work on  
 37 SFM are also based continuous relaxations of the problem in some form or another, leading algorithms  
 38 with strongly polynomial running time [LSW15] of  $O(n^3 \log^2 n)$  for which it was possible to drop the  
 39 log factors leading to a complexity of  $O(n^3)$  in [Jia21], weakly-polynomial running time [LSW15] of  
 40  $\tilde{O}(n^2 \log M)$  (where  $M \geq \max_{S \subseteq V} |f(S)|$ ), pseudo-polynomial running time [ALS20; Cha+17] of  
 41  $\tilde{O}(nM^2)$ , and a  $\epsilon$ -approximate minimization with a linear running time [ALS20] of  $\tilde{O}(n/\epsilon^2)$ . There  
 42 have been other efforts to utilize parallelism to further improve SFM [BS20].

43  
 44  
 45  
 46 7. This is the same Wolfe as the Wolfe in Frank-Wolfe but not the same algorithm.  
 47

1    **6.11.5 Applications of Submodularity in Machine Learning and AI**

2    Submodularity arises naturally in applications in machine learning and artificial intelligence, but its  
3    utility has still not yet been as widely recognized and exploited as other techniques. For example,  
4    while information theoretic concepts like entropy and mutual information are extremely widely used  
5    in machine learning (e.g., the cross entropy loss for classification is ubiquitous), the submodularity  
6    property of entropy is not nearly as widely explored.  
7

8    Still, the last several decades, submodularity has been increasingly studied and utilized in the  
9    context of machine learning. The below we begin to provide only a brief survey of some of the major  
10   sub-areas within machine learning that have been touched by submodularity. The list is not meant  
11   to be exhaustive, or even extensive. It is hoped that the below should, at least, offer a reasonable  
12   introduction into how submodularity has been and can continue to be useful in machine learning and  
13   artificial intelligence.  
14

15

16    **6.11.6 Sketching, CoreSets, Distillation, and Data Subset & Feature Selection**  
17

18    A summary is a concise representation of a body of data that can be used as an effective and efficient  
19    substitute for that data. There are many types of summary and some are extremely simple. For  
20    example, the mean or median of a list of numbers summarizes some property (the central tendency)  
21    of that list. A random subset is also a form of summary.

22    Any given summary, however, is not guaranteed to do a good job serving all purposes. Moreover,  
23    a summary usually involves at least some degree of approximation and fidelity loss relative to the  
24    original, and different summaries are faithful to the original in different ways and for different tasks.  
25    For these and other reasons, the field of summarization is rich and diverse, and summarization  
26    procedures are often very specialized.

27    Several distinct names for summarization have been used over the past few decades, including  
28    “sketches”, “coresets”, (in the field of natural language processing) “summaries”, and “distillation.”

29    Sketches [Cor17; CY20; Cor+12], arose in the field of computer science and was based on the  
30    acknowledgment that data is often too large to fit in memory and too large for an algorithm to run  
31    on a given machine, something enabled by a much smaller but still representative, and provably  
32    approximate, representation of the data.

33    Coresets are similar to sketches and there are some properties that are more often associated  
34    with coresets than with sketches, but sometimes the distinction is a bit vague. The notion of a  
35    coreset [BHP02; AHP+05; BC08] comes from the field of computational geometry where one is  
36    interested in solving certain geometric problems based on a set of points in  $\mathbb{R}^d$ . For any geometric  
37    problem and a set of points, a coreset problem typically involves finding the smallest weighted subset  
38    of points so that when an algorithm is run on the weighted subset, it produces approximately the  
39    same answer as when it is run on the original large dataset. For example, given a set of points, one  
40    might wish to find the diameter of set, or the radius of the smallest enclosing sphere, or finding the  
41    narrowest annulus (ring) containing the points, or a subset of points whose  $k$ -center clustering is  
42    approximately the same as the  $k$ -center clustering of the whole [BHP02].

43    Document summarization became one of the most important problems in natural language process-  
44    ing (NLP) in the 1990s although the idea of computing a summary of a text goes back much further to  
45    the 1950s [Luh58; Edm69], also and coincidentally around the same time that the CliffsNotes [Wik21]  
46    organization began. There are two main forms of document summarization [YWX17]. With *extractive*  
47

1 summarization [NM12], a set of sentences (or phrases) are extracted from the documents needing to  
2 be summarized, and the resulting subset of sentences, perhaps appropriately ordered, comprises the  
3 summary.  
4

5 With abstractive summarization [LN19], on the other hand, the goal is to produce an “abstract” of  
6 the documents, where one is not constrained to have any of the sentences in the abstract correspond  
7 to any of the sentences in the original documents. With abstractive summarization, therefore, the  
8 goal is to synthesize a small set of new pseudo sentences that represent the original documents.  
9 CliffsNotes, for example, are abstractive summaries of the literature being represented.

10 Another form of summarization that has more recently become popular in the machine learning  
11 community is *data distillation* [SG06b; Wan+20c; Suc+20; BYH20; NCL20; SS21; Ngu+21] or  
12 equivalently *dataset condensation* [ZMB21; ZB21]. With data distillation<sup>8</sup>, the goal is to produce a  
13 small set of synthetic pseudo-samples that can be used, for example, to train a model. The key here  
14 is that in the reduced dataset, the samples are not compelled to be the same as, or a subset of, the  
15 original dataset.

16 All of the above should be contrasted with data *compression*, which in some sense is the most  
17 extreme data reduction method. With compression, either lossless or lossy, one is no longer under  
18 any obligation that the reduced form of the data need to be used or recognizable by any algorithm or  
19 entity other than the decoder, or uncompression, algorithm.

### 20 6.11.6.1 Summarization Algorithm Design Choices

21 It is the author’s contention that the notions of summarization, coresets, sketching, and distillation  
22 are certainly analogous and quite possibly synonymous, and they are all different from compression.  
23 The different names for summarization are simply different nomenclatures for the same language  
24 game. What matters is not what you call it but the choices one makes when designing a procedure  
25 for summarization. And indeed, there are many choices.  
26

27 Submodularity offers essentially an infinite number ways to perform data sketching and coresets.  
28 When we view the submodular function as an information function (as we discuss in Section 6.11.7),  
29 where  $f(X)$  is the information contained in set  $X$  and  $f(V)$  is the maximum available information,  
30 finding the small  $X$  that maximizes  $f(X)$  (i.e.,  $X^* \in \text{argmax}\{f(X) : |X| \leq k\}$ ), is a form of coresset  
31 computation that is parameterized by the function  $f$  which has  $2^n$  parameters since  $f$  lives in a  
32  $2^n$ -dimensional cone. Performing this maximization will then minimize the residual information  
33  $f(V \setminus X|X)$  about anything not present the summary  $V \setminus X$  since  $f(V) = f(X \cup V \setminus X) =$   
34  $f(V \setminus X|X) + f(X)$  so maximizing  $f(X)$  will minimize  $f(V \setminus X|X)$ . For every  $f$ , moreover, the same  
35 algorithm (e.g., the greedy algorithm) can be used to produce the summarization, and in every case  
36 there is an approximation guarantee relative to the current  $f$ , as mentioned in earlier sections, as long  
37 as  $f$  stays submodular. Hence, submodularity provides a universal framework for summarization,  
38 coresets, and sketches to the extent that the space of submodular functions itself is sufficiently  
39 diverse and spans over different coreset problems.  
40

41 Overall, the corset or sketching problem, when using submodular functions, therefore becomes  
42 a problem of “submodular design.” That is, how do we construct submodular function that, for a  
43 particular problem, acts as a good coresset producer when the function is maximized. There are three  
44 general approaches to produce an  $f$  that works well as a summarization objective: (1) a pragmatic

45 8. Data distillation is distinct from the notion of *knowledge distillation* [HVD14; BC14; BCNM06] or *model distillation*,  
46 where the “knowledge” contained in a large model is distilled or reduced down into a different smaller model.

1 approach where the function is constructed by hand and heuristics, (2) a learning approach where all  
 2 or part of the submodular function is inferred from an optimization procedure, and (3) a mathematical  
 3 approach where a given submodular function when optimized offers of a coresset property.  
 4

5 When the primary goal is a practical and scalable algorithm that can produce an extractive  
 6 summary that works well on a variety of different data types, and if one is comfortable with heuristics  
 7 that work well in practice, a good option is to specify a submodular function by hand. For example,  
 8 given a similarity matrix, it is easy to instantiate a facility location function and maximize it to  
 9 produce a summary. If there are multiple similarity matrices, one can construct multiple facility  
 10 location functions and maximize their convex combination. Such an approach is viable and practical  
 11 and has been used successfully many times in the past for producing good summaries. One of the  
 12 earliest examples of this the algorithm presented in [KKT03] that shows how a submodular model can  
 13 be used to select the most influential nodes in a social network. Perhaps the earliest example of this  
 14 approach used for data subset selection for machine learning is [LB09] which utilizes a submodular  
 15 facility location function based on Fisher kernels (gradients w.r.t. parameters of log probabilities)  
 16 and applies it to unsupervised speech selection to reduce transcription costs. Other examples of  
 17 this approach includes: [LB10a; LB11] which developed submodular functions for query-focused  
 18 document summarization; [KB14b] which computes a subset of training data in the context of  
 19 transductive learning in a statistical machine translation system; [LB10b; Wei+13; Wei+14] which  
 20 develops submodular functions for speech data subset selection (the former, incidentally, is the first  
 21 use of a deep submodular function and the latter does this in an unsupervised label-free fashion);  
 22 [SS18a] which is a form of robust submodularity for producing coresets for training CNNs; [Kau+19]  
 23 which uses a facility location to facilitate diversity selection in active learning; [Bai+15; CTN17]  
 24 which develops a mixture of submodular functions for document summarization where the mixture  
 25 coefficients are also included in the hyperparameter set; [Xu+15] uses a symmetrized submodular  
 26 function for the purposes of video summarization.

27 The learnability and identifiability of submodular functions has received a good amount of study  
 28 from a theoretical perspective. Starting with the strictest learning settings, the problem looks pretty  
 29 dire. For example, [SF08; Goe+09] shows that if one is restricted to making a polynomial number of  
 30 queries (i.e., training pairs of the form  $(S, f(S))$ ) of a monotone submodular function, then it is not  
 31 possible to approximate  $f$  with a multiplicative approximation factor better than  $\tilde{\Omega}(\sqrt{n})$ . In [BH11],  
 32 goodness is judged multiplicatively, meaning for a set  $A \subseteq V$  we wish that  $\tilde{f}(A) \leq f(A) \leq g(n)f(A)$   
 33 for some function  $g(n)$ , and this is typically a probabilistic condition (i.e., measured by distribution,  
 34 or  $\tilde{f}(A) \leq f(A) \leq g(n)f(A)$ , should happen on a fraction at least  $1 - \beta$  of the points). Alternatively,  
 35 goodness may also be measured by an additive approximation error, say by a norm. I.e., defining  
 36  $\text{err}_p(f, \tilde{f}) = \|f - \tilde{f}\|_p = (E_{A \sim \Pr}[|f(A) - \tilde{f}(A)|^p])^{1/p}$ , we may wish  $\text{err}_p(f, \tilde{f}) < \epsilon$  for  $p = 1$  or  $p = 2$ .  
 37 In the PAC (probably approximately correct) model, we probably ( $\delta > 0$ ) approximately ( $\epsilon > 0$  or  
 38  $g(n) > 1$ ) learn ( $\beta = 0$ ) with a sample or algorithmic complexity that depends on  $\delta$  and  $g(n)$ . In the  
 39 PMAC (probably mostly approximately correct) model [BH11], we also “mostly”  $\beta > 0$  learn. In some  
 40 cases we wish to learn the best submodular approximation to a non-submodular function. In other  
 41 cases, we are allowed to deviate from submodularity as long as the error is small. Learning special  
 42 cases includes coverage functions [FK14; FK13a], and low-degree polynomials [FV15], curvature  
 43 limited functions [IJB13], functions with a limited “goal” [DHK14; Bac+18], functions that are  
 44 Fourier sparse [Wen+20a], or that are of a family called “juntas” [FV16], or that come from families  
 45 other than submodular [DFF21], and still others [BRS17; FKV14; FKV17; FKV20; FKV13; YZ19].  
 46 Other results include that one can not minimize a submodular function by learning it first from  
 47

1 samples [BS17]. The essential strategy of learning is to attempt to construct a submodular function  
2 approximation  $\hat{f}$  from an underlying submodular function  $f$  querying the latter only a small number  
3 of times. The overall gist of these results is that it is hard to learn everywhere and accurately.  
4

5 In the machine learning community, learning can be performed extremely efficiently in practice,  
6 although there are not the types of guarantees as one finds above. For example, given a mixture  
7 of submodular components of the form  $f(A) = \sum_i \alpha_i f_i(A)$ , if each  $f_i$  is considered fixed, then the  
8 learning occurs only over the mixture coefficients  $\alpha_i$ . This can be solved as a linear regression problem  
9 where the optimal coefficients can be computed in a linear regression setting. Alternatively, such  
10 functions can be learnt in a max-margin setting where the goal is primarily to adjust  $\alpha_i$  to ensure  
11 that  $f(A)$  is large on certain subsets [SSJ12; LB12; Tsc+14]. Even here there are practical challenges,  
12 however, since it is in general hard in practice to obtain a training set of pairs  $\{(S_i, F(S_i))\}_i$ .  
13 Alternatively, one also “learn” a submodular function in a reinforcement learning setting [CKK17] by  
14 optimizing the implicit function directly from gain vectors queried from an environment. In general,  
15 such practical learning algorithms have been used for image summarization [Tsc+14], document  
16 summarization [LB12], and video summarization [GGG15; Vas+17a; Gon+14; SGS16; SLG17]. While  
17 none of these learning approaches claim to approximate some true underlying submodular function,  
18 in practice, they do perform better than the by-hand crafting of a submodular functions mentioned  
19 above.

20 By a submodularity based coresset, we mean one where the direct optimization of a submodular  
21 function offers a theoretical guarantee for some specific problem. This is distinct from above where  
22 the submodular function is used as a surrogate heuristic objective function and for which, even if the  
23 submodular function is learnt, optimizing it is only a heuristic for the original problem. In some  
24 limited cases, it can be shown that the function we wish to approximate is already submodular,  
25 e.g., in the case of certain naive Bayes and k-NN classifiers [WIB15] where the training accuracy,  
26 as a function of the training data subset, can be shown to be submodular. Hence, maximizing this  
27 function offers the same guarantee on the training accuracy as it does on the submodular function.  
28 Unfortunately, the accuracy function for many models are not submodular, although they do have a  
29 difference of submodular [NB05; IB12] decomposition.

30 In other cases, it can be shown that certain desirable coresset objectives are inherently submodular.  
31 For example, in [MBL20], it is shown that the normed difference between the overall gradient  
32 (from summing over all samples in the training data) and an approximate gradient (from summing  
33 over only samples in a summary) can be upper bounded with a supermodular function that, when  
34 converted to a submodular facility location function and maximized, will select a set that reduces  
35 this difference, and will lead to similar convergence rates to an approximate optimum solution in the  
36 convex case. A similar example of this in a DPP context is shown in [TBA19]. In other cases, subsets  
37 of the training data and training occur simultaneously using a continuous-discrete optimization  
38 framework, where the goal is to minimize the loss on diverse and challenging samples measured by a  
39 submodular objective [ZB18]. In still other cases, bi-level objectives related to but not guaranteed to  
40 be submodular can be formed where a set is selected from a training set with the deliberate purpose  
41 of doing well on a validation set [Kil+20; BMK20].

42 The methods above have focused on reducing the number of samples in a training data. Considering  
43 the transpose of a design matrix, however, all of the above methods can be used for reducing the  
44 features of a machine learning procedure as well. Specifically, any of the extractive summarization,  
45 subset selection, or coresset methods can be seen as feature selection while any of the abstract  
46 summarization, sketching, or distillation approaches can be seen as dimensionality reduction.

---

1    **6.11.7 Combinatorial Information Functions**

2

3

4 The entropy function over a set of random variables  $X_1, X_2, \dots, X_n$  is defined as  $H(X_1, X_2, \dots, X_n) =$   
5  $-\sum_{x_1, x_2, \dots, x_n} p(x_1, \dots, x_n) \log p(x_1, \dots, x_n)$ . From this we can define three set-argument conditional  
6 mutual information functions as  $I_H(A; B|C) = I(X_A; X_B|X_C)$  where the latter is the mutual  
7 information between variables indexed by  $A$  and  $B$  given variables indexed by  $C$ . This mutual  
8 information expresses the residual information between  $X_A$  and  $X_B$  that is not explained by their  
9 common information with  $X_C$ .

10 As mentioned above, we may view any polymatroid function as a type of information function over  
11 subsets of  $V$ . That is,  $f(A)$  is the information in set  $A$  — to the extent that this is true, this property  
12 justifies  $f$ 's use as a summarization objective as mentioned above. The reason  $f$  may be viewed as an  
13 information function stems from  $f$  being normalized,  $f$ 's non-negativity,  $f$ 's monotonicity, and the  
14 property that further conditioning reduces valuation (i.e.,  $f(A|B) \geq f(A|B, C)$  which is identical to  
15 the submodularity property). These properties were outlined as being essential to the entropy function  
16 in Shannon's original work [Sha48] but are true of any polymatroid function as well. Hence, given any  
17 polymatroid function  $f$ , it is possible to define a combinatorial mutual information function [Iye+21]  
18 in a similar way. Specifically, we can define the combinatorial (submodular) conditional mutual  
19 information (CCMI) as  $I_f(A; B|C) = f(A + C) + f(B + C) - f(C) - f(A + B + C)$ , which has been  
20 known as the connectivity function [Cun83] amongst other names. If  $f$  is the entropy function then  
21 this yields the standard entropic mutual information but here the mutual information can be defined  
22 for any submodular information measure  $f$ . For an arbitrary polymatroid  $f$ , therefore,  $I_f(A; B|C)$   
23 can be seen as an  $A, B$  set-pair similarity score that ignores, neglects, or discounts any common  
24 similarity between the  $A, B$  pair that is due to  $C$ .

25 Historical use of a special case of CCMI, i.e.,  $I_f(A; B)$  where  $C = \emptyset$ , occurred in a number of  
26 circumstances. For example, in [GKS05] the function  $g(A) = I_f(A; V \setminus A)$  (which, incidentally  
27 is both symmetric ( $g(A) = g(V \setminus A)$  for all  $A$ ) and submodular was optimized using the greedy  
28 procedure which has a guarantee as long as  $g(A)$  is monotone up  $2k$  elements whenever one wishes  
29 for a summary of size  $k$ . This was done for  $f$  being the entropy function, but it can be used for  
30 any polymatroid function. In similar work where  $f$  is Shannon entropy, [KG05] demonstrated that  
31  $g_C(A) = I_f(A; C)$  for a fixed set  $C$  is not submodular in  $A$  but if it is the case that the elements of  
32  $V$  are independent given  $C$  then submodularity is preserved. This can be seen quickly easily by the  
33 consequence of the assumption which states that  $I_f(A; C) = f(A) - f(A|C) = f(A) - \sum_{a \in A} f(a|C)$   
34 where the second equality is due to the conditional independence property. In this case,  $I_f$  is the  
35 difference between a submodular and a modular function which preserves submodularity for any  
36 polymatroid  $f$ .

37 On the other hand, it would be useful for  $g_{B,C}(A) = I_f(A; B|C)$ , where  $B$  and  $C$  are fixed, to be  
38 possible to optimize in terms of  $A$ . One can view this function as one that, when it is maximized,  
39 chooses  $A$  to be similar to  $B$  in a way that neglects or discounts any common similarity that  $A$  and  
40  $B$  have with  $C$ . One option to optimize this function to utilize difference of submodular [NB05; IB12]  
41 optimization as mentioned earlier. A more recent result shows that in some cases  $g_{B,C}(A)$  is still  
42 submodular in  $A$ . Define the second order partial derivative of a submodular function  $f$  as follows  
43  $f(i, j|S) \triangleq f(j|S + i) - f(j|S)$ . Then if it is the case that  $f(i, j|S)$  is monotone non-decreasing in  $S$   
44 for  $S \subseteq V \setminus \{i, j\}$  then  $I_f(A; B|C)$  is submodular in  $A$  for fixed  $B$  and  $C$ . It may be thought that  
45 only esoteric functions have this property but in fact [Iye+21] shows that this is true for a number of  
46 widely used submodular functions in practice, including the facility location function which results  
47

in the form  $I_f(A; B|C) = \sum_{v \in V} \max\left(\min\left(\sum_{a \in A} \text{sim}(v, a), \max_{b \in B} \text{sim}(v, b)\right) - \max_{c \in C} \text{sim}(v, c), 0\right)$ . This function was used [Kot+22] to produce summaries  $A$  that were particularly relevant to a query given by  $B$  but that should neglect information in  $C$  that can be considered “private” information to avoid.

### 6.11.8 Clustering, Data Partitioning, and Parallel Machine Learning

There are an almost limited number of clustering algorithms and a plethora of reviews on their variants. Any given submodular function can also instantiate a clustering procedure as well, and there are several ways to do this. Here we offer only a brief outline of the approach. In the last section, we defined  $I_f(A; V \setminus A)$  as the CCMI between  $A$  and everything but  $A$ . When we view this as a function of  $A$ , then  $g(A) = I_f(A; V \setminus A)$  and  $g(A)$  is a symmetric submodular function that can be minimized using Queyranne’s algorithm [Que98; NI92]. Once this is done, the resulting  $A$  is such that it is least similar to  $V \setminus A$ , according to  $I_f(A; V \setminus A)$  and hence forms a 2-clustering. This process can then be recursively applied where we form two new functions  $g_A(B) = I_f(B; A \setminus B)$  for  $B \subseteq A$  and  $g_{V \setminus A}(B) = I_f(B; (V \setminus A) \setminus B)$  for  $B \subseteq V \setminus A$ . These are two symmetric submodular functions on different ground sets that also can be minimized using Queyranne’s algorithm. This recursive bisection algorithm then repeats until the desired number of clusters is formed. Hence, the CCMI function can be used as a top-down recursive bisection clustering procedure and has been called Q-clustering [NJB05; NB06]. It should be noted that such forms of clustering often generalizes forming a multi-way cut in an undirected graph in which case the objective becomes the graph-cut function that, as we saw above, is also submodular. In some cases, the number of clusters need not be specified in advance [NKI10]. Another submodular approach to clustering can be found in [Wei+15b] where the goal is to minimize the maximum valued block in a partitioning which can lead to submodular load balancing or minimum makespan scheduling [HS88; LST90].

Yet another form of clustering can be seen via the simple cardinality constrained submodular maximization process itself which can be compared to a  $k$ -medoids process whenever the objective  $f$  is the facility location function. Hence, any such submodular function can be seen as a submodular-function-parameterized form of finding the  $k$  “centers” among a set of data items. There have been numerous applications of submodular clustering. For example, using these techniques it is possible to identify parcellations of the human brain [Sal+17a]. Other applications include partitioning data for more effective and accurate and lower variance distributed machine learning training [Wei+15a] and also for more ideal mini-batch construction for training deep neural networks [Wan+19b].

### 6.11.9 Active and Semi-Supervised Learning

Suppose we are given data set  $\{x_i, y_i\}_{i \in V}$  consisting of  $|V| = n$  samples of  $x, y$  pairs but where the labels are unknown. Samples are labeled one at a time or one mini-batch at a time, and after each labeling step  $t$  each remaining unlabeled sample is given a score  $s_t(x_i)$  that indicates the potential benefit of acquiring a label for that sample. Examples include the entropy of the model’s output distribution on  $x_i$ , or a margin based score consisting of the difference between the top and the second-from-the-top posterior probability. This produces a modular function on the unlabeled samples,  $m_t(A) = \sum_{a \in A} s(x_a)$  where  $A \subseteq V$ . It is simple to use this modular function to produce a mini-batch active learning procedure where at each stage we form  $A_t \in \operatorname{argmax}_{A \subseteq U_t: |A|=k} m_t(A)$  where  $U_t$  is the set of labeled samples at stage  $t$ . Then  $A_t$  is a set of size  $k$  that gets labeled, we form

1  $U_t = U_t \setminus A_t$ , update  $s_t(a)$  for  $a \in U_t$  and repeat. This is called **active learning** (Section 34.7).

2 The reason for using active learning with mini-batches of size greater than one is that it is often  
3 inefficient to ask for single label at a time. The problem with such a minibatch strategy, however,  
4 is that the set  $A_t$  can be redundant. The reason is that the uncertainty about every sample in  $A_t$   
5 could be owing to the same underlying cause — even though the model is most uncertain about  
6 samples in  $A_t$ , once one sample in  $A_t$  is labeled, it may not be optimal to label the remaining samples  
7 in  $A_t$  due to this redundancy. Utilizing submodularity, therefore, can help reduce this redundancy.  
8 Suppose  $f_t(A)$  is a submodular diversity model over samples at step  $t$ . At each stage, choosing  
9 the set of samples to label becomes  $A_t \in \operatorname{argmax}_{A \subseteq U_t: |A|=k} m_t(A) + f_t(A)$  —  $A_t$  is selected based  
10 on a combination of both uncertainty (via  $m_t(A)$ ) and diversity (via  $f_t(A)$ ). This is precisely the  
11 submodular active learning approach taken in [WIB15; Kau+19].

12 Another quite different approach to a form of submodular “batch” active learning setting where a  
13 batch  $L$  of labeled samples are selected all at once and then used to label the rest of the unlabeled  
14 samples. This also allows the remaining unlabeled samples to be utilized in a semi-supervised  
15 framework [GB09; GB11]. In this setting, we start with a graph  $G = (V, E)$  where the nodes  $V$   
16 need to be given a binary  $\{0, 1\}$ -valued label,  $y \in \{0, 1\}^V$ . For any  $A \subseteq V$  let  $y_A \in \{0, 1\}^A$  be  
17 the labels just for node set  $A$ . We also define  $V(y) \subseteq V$  as  $V(y) = \{v \in V : y_v = 1\}$ . Hence  
18  $V(y)$  are the graph nodes labeled 1 by  $y$  and  $V \setminus V(y)$  are the nodes labeled 0. Given submodular  
19 objective  $f$ , we form its symmetric CCMI variant  $I_f(A) \triangleq I_f(A; V \setminus A)$  — note that  $I_f(A)$  is always  
20 submodular in  $A$ . This allows  $I_f(V(y))$  to determine the “smoothness” of a given candidate labeling  
21  $y$ . For example, if  $I_f$  is the weighted graph cut function where each weight corresponds to an affinity  
22 between the corresponding two nodes, then  $I_f(V(y))$  would be small if  $V(y)$  (the 1-labeled nodes)  
23 do not have strong affinity with  $V \setminus V(y)$  (the 0-labeled nodes). In general, however,  $I_f$  can be any  
24 symmetric submodular function. Let  $L \subseteq V$  be any candidate set of nodes to be labeled, and define  
25  $\Psi(L) \triangleq \min_{T \subseteq (V \setminus L): T \neq \emptyset} I_f(T)/|T|$ . Then  $\Psi(L)$  measures the “strength” of  $L$  in that if  $\Psi(L)$  is small,  
26 an adversary can label nodes other than  $L$  without being too unsatisfactory according to  $I_f$ , while if  
27  $\Psi(L)$  is large, an adversary can do no such thing. Then [GB11] showed that given a node set  $L$  to be  
28 queried, and the corresponding correct labels  $y_L$  that are completed (in a semi-supervised fashion)  
29 according to the following  $y' = \operatorname{argmin}_{\hat{y} \in \{0, 1\}^V: \hat{y}_L = y_L} I_f(V(\hat{y}))$ , then this results in the following  
30 bound on the true labeling  $\|y - y'\|^2 \leq 2I_f(V(y))/\Psi(L)$  suggesting that we can find a good set  
31 to query by maximizing  $L$  in  $\Psi(L)$ , and this holds for any submodular function. Of course it is  
32 necessary to find an underlying submodular function  $f$  that fits a given problem, and this is discussed  
33 in Section 6.11.6.

34

### 35 36 6.11.10 Probabilistic Modeling

37 Graphical models are often used to describe factorization requirements on families of probability  
38 distributions. Factorization is not the only way, however, to describe restrictions on such families.  
39 In a graphical model, graphs describe only which random variable may directly interact with other  
40 random variable. An entirely different strategy for producing families of often-tractable probabilistic  
41 models can be produced without requiring any factorization property at all. Considering an energy  
42 function  $E(x)$  where  $p(x) \propto \exp(E(x))$ , factorizations correspond to there being cliques in the graph  
43 such that the graph’s tree-width often is limited. On the other hand, finding  $\max_x p(x)$  is the same  
44 as finding  $\min_x E(x)$ , something that can be done if  $E(x) = f(V(x))$  is a submodular function (using  
45 the earlier used notation  $V(x)$  to map from binary vectors to subsets of  $V$ ). Even a submodular  
46

47

function as simple as  $f(A) = \sqrt{|A|} - m(A)$  where  $m$  is modular has tree-width of  $n - 1$ , and this leads to an energy function  $E(x)$  that allows  $\max_x p(x)$  to be solved in polynomial time using submodular function minimization (see Section 6.11.4.3). Such restrictions to  $E(x)$  therefore are not of the form *amongst the random variables, who is allowed to directly interact with whom*, but rather *amongst the random variables, what is the manner that they interact*. Such potential function restrictions can also combine with direct interaction restrictions as well and this has been widely used in computer vision, leading to cases where graph-cut and graph-cut like “move making” algorithms (such as *alpha-beta* swap and *alpha*-expansion algorithms) used in attractive models [BVZ99; BK01; BVZ01; SWW08]. In fact, the culmination of these efforts [KZ02] lead to a rediscovery of the submodularity (or the “regular” property) as being the essential ingredient for when Markov random fields can be solved using graph cut minimization, which is a special case of submodular function minimization.

The above model can be seen as log-supermodular since  $\log p(x) = -E(x) + \log 1/Z$  is a supermodular function. These are all distributions that put high probability on configurations that yield small valuation by a submodular function. Therefore, these distributions have high probability when  $x$  consists of a homogeneous and for this reason they are useful for computer vision segmentation problems (e.g., in a segment of an image, the nearby pixels should roughly be homogeneous as that is often what defines an object). The DPPs we saw above, however, are an example of a log-submodular probability distribution since  $f(X) = \log \det(\mathbf{M}_X)$  is submodular. These models have high probability for diverse sets.

More generally,  $E(x)$  being either a submodular or supermodular function can produce log-submodular or log-supermodular distributions, covering both cases above where the partition function takes the form  $Z = \sum_{A \subseteq V} \exp(f(A))$  for objective  $f$ . Moreover, we often wish to perform tasks much more than just finding the most probable random variable assignments. This includes marginalization, computing the partition function, constrained maximization, and so on. Unfortunately, many of these more general probabilistic inference problems do not have polynomial time solutions even though the objectives are submodular or supermodular. On the other hand, such structure has opened the doors to an assortment of new probabilistic inference procedures that exploit this structure [DK14; DK15a; DTK16; ZDK15; DJK18]. Most of these methods were of the variational sort and offered bounds on the partition function  $Z$ , sometimes making use of the fact that submodular functions have easily computable semi-gradients [IB15; Fuj05] which are modular upper and lower bounds on a submodular or supermodular function that are tight at one or more subsets. Given a submodular (or supermodular) function  $f$  and a set  $A$ , it is possible to easily construct (in linear time) a modular function upperbound  $m^A : 2^V \rightarrow \mathbb{R}$  and a modular function lower bound  $m_A : 2^V \rightarrow \mathbb{R}$  having the properties that  $m_A(X) \leq f(X) \leq m^A(X)$  for all  $X \subseteq V$  and that is tight at  $X = A$  meaning  $m_A(A) = f(A) = m^A(A)$  [IB15]. For any modular function  $m$ , the probability function for a characteristic vector  $x = \mathbf{1}_A$  becomes  $p(\mathbf{1}_A) = 1/Z \exp(E(\mathbf{1}_A)) = \prod_{a \in A} \sigma(m(a)) \prod_{a \notin A} \sigma(-m(a))$  where  $\sigma$  is the logistic function. Thus, a modular approximation of a submodular function is like a mean-field approximation of the distribution, and makes the assumption that all random variables are independent. Such an approximation can then be used to compute quantities such as upper and lower bounds on the partition function, and much else.

### 6.11.11 Structured Norms and Loss Functions

Convex norms are used ubiquitously in machine learning, often as complexity penalizing regularizers (e.g., the ubiquitous  $p$ -norms for  $p \geq 1$ ) and also sometimes as losses (e.g., squared error). Identifying

1 new useful structured and possibly learnable sparse norms is an interesting and useful endeavor,  
2 and submodularity can help here as well. Firstly, recall the  $\ell_0$  or counting norm  $\|x\|_0$  simply counts  
3 the number of nonzero entries in  $x$ . When we wish for a sparse solution, we may wish to regularize  
4 using  $\|x\|_0$  but it both leads to an intractable combinatorial optimization problem and it leads to an  
5 object that is not differentiable. The usual approach is to find the closest convex relaxation of this  
6 norm and that is the one norm or  $\|x\|_1$ . This is convex in  $x$  and has a sub-gradient structure and  
7 hence can be combined with a loss function to produce an optimizable machine learning objective,  
8 for example the lasso. On the other hand  $\|x\|_1$  has no structure, as each element of  $x$  is penalized  
9 based on its absolute value irrespective of the state of any of the other elements. There have thus  
10 been efforts to develop group norms that penalize groups or subsets of elements of  $x$  together, such  
11 as group lasso [HTW15].  
12

13 It turns out that there is a way to utilize a submodular function as the regularizer. Penalizing  
14  $x$  via  $\|x\|_0$  is identical to penalizing it via  $|V(x)|$  and note that  $m(A) = |A|$  is a modular function.  
15 Instead, we could penalize  $x$  via  $f(V(x))$  for a submodular function  $f$ . Here, any element of  $x$   
16 being non-zero would allow for a diminishing penalty of other elements of  $x$  being zero all according  
17 to the submodular function, and such cooperative penalties can be obtained via a submodular  
18 parameterization. Like when using the zero-norm  $\|x\|_0$ , this leads to the same combinatorial problem  
19 due to continuous optimization of  $x$  with a penalty term of the form  $f(V(x))$ . To address this, we can  
20 use the Lovász extension  $\check{f}(x)$  on a vector  $x$ . This function is convex but it is not a norm, but if we  
21 consider the construct defined as  $\|x\|_f = \check{f}(|x|)$ , it can be shown that this satisfies all the properties  
22 of a norm for all non-trivial submodular functions [PG98; Bac+13] (i.e., those normalized submodular  
23 functions for which  $f(v) > 0$  for all  $v$ ). In fact, the group lasso mentioned above is a special case  
24 for a particularly simple feature-based submodular function (a sum of min-truncated cardinality  
25 functions). But in principle, the same submodular design strategies mentioned in Section 6.11.6 can  
26 be used to produce a submodular function to instantiate an appropriate convex structured norm for  
27 a given machine learning problem.  
28

### 29 6.11.12 Conclusions

30 We have only barely touched the surface of submodularity and how it applies to and can benefit  
31 machine learning. For more details, see [Bil22] and the many references contained therein. Considering  
32 once again the innocuous looking submodular inequality, then very much like the definition of  
33 convexity, we observe something that belies much of its complexity while opening the gates to wide  
34 and worthwhile avenues for machine learning exploration.  
35

36

37

38

39

40

41

42

43

44

45

46

47

PART II

## Inference



# 7

## Inference algorithms: an overview

### 7.1 Introduction

In the probabilistic approach to machine learning, all unknown quantities — be they predictions about the future, hidden states of a system, or parameters of a model — are treated as random variables, and endowed with probability distributions. The process of **inference** corresponds to computing the posterior distribution over these quantities, conditioning on whatever data is available.

Let  $\mathbf{h}$  represent the unknown variables, and  $\mathcal{D}$  represent the known variables. Given a likelihood  $p(\mathcal{D}|\mathbf{h})$  and a prior  $p(\mathbf{h})$ , we can compute the posterior  $p(\mathbf{h}|\mathcal{D})$  using Bayes' rule:

$$p(\mathbf{h}|\mathcal{D}) = \frac{p(\mathbf{h})p(\mathcal{D}|\mathbf{h})}{p(\mathcal{D})} \quad (7.1)$$

The main computational bottleneck is computing the normalization constant in the denominator, which requires solving the following high dimensional integral:

$$p(\mathcal{D}) = \int p(\mathcal{D}|\mathbf{h})p(\mathbf{h})d\mathbf{h} \quad (7.2)$$

This is needed to convert the unnormalized joint probability of some parameter value,  $p(\mathbf{h}, \mathcal{D})$ , to a normalized probability,  $p(\mathbf{h}|\mathcal{D})$ , which takes into account all the other plausible values that  $\mathbf{h}$  could have. Similarly, computing posterior marginals also requires computing integrals:

$$p(h_i|\mathcal{D}) = \int p(h_i, \mathbf{h}_{-i}|\mathcal{D})d\mathbf{h}_{-i} \quad (7.3)$$

Thus integration is at the heart of Bayesian inference, whereas differentiation is at the heart of optimization.

In this chapter, we give a high level summary of algorithmic techniques for computing (approximate) posteriors. We will give more details in the following chapters. Note that most of these methods are independent of the specific model. This allows problem solvers to focus on creating the best model possible for the task, and then relying on some inference engine to do the rest of the work — this latter process is sometimes called “**turning the Bayesian crank**”. For more details, see e.g., [Gel+14a; MKL21; MFR20].

### 7.2 Common inference patterns

There are kinds of posterior we may want to compute, but we can identify 3 main patterns, as we discuss below. These give rise to different types of inference algorithm, as we will see in later chapters.

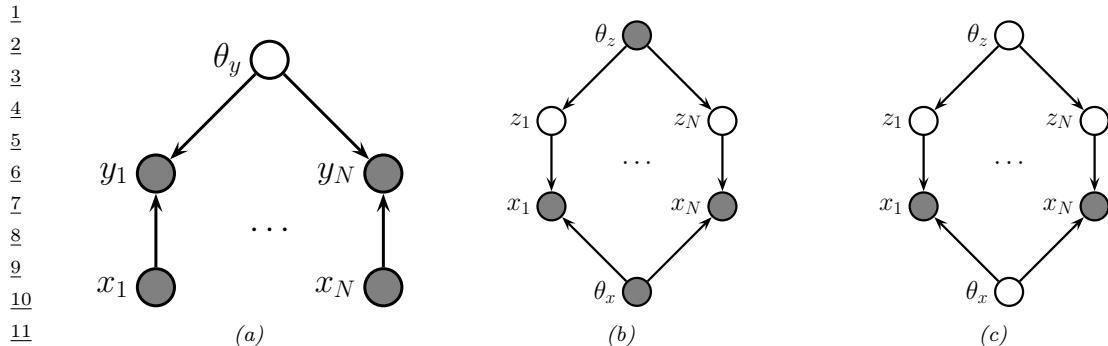


Figure 7.1: Graphical models with (a) Global hidden variables for representing the Bayesian discriminative model  $p(\mathbf{y}_{1:N}, \boldsymbol{\theta}_y | \mathbf{x}_{1:N}) = p(\boldsymbol{\theta}_y) \prod_{n=1}^N p(\mathbf{y}_n | \mathbf{x}_n; \boldsymbol{\theta}_y)$ ; (b) Local hidden variables for representing the generative model  $p(\mathbf{x}_{1:N}, \mathbf{z}_{1:N} | \boldsymbol{\theta}) = \prod_{n=1}^N p(\mathbf{z}_n | \boldsymbol{\theta}_z) p(\mathbf{x}_n | \mathbf{z}_n, \boldsymbol{\theta}_x)$ ; (c) Local and global hidden variables for representing the Bayesian generative model  $p(\mathbf{x}_{1:N}, \mathbf{z}_{1:N}, \boldsymbol{\theta}) = p(\boldsymbol{\theta}_z) p(\boldsymbol{\theta}_x) \prod_{n=1}^N p(\mathbf{z}_n | \boldsymbol{\theta}_z) p(\mathbf{x}_n | \mathbf{z}_n, \boldsymbol{\theta}_x)$ . Shaded nodes are assumed to be known (observed), unshaded nodes are hidden.

### 7.2.1 Global latents

The first pattern arises when we need to perform inference in models which have **global latent variables**, such as parameters of a model  $\boldsymbol{\theta}$ , which are shared across all  $N$  observed training cases. This is shown in Figure 7.1a, and corresponds to the usual setting for supervised or discriminative learning, where the joint distribution has the form

$$p(\mathbf{y}_{1:N}, \boldsymbol{\theta} | \mathbf{x}_{1:N}) = p(\boldsymbol{\theta}) \left[ \prod_{n=1}^N p(\mathbf{y}_n | \mathbf{x}_n, \boldsymbol{\theta}) \right] \quad (7.4)$$

The goal is to compute the posterior  $p(\boldsymbol{\theta} | \mathbf{x}_{1:N}, \mathbf{y}_{1:N})$ . Most of the Bayesian supervised learning models discussed in Part III follow this pattern.

### 7.2.2 Local latents

The second pattern arises when we need to perform inference in models which have **local latent variables**, such as hidden states  $\mathbf{z}_{1:N}$ ; we assume the model parameters  $\boldsymbol{\theta}$  are known. This is shown in Figure 7.1b. Now the joint distribution has the form

$$p(\mathbf{x}_{1:N}, \mathbf{z}_{1:N} | \boldsymbol{\theta}) = \left[ \prod_{n=1}^N p(\mathbf{x}_n | \mathbf{z}_n, \boldsymbol{\theta}_x) p(\mathbf{z}_n | \boldsymbol{\theta}_z) \right] \quad (7.5)$$

The goal is to compute  $p(\mathbf{z}_n | \mathbf{x}_n, \boldsymbol{\theta})$  for each  $n$ . This is the setting we consider for most of the PGM inference methods in Chapter 9.

If the parameters are not known (which is the case for most latent variable models, such as mixture models), we may choose to estimate them by some method (e.g., maximum likelihood), and then plugin this point estimate. The advantage of this approach is that, conditional on  $\boldsymbol{\theta}$ , all the latent variables are conditionally independent, so we can perform inference in parallel across the data. This

lets us use methods such as expectation maximization (Section 6.6.3), in which we infer  $p(\mathbf{z}_n | \mathbf{x}_n, \boldsymbol{\theta}_t)$  in the E step for all  $n$  simultaneously, and then update  $\boldsymbol{\theta}_t$  in the M step. If the inference of  $\mathbf{z}_n$  cannot be done exactly, we can use variational inference, a combination known as variational EM (Section 6.6.6.1).

Alternatively, we can use a minibatch approximation to the likelihood, marginalizing out  $\mathbf{z}_n$  for each example in the minibatch to get

$$\log p(\mathcal{D}_t | \boldsymbol{\theta}_t) = \sum_{n \in \mathcal{D}_t} \log \left[ \sum_{\mathbf{z}_n} p(\mathbf{x}_n, \mathbf{z}_n | \boldsymbol{\theta}_t) \right] \quad (7.6)$$

where  $\mathcal{D}_t$  is the minibatch at step  $t$ . If the marginalization cannot be done exactly, we can use variational inference, a combination known as stochastic variational inference or SVI (Section 10.3.1). We can also learn an inference network  $q_\phi(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta})$  to perform the inference for us, rather than running an inference engine for each example  $n$  in each batch  $t$ ; the cost of learning  $\phi$  can be amortized across the batches. This is called amortized SVI, and is commonly used to train deep latent variable models such as VAEs (Section 21.2).

### 7.2.3 Global and local latents

The third pattern arises when we need to perform inference in models which have **local and global latent variables**. This is shown in Figure 7.1c, and corresponds to the following joint distribution:

$$p(\mathbf{x}_{1:N}, \mathbf{z}_{1:N}, \boldsymbol{\theta}) = p(\boldsymbol{\theta}_x)p(\boldsymbol{\theta}_z) \left[ \prod_{n=1}^N p(\mathbf{x}_n | \mathbf{z}_n, \boldsymbol{\theta}_x)p(\mathbf{z}_n | \boldsymbol{\theta}_z) \right] \quad (7.7)$$

This is essentially a Bayesian version of the latent variable model in Figure 7.1b, where now we model uncertainty in both the local variables  $\mathbf{z}_n$  and the shared global variables  $\boldsymbol{\theta}$ . This approach is less common in the ML community, since it is often assumed that the uncertainty in the parameters  $\boldsymbol{\theta}$  is negligible compared to the uncertainty in the local variables  $\mathbf{z}_n$ . The reason for this is that the parameters are “informed” by all  $N$  data cases, whereas each local latent  $\mathbf{z}_n$  is only informed by a single data point, namely  $\mathbf{x}_n$ . Nevertheless, there are advantages to being “fully Bayesian”, and modeling uncertainty in both local and global variables. We will see some examples of this later in the book.

## 7.3 Exact inference algorithms

In some cases, we can perform exact posterior inference in a tractable manner. In particular, if the prior is **conjugate** to the likelihood, the posterior will be analytically tractable. In general, this will be the case when the prior and likelihood are from the same exponential family (Section 2.3). In particular, if the unknown variables are represented by  $\boldsymbol{\theta}$ , then we assume

$$p(\boldsymbol{\theta}) \propto \exp(\boldsymbol{\lambda}_0^\top \mathcal{T}(\boldsymbol{\theta})) \quad (7.8)$$

$$p(\mathbf{y}_i | \boldsymbol{\theta}) \propto \exp(\tilde{\boldsymbol{\lambda}}_i(\mathbf{y}_i)^\top \mathcal{T}(\boldsymbol{\theta})) \quad (7.9)$$

1 where  $\mathcal{T}(\boldsymbol{\theta})$  are the sufficient statistics, and  $\boldsymbol{\lambda}$  are the natural parameters. We can then compute the  
 2 posterior by just adding the natural parameters:  
 3

$$\frac{4}{5} p(\boldsymbol{\theta} | \mathbf{y}_{1:N}) = \exp(\boldsymbol{\lambda}_*^\top \mathcal{T}(\boldsymbol{\theta})) \quad (7.10)$$

$$\frac{6}{7} \boldsymbol{\lambda}_* = \boldsymbol{\lambda}_0 + \sum_{n=1}^N \tilde{\boldsymbol{\lambda}}_n(\mathbf{y}_n) \quad (7.11)$$

8 See Section 3.2 for details.

10 Another setting where we can compute the posterior exactly arises when the  $D$  unknown variables  
 11 are all discrete, each with  $K$  states; in this case, the integral for the normalizing constant becomes a  
 12 sum with  $K^D$  terms. In many cases,  $K^D$  will be too large to be tractable. However, if the distribution  
 13 satisfies certain conditional independence properties, as expressed by a probabilistic graphical model  
 14 (PGM), then we can write the joint as a product of local terms (see Chapter 4 and Chapter 4). This  
 15 lets us use dynamic programming to make the computation tractable. See Chapter 9 for details.

16

## 17 7.4 Approximate inference algorithms

19 For most probability models, we will not be able to compute marginals or posteriors exactly, so we  
 20 must resort to using **approximate inference**. There are many different algorithms, which trade off  
 21 speed, accuracy, simplicity, and generality. We briefly discuss some of these algorithms below. We  
 22 give more detail in the following chapters.

23

### 24 7.4.1 MAP estimation

25 The simplest approximate inference method is to compute the MAP estimate

$$\frac{27}{28} \hat{\boldsymbol{\theta}} = \operatorname{argmax} p(\boldsymbol{\theta} | \mathcal{D}) = \operatorname{argmax} \log p(\boldsymbol{\theta}) + \log p(\mathcal{D} | \boldsymbol{\theta}) \quad (7.12)$$

29 and then to assume that the posterior puts 100% of its probability on this single value:

$$\frac{30}{31} p(\boldsymbol{\theta} | \mathcal{D}) \approx \delta(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}) \quad (7.13)$$

32 The advantage of this approach is that we can compute the MAP estimate using a variety of  
 33 optimization algorithms, which we discuss in Chapter 6. The disadvantages of the MAP approximation  
 34 are discussed in Section 3.1.5.

35

### 36 7.4.2 Grid approximation

38 If we want to capture uncertainty, we need to allow for the fact that  $\boldsymbol{\theta}$  may have a range of possible  
 39 values, each with non-zero probability. The simplest way to capture this property is to partition  
 40 the space of possible values into a finite set of regions, call them  $\mathbf{r}_1, \dots, \mathbf{r}_K$ , each representing a  
 41 region of parameter space of volume  $\Delta$  centered on  $\boldsymbol{\theta}_k$ . This is called a **grid approximation**. The  
 42 probability of being in each region is given by  $p(\boldsymbol{\theta} \in \mathbf{r}_k | \mathcal{D}) \approx p_k \Delta$ , where

$$\frac{43}{44} p_k = \frac{\tilde{p}_k}{\sum_{k'=1}^K \tilde{p}_{k'}} \quad (7.14)$$

$$\frac{45}{46} \tilde{p}_k = p(\mathcal{D} | \boldsymbol{\theta}_k) p(\boldsymbol{\theta}_k) \quad (7.15)$$

47

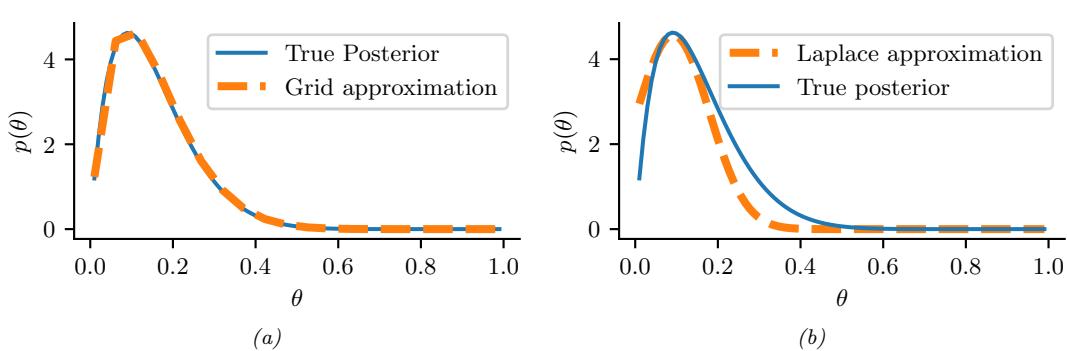


Figure 7.2: Approximating the posterior of a beta-Bernoulli model. (a) Grid approximation using 20 grid points. (b) Laplace approximation. Generated by [laplace\\_approx\\_beta\\_binom.ipynb](#).

As  $K$  increases, we decrease the size of each grid cell. Thus the denominator is just a simple numerical approximation of the integral

$$p(\mathcal{D}) = \int p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta} \approx \sum_{k=1}^K \Delta \tilde{p}_k \quad (7.16)$$

As a simple example, we will use the problem of approximating the posterior of a beta-Bernoulli model. Specifically, the goal is to approximate

$$p(\theta|\mathcal{D}) \propto \left[ \prod_{n=1}^N \text{Ber}(y_n|\theta) \right] \text{Beta}(1, 1) \quad (7.17)$$

where  $\mathcal{D}$  consists of 10 heads and 1 tail (so the total number of observations is  $N = 11$ ), with a uniform prior. Although we can compute this posterior exactly using the method discussed in Section 3.2.1, this serves as a useful pedagogical example since we can compare the approximation to the exact answer. Also, since the target distribution is just 1d, it is easy to visualize the results.

In Figure 7.2a, we illustrate the grid approximation applied to our 1d problem. We see that it is easily able to capture the skewed posterior (due to the use of an imbalanced sample of 10 heads and 1 tail). Unfortunately, this approach does not scale to problems in more than 2 or 3 dimensions, because the number of grid points grows exponentially with the number of dimensions.

### 7.4.3 Laplace (quadratic) approximation

In this section, we discuss a simple way to approximate the posterior using a multivariate Gaussian; this known as a **Laplace approximation** or **quadratic approximation** (see e.g., [TK86; RMC09]).

Suppose we write the posterior as follows:

$$p(\boldsymbol{\theta}|\mathcal{D}) = \frac{1}{Z} e^{-\mathcal{E}(\boldsymbol{\theta})} \quad (7.18)$$

where  $\mathcal{E}(\boldsymbol{\theta}) = -\log p(\boldsymbol{\theta}, \mathcal{D})$  is called an energy function, and  $Z = p(\mathcal{D})$  is the normalization constant. Performing a Taylor series expansion around the mode  $\hat{\boldsymbol{\theta}}$  (i.e., the lowest energy state) we get

$$\mathcal{E}(\boldsymbol{\theta}) \approx \mathcal{E}(\hat{\boldsymbol{\theta}}) + (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^T \mathbf{g} + \frac{1}{2} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^T \mathbf{H} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}) \quad (7.19)$$

where  $\mathbf{g}$  is the gradient at the mode, and  $\mathbf{H}$  is the Hessian. Since  $\hat{\boldsymbol{\theta}}$  is the mode, the gradient term is zero. Hence

$$\hat{p}(\boldsymbol{\theta}, \mathcal{D}) = e^{-\mathcal{E}(\hat{\boldsymbol{\theta}})} \exp \left[ -\frac{1}{2} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^T \mathbf{H} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}) \right] \quad (7.20)$$

$$\hat{p}(\boldsymbol{\theta} | \mathcal{D}) = \frac{1}{Z} \hat{p}(\boldsymbol{\theta}, \mathcal{D}) = \mathcal{N}(\boldsymbol{\theta} | \hat{\boldsymbol{\theta}}, \mathbf{H}^{-1}) \quad (7.21)$$

$$Z = e^{-\mathcal{E}(\hat{\boldsymbol{\theta}})} (2\pi)^{D/2} |\mathbf{H}|^{-\frac{1}{2}} \quad (7.22)$$

The last line follows from normalization constant of the multivariate Gaussian.

The Laplace approximation is easy to apply, since we can leverage existing optimization algorithms to compute the MAP estimate, and then we just have to compute the Hessian at the mode. (In high dimensional spaces, we can use a diagonal approximation.)

In Figure 7.2b, we illustrate this method applied to our 1d problem. Unfortunately we see that it is not a particularly good approximation. This is because the posterior is skewed, whereas a Gaussian is symmetric. In addition, the parameter of interest lies in the constrained interval  $\theta \in [0, 1]$ , whereas the Gaussian assumes an unconstrained space,  $\boldsymbol{\theta} \in \mathbb{R}^D$ . Fortunately, we can solve this latter problem by using a change of variable. For example, in this case we can apply the Laplace approximation to  $\alpha = \text{logit}(\theta)$ . This is a common trick to simplify the job of inference.

#### 7.4.4 Variational inference

In Section 7.4.3, we discussed the Laplace approximation, which uses an optimization procedure to find the MAP estimate, and then approximates the curvature of the posterior at that point based on the Hessian. In this section, we discuss **variational inference (VI)**, also called **variational Bayes (VB)**. This is another optimization-based approach to posterior inference, but which has much more modeling flexibility (and thus can give a much more accurate approximation).

VI attempts to approximate an intractable probability distribution, such as  $p(\boldsymbol{\theta} | \mathcal{D})$ , with one that is tractable,  $q(\boldsymbol{\theta})$ , so as to minimize some discrepancy  $D$  between the distributions:

$$q^* = \underset{q \in \mathcal{Q}}{\operatorname{argmin}} D(q, p) \quad (7.23)$$

where  $\mathcal{Q}$  is some tractable family of distributions (e.g., fully factorized distributions). Rather than optimizing over functions  $q$ , we typically optimize over the parameters of the function  $q$ ; we denote these **variational parameters** by  $\psi$ .

It is common to use the KL divergence (Section 5.1) as the discrepancy measure, which is given by

$$D(q, p) = D_{\text{KL}}(q(\boldsymbol{\theta} | \psi) \| p(\boldsymbol{\theta} | \mathcal{D})) = \int q(\boldsymbol{\theta} | \psi) \log \frac{q(\boldsymbol{\theta} | \psi)}{p(\boldsymbol{\theta} | \mathcal{D})} d\boldsymbol{\theta} \quad (7.24)$$

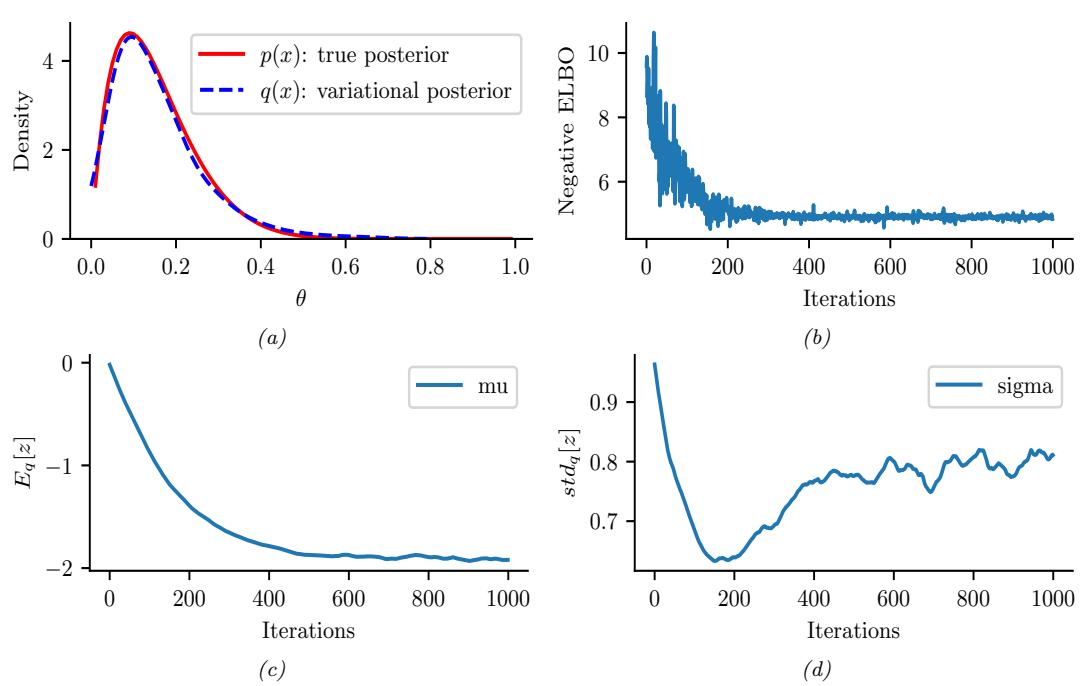


Figure 7.3: ADVI applied to the beta-Bernoulli model. (a) Approximate vs true posterior. (b) Negative ELBO over time. (c) Variational  $\mu$  parameter over time. (d) Variational  $\sigma$  parameter over time. Generated by `advi_beta_binom.ipynb`.

where  $p(\boldsymbol{\theta}|\mathcal{D}) = p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})/p(\mathcal{D})$ . The inference problem then reduces to the following optimization problem:

$$\psi^* = \operatorname{argmin}_{\psi} D_{\text{KL}}(q(\boldsymbol{\theta}|\psi) \| p(\boldsymbol{\theta}|\mathcal{D})) \quad (7.25)$$

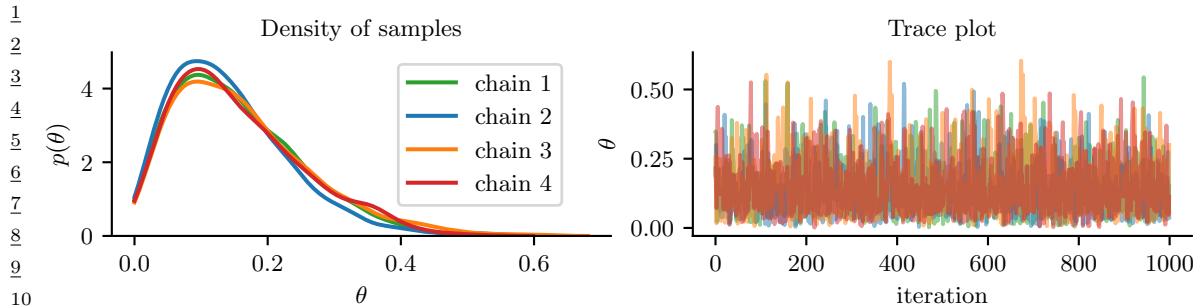
$$= \operatorname{argmin}_{\psi} \mathbb{E}_{q(\boldsymbol{\theta}|\psi)} \left[ \log q(\boldsymbol{\theta}|\psi) - \log \left( \frac{p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathcal{D})} \right) \right] \quad (7.26)$$

$$= \operatorname{argmin}_{\psi} \underbrace{\mathbb{E}_{q(\boldsymbol{\theta}|\psi)} [-\log p(\mathcal{D}|\boldsymbol{\theta}) - \log p(\boldsymbol{\theta}) + \log q(\boldsymbol{\theta}|\psi)]}_{-\mathcal{L}(\psi)} + \log p(\mathcal{D}) \quad (7.27)$$

Note that  $\log p(\mathcal{D})$  is independent of  $\psi$ , so we can ignore it when fitting the approximate posterior, and just focus on maximizing the term

$$\mathcal{L}(\psi) \triangleq \mathbb{E}_{q(\boldsymbol{\theta}|\psi)} [\log p(\mathcal{D}|\boldsymbol{\theta}) + \log p(\boldsymbol{\theta}) - \log q(\boldsymbol{\theta}|\psi)] \quad (7.28)$$

Since we have  $D_{\text{KL}}(q \| p) \geq 0$ , we have  $\mathcal{L}(\psi) \leq \log p(\mathcal{D})$ . The quantity  $\log p(\mathcal{D})$ , which is the log marginal likelihood, is also called the **evidence**. Hence  $\mathcal{L}(\psi)$  is known as the **evidence lower bound** or **ELBO**. By maximizing this bound, we are making the variational posterior closer to the true posterior. (See Section 10.1 for details.)



11 *Figure 7.4: Approximating the posterior of a beta-Bernoulli model using MCMC. (a) Kernel density estimate  
12 derived from samples from 4 independent chains. (b) Trace plot of the chains as they generate posterior  
13 samples. Generated by `hmc_beta_binom.ipynb`.*

14  
15  
16 We can chose any kind of approximate posterior that we like. For example, we may use a Gaussian,  
18  $q(\boldsymbol{\theta}|\boldsymbol{\psi}) = \mathcal{N}(\boldsymbol{\theta}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . This is different from the Laplace approximation, since in VI, we optimize  
19  $\boldsymbol{\Sigma}$ , rather than equating it to the Hessian. If  $\boldsymbol{\Sigma}$  is diagonal, we are assuming the posterior is fully  
20 factorized; this is called a **mean field** approximation.

21 A Gaussian approximation is not always suitable for all parameters. For example, in our 1d  
22 example we have the constraint that  $\theta \in [0, 1]$ . We could use a variational approximation of the form  
23  $q(\theta|\boldsymbol{\psi}) = \text{Beta}(\theta|a, b)$ , where  $\boldsymbol{\psi} = (a, b)$ . However choosing a suitable form of variational distribution  
24 requires some level of expertise. To create a more easily applicable, or “turn-key”, method, that works  
25 on a wide range of models, we can use a method called **automatic differentiation variational**  
26 **inference** or **ADVI** [Kuc+16]. This uses the change of variables method to convert the parameters  
27 to an unconstrained form, and then computes a Gaussian variational approximation. The method  
28 also uses automatic differentiation to derive the Jacobian term needed to compute the density of the  
29 transformed variables. See Section 10.3.5 for details.

30 We now apply ADVI to our 1d beta-Bernoulli model. Let  $\theta = \sigma(z)$ , where we replace  $p(\theta|\mathcal{D})$  with  
31  $q(z|\boldsymbol{\psi}) = \mathcal{N}(z|\mu, \sigma)$ , where  $\boldsymbol{\psi} = (\mu, \sigma)$ . We optimize a stochastic approximation to the ELBO using  
32 SGD. The results are shown in Figure 7.3 and seem reasonable.

33

#### 34 7.4.5 Markov Chain Monte Carlo (MCMC)

35  
36 Although VI is fast, it can give a biased approximation to the posterior, since it is restricted to a  
37 specific function form  $q \in \mathcal{Q}$ . A more flexible approach is to use a non-parametric approximation in  
38 terms of a set of samples,  $q(\boldsymbol{\theta}) \approx \frac{1}{S} \sum_{s=1}^S \delta(\boldsymbol{\theta} - \boldsymbol{\theta}^s)$ . This is called a **Monte Carlo approximation**.  
39 The key issue is how to create the posterior samples  $\boldsymbol{\theta}^s \sim p(\boldsymbol{\theta}|\mathcal{D})$  efficiently, without having to  
40 evaluate the normalization constant  $p(\mathcal{D}) = \int p(\boldsymbol{\theta}, \mathcal{D}) d\boldsymbol{\theta}$ .

41 For low dimensional problems, we can use methods such as **importance sampling**, which we  
42 discuss in Section 11.5. However, for high dimensional problems, it is more common to use **Markov**  
43 **chain Monte Carlo** or **MCMC**. We give the details in Chapter 12, but give a brief introduction  
44 here.

45 The most common kind of MCMC is known as the **Metropolis Hastings algorithm**. The basic  
46 idea behind MH is as follows: we start at a random point in parameter space, and then perform a  
47

random walk, by sampling new states (parameters) from a **proposal distribution**  $q(\boldsymbol{\theta}'|\boldsymbol{\theta})$ . If  $q$  is chosen carefully, the resulting Markov chain distribution will satisfy the property that the fraction of time we visit each point in space is proportional to the posterior probability. The key point is that to decide whether to move to a newly proposed point  $\boldsymbol{\theta}'$  or to stay in the current point  $\boldsymbol{\theta}$ , we only need to evaluate the unnormalized density ratio

$$\frac{p(\boldsymbol{\theta}|\mathcal{D})}{p(\boldsymbol{\theta}'|\mathcal{D})} = \frac{p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})/p(\mathcal{D})}{p(\mathcal{D}|\boldsymbol{\theta}')p(\boldsymbol{\theta}')/p(\mathcal{D})} = \frac{p(\mathcal{D}, \boldsymbol{\theta})}{p(\mathcal{D}, \boldsymbol{\theta}')} \quad (7.29)$$

This avoids the need to compute the normalization constant  $p(\mathcal{D})$ . (In practice we usually work with log probabilities, instead of joint probabilities, to avoid numerical issues.)

We see that the input to the algorithm is just a function that computes the log joint density,  $\log p(\boldsymbol{\theta}, \mathcal{D})$ , as well as a proposal distribution  $q(\boldsymbol{\theta}'|\boldsymbol{\theta})$  for deciding which states to visit next. It is common to use a Gaussian distribution for the proposal,  $q(\boldsymbol{\theta}'|\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}'|\boldsymbol{\theta}, \sigma\mathbf{I})$ ; this is called the **random walk Metropolis** algorithm. However, this can be very inefficient, since it is blindly walking through the space, in the hopes of finding higher probability regions.

In models that have conditional independence structure, it is often easy to compute the **full conditionals**  $p(\boldsymbol{\theta}_d|\boldsymbol{\theta}_{-d}, \mathcal{D})$  for each variable  $d$ , one at a time, and then sample from them. This is like a stochastic analog of coordinate ascent, and is called **Gibbs sampling** (see Section 12.3 for details).

For models where all unknown variables are continuous, we can often compute the gradient of the log joint,  $\nabla_{\boldsymbol{\theta}} \log p(\boldsymbol{\theta}, \mathcal{D})$ . We can use this gradient information to guide the proposals into regions of space with higher probability. This approach is called **Hamiltonian Monte Carlo** or **HMC**, and is one of the most widely used MCMC algorithms due to its speed. For details, see Section 12.5.

We apply HMC to our beta-Bernoulli model in Figure 7.4. (We use a logit transformation for the parameter.) In panel b, we show samples generated by the algorithm from 4 parallel Markov chains. We see that they oscillate around the true posterior, as desired. In panel a, we compute a kernel density estimate from the posterior samples from each chain; we see that the result is a good approximation to the true posterior in Figure 7.2.

### 7.4.6 Sequential Monte Carlo

MCMC is like a stochastic local search algorithm, in that it makes moves through the state space of the posterior distribution, comparing the current value to proposed neighboring values. An alternative approach is to use perform inference using a sequence of different distributions, from simpler to more complex, with the final distribution being equal to the target posterior. This is called **sequential Monte Carlo** or **SMC**. This approach, which is more similar to tree search than local search, has various advantages over MCMC, which we discuss in Chapter 13.

A common application of SMC is to **sequential Bayesian inference**, in which we recursively compute (i.e., in an online fashion) the posterior  $p(\boldsymbol{\theta}_t|\mathcal{D}_{1:t})$ , where  $\mathcal{D}_{1:t} = \{(\mathbf{x}_n, y_n) : n = 1 : t\}$  is all the data we have seen so far. This sequence of distributions converges to the full batch posterior  $p(\boldsymbol{\theta}|\mathcal{D})$  once all the data has been seen. However, the approach can also be used when the data is arriving in a continual, unending stream, as in state-space models (see Chapter 29). The application of SMC to such dynamical models is known as **particle filtering**. See Section 13.2 for details.

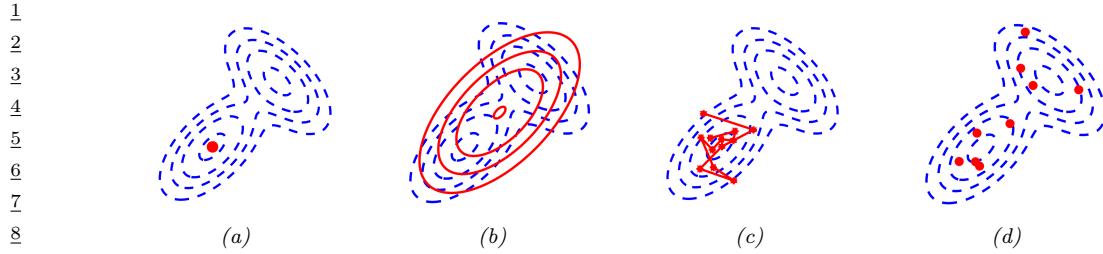


Figure 7.5: Different approximations to a bimodal 2d distribution. (a) Local MAP estimate. (b) Parametric Gaussian approximation. (c) Correlated samples from near one mode. (d) Independent samples from the distribution. Adapted from Figure 2 of [PY14]. Used with kind permission of George Panadreou.

#### 7.4.7 Challenging posteriors

In many applications, the posterior can be high dimensional and multimodal. Approximating such distributions can be quite challenging. In Figure 7.5, we give a simple 2d example. We compare MAP estimation (which does not capture any uncertainty), a Gaussian parametric approximation such as the Laplace approximation or variational inference, and a nonparametric approximation in terms of samples. If the samples are generated from MCMC, they are serially correlated, and may only explore a local model (see panel b). However, ideally we can draw independent samples from the entire support of the distribution, as shown in panel d. We may also be able to fit a local parametric approximation around each such sample (c.f., Section 17.3.9.1), to get a semi-parametric approximation to the posterior.

### 7.5 Evaluating approximate inference algorithms

There are many different inference algorithms each of which make different tradeoffs between speed, accuracy, generality, simplicity, etc. This makes it hard to compare them on an equal footing. However, a common approach is to evaluate the accuracy of the approximation as a function of compute time.

We give an example of this in Figure 7.6, where we plot performance for several different inference algorithms applied to the following simple 1d model:

$$p(z) = \mathcal{N}(z|0, 100) \quad (7.30)$$

$$p(x_i|z) = 0.5\mathcal{N}(x_i|z, 1) + 0.5\mathcal{N}(x_i|0, 10) \quad (7.31)$$

That is, each measurement  $x_i$  is either a noisy copy of the hidden quantity of interest,  $z$ , or is an outlier coming from a uniform background noise model, approximated by a Gaussian with a large variance,  $\mathcal{N}(0, 10)$ . The goal is to compute  $p(z|\mathbf{x}_{1:N})$ . (Tom Minka (who created this example) calls this the **clutter problem**.)

Since this is a 1d problem, we can compute the exact answer using numerical integration. In Figure 7.6, we plot the accuracy of the posterior mean estimate using various approximate inference methods, namely: the Laplace approximation (Section 7.4.3), variational Bayes (Section 10.2.3), expectation propagation (Section 10.7), importance sampling (Section 11.5), and Gibbs sampling (Section 12.3). We see that the error smoothly decreases for the 3 deterministic methods (Laplace,

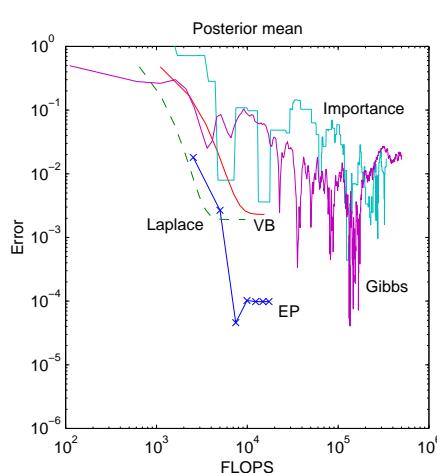


Figure 7.6: Accuracy vs compute time for different inference methods. Code source: <https://github.com/tminka/ep-clutter-example>. From Figure 1 of [Min01b]. Used with kind permission of Tom Minka.

VB and EP). For the 2 stochastic methods (IS and Gibbs), the error decreases on average, but the performance is quite noisy.

In principle, the Monte Carlo methods will converge to a zero overall error, since they are unbiased estimators, but it is clear that their variance is quite high. The deterministic methods, by comparison, converge to a finite error, so they are biased, but their variance is much lower. We can trade off bias and variance by creating hybrid algorithms, that combine different techniques. We will see examples of this in later chapters.

When we cannot compute the true target posterior distribution to compare to, evaluation is harder, but there are some approaches than can be used in certain cases. For some methods for assessing variational inference, see e.g., [Yao+18b; Hug+20]. For some methods for assessing Monte Carlo methods, see [CGR06; CTM17; GAR16]. For assessing posterior predictive distributions, see Section 14.2.3.



# 8 Inference for state-space models

## 8.1 Introduction

In this chapter, we consider the task of posterior inference in **state-space models** (SSM). We discuss SSMs in more detail in Chapter 29, but we can think of them as latent variable sequence models with the conditional independencies shown by the chain-structured graphical model Figure 8.1. The corresponding joint distribution has the form

$$p(\mathbf{y}_{1:T}, \mathbf{z}_{1:T} | \mathbf{u}_{1:T}) = \left[ p(\mathbf{z}_1 | \mathbf{u}_1) \prod_{t=2}^T p(\mathbf{z}_t | \mathbf{z}_{t-1}, \mathbf{u}_t) \right] \left[ \prod_{t=1}^T p(\mathbf{y}_t | \mathbf{z}_t, \mathbf{u}_t) \right] \quad (8.1)$$

where  $\mathbf{z}_t$  are the hidden variables at time  $t$ ,  $\mathbf{y}_t$  are the observations (outputs), and  $\mathbf{u}_t$  are the optional inputs.

Given the sequence of observations, and a known model, one of the main tasks with SSMs is to perform posterior inference about the hidden states; this is also called **state estimation**. At each time step  $t$ , there are multiple forms of posterior we may be interested in computing, including the **filtering distribution**  $p(\mathbf{z}_t | \mathbf{y}_{1:t})$ , the **smoothing distribution**  $p(\mathbf{z}_t | \mathbf{y}_{1:T})$  (note that this conditions on future data  $T \geq t$ ), and the **fixed-lag smoothing distribution**  $p(\mathbf{z}_{t-\ell} | \mathbf{y}_{1:t})$  (note that this infers  $\ell$  steps in the past given data up to the present). We may also want to compute the **predictive distribution**  $h$  steps into the future:

$$p(\mathbf{y}_{t+h} | \mathbf{y}_{1:t}) = \sum_{\mathbf{z}_{t+h}} p(\mathbf{y}_{t+h} | \mathbf{z}_{t+h}) p(\mathbf{z}_{t+h} | \mathbf{y}_{1:t}) \quad (8.2)$$

where the hidden state predictive distribution is

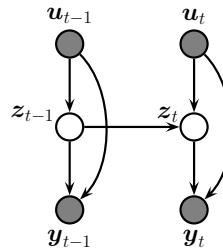
$$p(\mathbf{z}_{t+h} | \mathbf{y}_{1:t}) = \sum_{\mathbf{z}_{t:t+h-1}} p(\mathbf{z}_t | \mathbf{y}_{1:t}) p(\mathbf{z}_{t+1} | \mathbf{z}_t) p(\mathbf{z}_{t+2} | \mathbf{z}_{t+1}) \cdots p(\mathbf{z}_{t+h} | \mathbf{z}_{t+h-1}) \quad (8.3)$$

See Figure 8.2 for a summary of these distributions. In addition to computing posterior marginals, we may want to compute the most probable hidden sequence,  $\text{argmax}_{\mathbf{z}_{1:T}} p(\mathbf{z}_{1:T} | \mathbf{y}_{1:T})$ , or sample sequences from the posterior,  $\mathbf{z}_{1:T} \sim p(\mathbf{z}_{1:T} | \mathbf{y}_{1:T})$ . We discuss algorithms for all these tasks later in this chapter. See also [Sar13; Tri21].

## 8.2 Inference based on the HMM filter

In this section, we consider inference for SSMs where all the hidden variables are discrete. This is known as a **hidden Markov model** or **HMM**. This is discussed in detail in Section 29.2, but

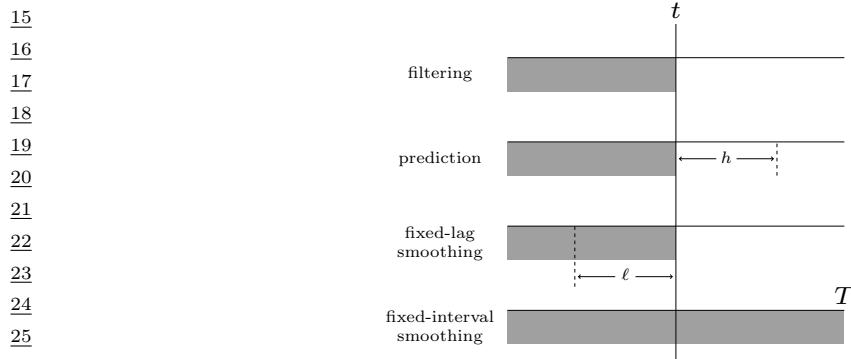
1  
2  
3  
4  
5  
6  
7  
8  
9  
10



11 *Figure 8.1: A state-space model represented as a graphical model.  $z_t$  are the hidden variables at time  $t$ ,  $y_t$  are  
12 the observations (outputs), and  $u_t$  are the optional inputs.*

13

14



27 *Figure 8.2: The main kinds of inference for state-space models. The shaded region is the interval for which  
28 we have data. The arrow represents the time step at which we want to perform inference.  $t$  is the current  
29 time,  $T$  is the sequence length,  $\ell$  is the lag and  $h$  is the prediction horizon. Used with kind permission of  
30 Peter Chang.*

31

32

33 in brief, it corresponds to the model in Equation (8.1) where  $p(z_1 = k) = \pi_k$  is the initial state  
34 distribution,  $p(z_t = k | z_{t-1} = j) = A_{jk}$  is the state transition matrix (assumed stationary), and  
35  $p(y_t | z_t = k)$  is some conditional distribution over observations. For notational simplicity, we assume  
36 all conditional distributions are stationary (the same over time), and we ignore inputs. Note, however,  
37 that the algorithms we discuss can be generalized to the non-stationary case, and also to 1d  
38 undirected graphical models, such as linear-chain CRFs (Section 4.4).

39

#### 40 8.2.1 Example: casino HMM

41

42 Before explaining the algorithms, we introduce a simple example of an HMM, which we will use to  
43 illustrate the methods.

44 Suppose we are in a casino and observe a series of die rolls,  $y_t \in \{1, 2, \dots, 6\}$ . Being a keen-eyed  
45 statistician, we notice that the distribution of values is not what we expect from a fair die: it seems  
46 that there are occasional “streaks”, in which 6s seem to show up more often than other values. (This  
47

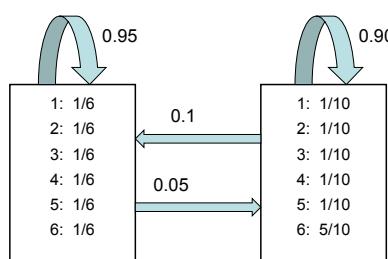


Figure 8.3: The state transition matrix  $\mathbf{A}$  and observation matrix  $\mathbf{B}$  for the casino HMM. Adapted from [Dur+98, p54].

example is from [Dur+98], who call this setup the **occasionally dishonest casino**.) We would like to **segment** the time series into regimes corresponding to the use of a fair die and a loaded die.

Let  $A_{jk} = p(z_t = k | z_{t-1} = j)$  be the state transition matrix, and  $B_{kl} = p(y_t = l | z_t = k)$  be the observation matrix corresponding to a categorical distribution over values of the die face. Most of the time the casino uses a fair die,  $z = 1$ , but occasionally it switches to a loaded die,  $z = 2$ , for a short period. If  $z = 1$  the observation distribution is a uniform categorical distribution over the symbols  $\{1, \dots, 6\}$ . If  $z = 2$ , the observation distribution is skewed towards face 6. That is,

$$p(y_t | z_t = 1) = \text{Cat}(y_t | [1/6, \dots, 1/6]) \quad (8.4)$$

$$p(y_t | z_t = 2) = \text{Cat}(y_t | [1/10, 1/10, 1/10, 1/10, 1/10, 5/10]) \quad (8.5)$$

See Figure 8.3 for an illustration of the state transition diagram  $\mathbf{A}$  and the emission distributions  $\mathbf{B}$ . If we sample from this model, we may generate data such as the following:

```

hid: 11111111122221111111111111111111222222222122221111111111111111111111111111111
obs: 1355534526553366316351551526232112113462221263264265422344645323242361

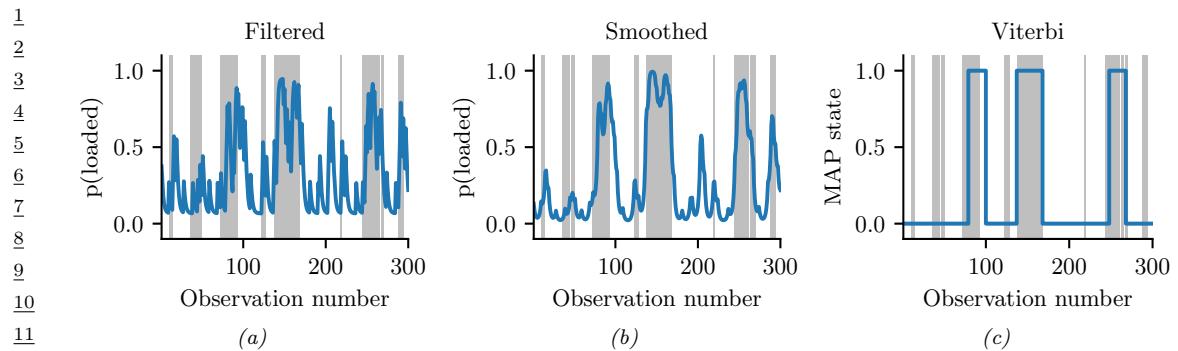
```

Here obs refers to the observation and hid refers to the hidden state (1 is fair and 2 is loaded). In the full sequence of length 300, we find the empirical fraction of times that we observe a 6 in hidden state 1 to be 0.149, and in state 2 to be 0.472, which is very close to the expected fractions. (See [casino\\_hmm.ipynb](#) for the code.)

Of course, when using an HMM in practice, we just see the observed sequence, and need to infer the hidden sequence, as we discuss below.

## 8.2.2 Forwards filtering

The **Bayes filter** is an algorithm for recursively computing the **belief state**  $p(z_t | y_{1:t})$  given the prior belief from the previous step,  $p(z_{t-1} | y_{1:t-1})$ , the new observation  $y_t$ , and the model. This can be done using **sequential Bayesian updating**. For a dynamical model, this reduces to the **predict-update** cycle described below.



*Figure 8.4: Inference in the dishonest casino. Vertical gray bars denote times when the hidden state corresponded to the loaded die. Blue lines represent the posterior probability of being in that state given different subsets of observed data. If we recover the true state exactly, the blue curve will transition at the same time as the gray bars. (a) Filtered estimate. (b) Smoothed estimates. (c) MAP trajectory. Generated by [casino\\_hmm.ipynb](#).*

### 8.2.2.1 Prediction step

The **prediction step** is just the **Chapman-Kolmogorov equation**:

$$p(\mathbf{z}_t | \mathbf{y}_{1:t-1}) = \int p(\mathbf{z}_t | \mathbf{z}_{t-1}) p(\mathbf{z}_{t-1} | \mathbf{y}_{1:t-1}) d\mathbf{z}_{t-1} \quad (8.6)$$

The prediction step computes the one-step-ahead predictive distribution for the latent state, which updates the posterior from the previous time step into the prior for the current step.<sup>1</sup>

### 8.2.2.2 Update step

The **update step** is just Bayes rule:

$$p(\mathbf{z}_t | \mathbf{y}_{1:t}) = \frac{1}{Z_t} p(\mathbf{y}_t | \mathbf{z}_t) p(\mathbf{z}_t | \mathbf{y}_{1:t-1}) \quad (8.7)$$

where the normalization constant is

$$Z_t = \int p(\mathbf{y}_t | \mathbf{z}_t) p(\mathbf{z}_t | \mathbf{y}_{1:t-1}) d\mathbf{z}_t = p(\mathbf{y}_t | \mathbf{y}_{1:t-1}) \quad (8.8)$$

We can use the normalization constants to compute the log likelihood of the sequence as follows:

$$\log p(\mathbf{y}_{1:T}) = \sum_{t=1}^T \log p(\mathbf{y}_t | \mathbf{y}_{1:t-1}) = \sum_{t=1}^T \log Z_t \quad (8.9)$$

where we define  $p(\mathbf{y}_1 | \mathbf{y}_0) = p(\mathbf{y}_1)$ . This is useful for computing the MLE of the parameters (see Section 29.4.2).

<sup>1</sup> 1. The prediction step is not needed at  $t = 1$  if  $p(\mathbf{z}_1)$  is provided as input to the model. However, if we just provide  $p(\mathbf{z}_0)$ , we need to compute  $p(\mathbf{z}_1 | \mathbf{y}_{1:0}) = p(\mathbf{z}_1)$  by applying the prediction step.

1

### 8.2.2.3 Implementation

2

When the latent states  $\mathbf{z}_t$  are discrete, the above integrals become sums which can be implemented  
3 as follows. First we define the belief state as  $\alpha_t(j) \triangleq p(z_t = j | \mathbf{y}_{1:t})$ , the local evidence as  $\lambda_t(j) \triangleq$   
4  $p(\mathbf{y}_t | z_t = j)$ , and the transition matrix  $A_{i,j} = p(z_t = j | z_{t-1} = i)$ . Then the predict step becomes  
5

6

$$\alpha_{t|t-1}(j) \triangleq p(z_t = j | \mathbf{y}_{1:t-1}) = \sum_i \alpha_{t-1}(i) A_{i,j} \quad (8.10)$$

7

8

and the update step becomes

9

10

$$\alpha_t(j) = \frac{1}{Z_t} \lambda_t(j) \alpha_{t|t-1}(j) = \frac{1}{Z_t} \lambda_t(j) \left[ \sum_i \alpha_{t-1}(i) A_{i,j} \right] \quad (8.11)$$

11

12

where the normalization constant for each time step is given by

13

14

$$Z_t \triangleq p(\mathbf{y}_t | \mathbf{y}_{1:t-1}) = \sum_{j=1}^K p(\mathbf{y}_t | z_t = j) p(z_t = j | \mathbf{y}_{1:t-1}) = \sum_{j=1}^K \lambda_t(j) \alpha_{t|t-1}(j) \quad (8.12)$$

15

16

We can write the update equation in matrix-vector notation as follows:

17

18

$$\boldsymbol{\alpha}_t = \text{normalize} (\boldsymbol{\lambda}_t \odot (\mathbf{A}^\top \boldsymbol{\alpha}_{t-1})) \quad (8.13)$$

19

20

where  $\odot$  represents elementwise vector multiplication, and the normalize function just ensures its  
21 argument sums to one. (See Section 8.2.5 for more discussion on normalization.)

22

23

### 8.2.2.4 Example

24

25

Figure 8.4(a) illustrates filtering for the casino HMM, applied to a random sequence  $\mathbf{y}_{1:T}$  of length  
26  $T = 300$ . In blue, we plot the probability that the die is in the loaded (vs fair) state, based on  
27 the evidence seen so far. The gray bars indicate time intervals during which the generative process  
28 actually switched to the loaded die. We see that the probability generally increases in the right  
29 places.

30

31

### 8.2.3 Backwards smoothing

32

33

In the offline setting, we want to compute  $p(\mathbf{z}_t | \mathbf{y}_{1:T})$ , which is the belief about the hidden state at  
34 time  $t$  given all the data, both past and future. This is called (fixed interval) **smoothing**. We first  
35 perform the forwards or filtering pass, and then compute the smoothed belief states by working  
36 backwards, from right (time  $t = T$ ) to left (to  $t = 1$ ), as we explain below. Hence this method is also  
37 called **forwards filtering backwards smoothing** or **FFBS**.

38

39

#### 8.2.3.1 Backwards recursion

40

41

Suppose, by induction, that we have already computed  $p(\mathbf{z}_{t+1} | \mathbf{y}_{1:T})$ . We can convert this into a joint  
42 smoothed distribution over two consecutive time steps using

43

44

$$p(\mathbf{z}_t, \mathbf{z}_{t+1} | \mathbf{y}_{1:T}) = p(\mathbf{z}_t | \mathbf{z}_{t+1}, \mathbf{y}_{1:T}) p(\mathbf{z}_{t+1} | \mathbf{y}_{1:T}) \quad (8.14)$$

45

1 To derive the first term, note that from the Markov properties of the model, and Bayes rule, we have  
2

$$\underline{3} \quad p(\mathbf{z}_t = i | \mathbf{z}_{t+1} = j, \mathbf{y}_{1:T}) = p(\mathbf{z}_t = i | \mathbf{z}_{t+1} = j, \mathbf{y}_{1:t}, \mathbf{y}_{t+1:T}) \quad (8.15)$$

$$\underline{4} \quad = \frac{p(\mathbf{z}_t = i, \mathbf{z}_{t+1} = j | \mathbf{y}_{1:t})}{p(\mathbf{z}_{t+1} = j | \mathbf{y}_{1:t})} \quad (8.16)$$

$$\underline{5} \quad = \frac{p(\mathbf{z}_{t+1} = j | \mathbf{z}_t = i) p(\mathbf{z}_t = i | \mathbf{y}_{1:t})}{p(\mathbf{z}_{t+1} = j | \mathbf{y}_{1:t})} \quad (8.17)$$

6 Thus the joint distribution over two consecutive time steps is given by  
7

$$\underline{8} \quad p(\mathbf{z}_t, \mathbf{z}_{t+1} | \mathbf{y}_{1:T}) = p(\mathbf{z}_t | \mathbf{z}_{t+1}, \mathbf{y}_{1:t}) p(\mathbf{z}_{t+1} | \mathbf{y}_{1:T}) = \frac{p(\mathbf{z}_{t+1} | \mathbf{z}_t) p(\mathbf{z}_t | \mathbf{y}_{1:t}) p(\mathbf{z}_{t+1} | \mathbf{y}_{1:T})}{p(\mathbf{z}_{t+1} | \mathbf{y}_{1:t})} \quad (8.18)$$

9 from which we get the new smoothed marginal distribution:  
10

$$\underline{11} \quad p(\mathbf{z}_t | \mathbf{y}_{1:T}) = p(\mathbf{z}_t | \mathbf{y}_{1:t}) \int \left[ \frac{p(\mathbf{z}_{t+1} | \mathbf{z}_t) p(\mathbf{z}_{t+1} | \mathbf{y}_{1:T})}{p(\mathbf{z}_{t+1} | \mathbf{y}_{1:t})} \right] d\mathbf{z}_{t+1} \quad (8.19)$$

$$\underline{12} \quad = \int p(\mathbf{z}_t, \mathbf{z}_{t+1} | \mathbf{y}_{1:t}) \frac{p(\mathbf{z}_{t+1} | \mathbf{y}_{1:T})}{p(\mathbf{z}_{t+1} | \mathbf{y}_{1:t})} d\mathbf{z}_{t+1} \quad (8.20)$$

13 Intuitively we can interpret this as follows: we start with the two-slice filtered distribution,  
14  $p(\mathbf{z}_t, \mathbf{z}_{t+1} | \mathbf{y}_{1:t})$ , and then we divide out the old  $p(\mathbf{z}_{t+1} | \mathbf{y}_{1:t})$  and multiply in the new  $p(\mathbf{z}_{t+1} | \mathbf{y}_{1:T})$ ,  
15 and then marginalize out  $\mathbf{z}_{t+1}$ .  
16

### 17 8.2.3.2 Implementation

18 Let us define the two-slice marginal using the following notation:  
19

$$\underline{20} \quad \xi_{t,t+1}(i, j) \triangleq p(\mathbf{z}_t = i, \mathbf{z}_{t+1} = j | \mathbf{y}_{1:T}) \quad (8.21)$$

21 We can then rewrite Equation (8.18) as follows:  
22

$$\underline{23} \quad \xi_{t,t+1}(i, j) = \alpha_t(i) A_{i,j} \frac{\gamma_{t+1}(j)}{\alpha_{t+1|t}(j)} \quad (8.22)$$

24 where  
25

$$\underline{26} \quad \alpha_{t+1|t}(j) = p(\mathbf{z}_{t+1} = j | \mathbf{y}_{1:t}) = \sum_{i'} A(i', j) \alpha_t(i') \quad (8.23)$$

27 is the one-step-ahead predictive distribution. We can interpret the ratio in Equation (8.22) as dividing  
28 out the old estimate of  $\mathbf{z}_{t+1}$  given  $\mathbf{y}_{1:t}$ , namely  $\alpha_{t+1|t}$ , and multiplying in the new estimate given  
29  $\mathbf{y}_{1:T}$ , namely  $\gamma_{t+1}$ .  
30

31 Finally we can recover the one-slice posterior marginal by marginalizing:  
32

$$\underline{33} \quad \gamma_t(i) \triangleq p(\mathbf{z}_t = i | \mathbf{y}_{1:T}) = \sum_j \xi_{t,t+1}(i, j) = \alpha_t(i) \sum_j \left[ A_{i,j} \frac{\gamma_{t+1}(j)}{\alpha_{t+1|t}(j)} \right] \quad (8.24)$$

34 We initialize the recursion using  $\gamma_T(j) = \alpha_T(j) = p(\mathbf{z}_T = j | \mathbf{y}_{1:T})$ .  
35

---

### 8.2.3.3 Example

In Figure 8.4(a-b), we compare filtering and smoothing for the casino HMM. We see that the posterior distributions when conditioned on all the data (past and future) are indeed smoother than when just conditioned on the past (filtering).

To understand this behavior intuitively, consider a detective trying to figure out who committed a crime. As they move through the crime scene, his uncertainty is high until he finds the key clue; then he has an “aha” moment, his uncertainty is reduced, and all the previously confusing observations are, in **hindsight**, easy to explain. Thus we see that, given all the data (including finding the clue), it is much easier to infer the state of the world.

### 8.2.4 The forwards-backwards algorithm

In this section, we present a more common approach to smoothing in HMMs known as the **forwards-backwards** or **FB** algorithm [Rab89]. In the forwards pass, we compute  $\alpha_t(j) = p(\mathbf{z}_t = j | \mathbf{y}_{1:t})$  as before. In the backwards pass, we compute the conditional likelihood

$$\beta_t(j) \triangleq p(\mathbf{y}_{t+1:T} | \mathbf{z}_t = j) \quad (8.25)$$

We then combine these using

$$\gamma_t(j) = p(\mathbf{z}_t = j | \mathbf{y}_{t+1:T}, \mathbf{y}_{1:t}) \propto p(\mathbf{z}_t = j, \mathbf{y}_{t+1:T} | \mathbf{y}_{1:t}) \quad (8.26)$$

$$= p(\mathbf{z}_t = j | \mathbf{y}_{1:t}) p(\mathbf{y}_{t+1:T} | \mathbf{z}_t = j, \mathbf{y}_{1:t}) = \alpha_t(j) \beta_t(j) \quad (8.27)$$

In matrix notation, this becomes

$$\boldsymbol{\gamma}_t = \text{normalize}(\boldsymbol{\alpha}_t \odot \boldsymbol{\beta}_t) \quad (8.28)$$

Note that the forwards and backwards passes can be computed independently, but both need access to the local evidence  $p(\mathbf{y}_t | \mathbf{z}_t)$ . The results are only combined at the end. This is therefore called **two-filter smoothing** [Kit04].

#### 8.2.4.1 Backwards recursion

We can recursively compute the  $\beta$ 's in a right-to-left fashion as follows:

$$\beta_{t-1}(i) = p(\mathbf{y}_{t:T} | \mathbf{z}_{t-1} = i) \quad (8.29)$$

$$= \sum_j p(\mathbf{z}_t = j, \mathbf{y}_t, \mathbf{y}_{t+1:T} | \mathbf{z}_{t-1} = i) \quad (8.30)$$

$$= \sum_j p(\mathbf{y}_{t+1:T} | \mathbf{z}_t = j, \mathbf{y}_t, \mathbf{z}_{t-1} = i) p(\mathbf{z}_t = j, \mathbf{y}_t | \mathbf{z}_{t-1} = i) \quad (8.31)$$

$$= \sum_j p(\mathbf{y}_{t+1:T} | \mathbf{z}_t = j) p(\mathbf{y}_t | \mathbf{z}_t = j, \mathbf{z}_{t-1} = i) p(\mathbf{z}_t = j | \mathbf{z}_{t-1} = i) \quad (8.32)$$

$$= \sum_j \beta_t(j) \lambda_t(j) A_{i,j} \quad (8.33)$$

1 We can write the resulting equation in matrix-vector form as  
2

$$\underline{3} \quad \beta_{t-1} = \mathbf{A}(\boldsymbol{\lambda}_t \odot \boldsymbol{\beta}_t) \quad (8.34)$$

4 The base case is  
5

$$\underline{6} \quad \beta_T(i) = p(\mathbf{y}_{T+1:T} | \mathbf{z}_T = i) = p(\emptyset | \mathbf{z}_T = i) = 1 \quad (8.35)$$

7 which is the probability of a non-event.  
8

9 Note that  $\boldsymbol{\beta}_t$  is not a probability distribution over states, since it does not need to satisfy  
10  $\sum_j \beta_t(j) = 1$ . However, we usually normalize it to avoid numerical underflow (see Section 8.2.5).  
11

#### 12 8.2.4.2 Two-slice smoothed marginals

13 We can compute the two-slice marginals using the output of the forwards-backwards algorithm as  
14 follows:  
15

$$\underline{16} \quad p(\mathbf{z}_t, \mathbf{z}_{t+1} | \mathbf{y}_{1:T}) = p(\mathbf{z}_t, \mathbf{z}_{t+1} | \mathbf{y}_{1:t}, \mathbf{y}_{t+1:T}) \quad (8.36)$$

$$\underline{17} \quad \propto p(\mathbf{y}_{t+1:T} | \mathbf{z}_t, \mathbf{z}_{t+1}, \mathbf{y}_{1:t}) p(\mathbf{z}_t, \mathbf{z}_{t+1} | \mathbf{y}_{1:t}) \quad (8.37)$$

$$\underline{18} \quad = p(\mathbf{y}_{t+1:T} | \mathbf{z}_{t+1}) p(\mathbf{z}_t, \mathbf{z}_{t+1} | \mathbf{y}_{1:t}) \quad (8.38)$$

$$\underline{19} \quad = p(\mathbf{y}_{t+1:T} | \mathbf{z}_{t+1}) p(\mathbf{z}_t | \mathbf{y}_{1:t}) p(\mathbf{z}_{t+1} | \mathbf{z}_t) \quad (8.39)$$

$$\underline{20} \quad = p(\mathbf{y}_{t+1}, \mathbf{y}_{t+2:T} | \mathbf{z}_{t+1}) p(\mathbf{z}_t | \mathbf{y}_{1:t}) p(\mathbf{z}_{t+1} | \mathbf{z}_t) \quad (8.40)$$

$$\underline{21} \quad = p(\mathbf{y}_{t+1} | \mathbf{z}_{t+1}) p(\mathbf{y}_{t+2:T} | \mathbf{z}_{t+1}, \mathbf{y}_{t+1}) p(\mathbf{z}_t | \mathbf{y}_{1:t}) p(\mathbf{z}_{t+1} | \mathbf{z}_t) \quad (8.41)$$

$$\underline{22} \quad = p(\mathbf{y}_{t+1} | \mathbf{z}_{t+1}) p(\mathbf{y}_{t+2:T} | \mathbf{z}_{t+1}) p(\mathbf{z}_t | \mathbf{y}_{1:t}) p(\mathbf{z}_{t+1} | \mathbf{z}_t) \quad (8.42)$$

23 We can rewrite this in terms of the already computed quantities as follows:  
24

$$\underline{25} \quad \xi_{t,t+1}(i, j) \propto \lambda_{t+1}(j) \beta_{t+1}(j) \alpha_t(i) A_{i,j} \quad (8.43)$$

26 Or in matrix-vector form:  
27

$$\underline{28} \quad \xi_{t,t+1} \propto \mathbf{A} \odot [\boldsymbol{\alpha}_t (\boldsymbol{\lambda}_{t+1} \odot \boldsymbol{\beta}_{t+1})^\top] \quad (8.44)$$

29 Since  $\boldsymbol{\alpha}_t \propto \boldsymbol{\lambda}_t \odot \boldsymbol{\alpha}_{t|t-1}$ , we can also write the above equation as follows:  
30

$$\underline{31} \quad \xi_{t,t+1} \propto \mathbf{A} \odot (\boldsymbol{\lambda}_t \odot \boldsymbol{\alpha}_{t|t-1}) \odot (\boldsymbol{\lambda}_{t+1} \odot \boldsymbol{\beta}_{t+1})^\top \quad (8.45)$$

32 This can be interpreted as a product of incoming messages and local factors, as shown in Figure 8.5.  
33

34 In particular, we combine the factors  $\boldsymbol{\alpha}_{t|t-1} = p(\mathbf{z}_t | \mathbf{y}_{1:t-1})$ ,  $\mathbf{A} = p(\mathbf{z}_{t+1} | \mathbf{z}_t)$ ,  $\boldsymbol{\lambda}_t \propto p(\mathbf{y}_t | \mathbf{z}_t)$ ,  $\boldsymbol{\lambda}_{t+1} \propto p(\mathbf{y}_{t+1} | \mathbf{z}_{t+1})$ , and  $\boldsymbol{\beta}_{t+1} \propto p(\mathbf{y}_{t+2:T} | \mathbf{z}_{t+1})$  to get  $p(\mathbf{z}_t, \mathbf{z}_{t+1}, \mathbf{y}_t, \mathbf{y}_{t+1}, \mathbf{y}_{t+2:T} | \mathbf{y}_{1:t-1})$ , which we can then  
35 normalize.  
36

#### 37 8.2.5 Numerically stable implementation

38 In most publications on HMMs, such as [Rab89], the forwards message is defined as the following  
39 unnormalized joint probability:  
40

$$\underline{41} \quad \alpha'_t(j) = p(\mathbf{z}_t = j, \mathbf{y}_{1:t}) = \lambda_t(j) \left[ \sum_i \alpha'_{t-1}(i) A_{i,j} \right] \quad (8.46)$$

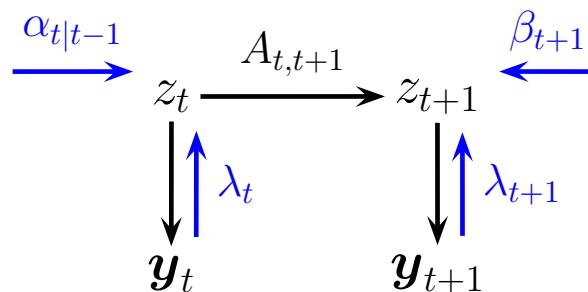


Figure 8.5: Computing the two-slice joint distribution for an HMM from the forwards messages, backwards messages, and local evidence messages.

We instead define the forwards message as the normalized *conditional* probability

$$\alpha_t(j) = p(z_t = j | y_{1:t}) = \frac{1}{Z_t} \lambda_t(j) \left[ \sum_i \alpha_{t-1}(i) A_{i,j} \right] \quad (8.47)$$

The unnormalized (joint) form has several problems. First, it rapidly suffers from numerical underflow, since the probability of the joint event that  $(z_t = j, y_{1:t})$  is vanishingly small.<sup>2</sup> Second, it is less interpretable, since it is not a distribution over states. Third, it precludes the use of approximate inference methods that try to approximate posterior distributions (we will see such methods later). We therefore always used the normalized (conditional) form.

Of course, the two definitions only differ by a multiplicative constant, since  $p(z_t = j | y_{1:t}) = p(z_t = j, y_{1:t}) / p(y_{1:t})$  [Dev85]. So the *algorithmic* difference is just one line of code (namely the presence or absence of a call to the `normalize` function). Nevertheless, we feel it is better to present the normalized version, since it will encourage readers to implement the method properly (i.e., normalizing after each step to avoid underflow).

In practice it is more numerically stable to compute the log probabilities  $\ell_t(j) = \log p(y_t | z_t = j)$  of the evidence, rather than the probabilities  $\lambda_t(j) = p(y_t | z_t = j)$ . We can combine the state conditional log likelihoods  $\lambda_t(j)$  with the state prior  $p(z_t = j | y_{1:t-1})$  by using the log-sum-exp trick, as in Equation (28.30).

### 8.2.6 Time and space complexity

It is clear that a straightforward implementation of the forwards-backwards algorithm takes  $O(K^2T)$  time, since we must perform a  $K \times K$  matrix multiplication at each step. For some applications, such as speech recognition,  $K$  is very large, so the  $O(K^2)$  term becomes prohibitive. Fortunately, if the transition matrix is sparse, we can reduce this substantially. For example, in a sparse left-to-right transition matrix (e.g., Figure 8.7(a)), the algorithm takes  $O(TK)$  time.

In some cases, we can exploit special properties of the state space, even if the transition matrix is not sparse. In particular, suppose the states represent a discretization of an underlying continuous

<sup>2</sup> For example, if the observations are independent of the states, we have  $p(z_t = j, y_{1:t}) = p(z_t = j) \prod_{i=1}^t p(y_i)$ , which becomes exponentially small with  $t$ .

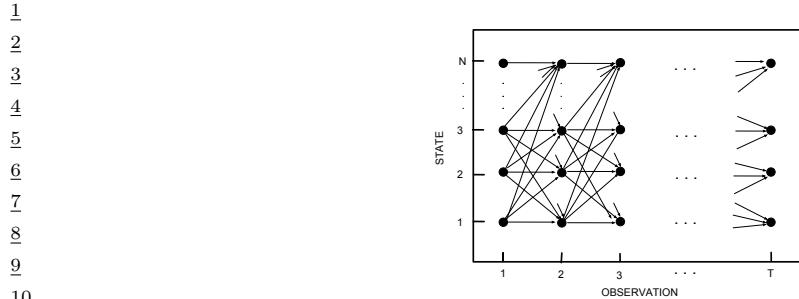


Figure 8.6: The trellis of states vs time for a Markov chain. Adapted from [Rab89].

state-space, and the transition matrix has the form  $A_{i,j} \propto \rho(\mathbf{z}_j - \mathbf{z}_i)$ , where  $\mathbf{z}_i$  is the continuous vector represented by state  $i$  and  $\rho(u)$  is some scalar cost function, such as Euclidean distance. Then one can implement the forwards-backwards algorithm in  $O(TK \log K)$  time. The key is to rewrite Equation (8.10) as a convolution,

$$\alpha_{t|t-1}(j) = p(z_t = j | \mathbf{y}_{1:t-1}) = \sum_i \alpha_{t-1}(i) A_{i,j} = \sum_i \alpha_{t-1}(i) \rho(j - i) \quad (8.48)$$

and then to apply the Fast Fourier Transform. (A similar transformation can be applied in the backwards pass.) This is very useful for models with large state spaces. See [FKH03] for details.

We can also reduce inference to  $O(\log T)$  time by using a **parallel prefix scan** operator that can be run efficiently on GPUs. For details, see [HSGF21].

In some cases, the bottleneck is memory, not time. In particular, to compute the posteriors  $\gamma_t$ , we must store the filtered distributions  $\alpha_t$  for  $t = 1, \dots, T$  until we do the backwards pass. It is possible to devise a simple divide-and-conquer algorithm that reduces the space complexity from  $O(KT)$  to  $O(K \log T)$  at the cost of increasing the running time from  $O(K^2T)$  to  $O(K^2T \log T)$ . The basic idea is to store  $\alpha_t$  and  $\beta_t$  vectors at a logarithmic number of intermediate checkpoints, and then recompute the missing messages on demand from these checkpoints. See [BMR97; ZP00] for details.

### 8.2.7 The Viterbi algorithm

The MAP estimate is (one of) the sequences with maximum posterior probability:

$$\mathbf{z}_{1:T}^* = \operatorname{argmax}_{\mathbf{z}_{1:T}} p(\mathbf{z}_{1:T} | \mathbf{y}_{1:T}) = \operatorname{argmax}_{\mathbf{z}_{1:T}} \log p(\mathbf{z}_{1:T} | \mathbf{y}_{1:T}) \quad (8.49)$$

$$= \operatorname{argmax}_{\mathbf{z}_{1:T}} \log \pi_1(\mathbf{z}_1) + \log \lambda_1(\mathbf{z}_1) + \sum_{t=2}^T [\log A(\mathbf{z}_{t-1}, \mathbf{z}_t) + \log \lambda_t(\mathbf{z}_t)] \quad (8.50)$$

This is equivalent to computing a shortest path through the **trellis diagram** in Figure 8.6, where the nodes are possible states at each time step, and the node and edge weights are log probabilities. This can be computed in  $O(TK^2)$  time using the **Viterbi algorithm** [Vit67], as we explain below.

1    **8.2.7.1 Forwards pass**

2    Recall the (unnormalized) forwards equation

3

$$\alpha'_t(j) = p(\mathbf{z}_t = j, \mathbf{y}_{1:t}) = \sum_{\mathbf{z}_{1:t-1}} p(\mathbf{z}_{1:t-1}, \mathbf{z}_t = j, \mathbf{y}_{1:t}) \quad (8.51)$$

4

5    Now suppose we replace sum with max to get

6

$$\delta_t(j) \triangleq \max_{\mathbf{z}_{1:t-1}} p(\mathbf{z}_{1:t-1}, \mathbf{z}_t = j, \mathbf{y}_{1:t}) \quad (8.52)$$

7

8    This is the maximum probability we can assign to the data so far if we end up in state  $j$ . The key  
9    insight is that the most probable path to state  $j$  at time  $t$  must consist of the most probable path to  
10   some other state  $i$  at time  $t-1$ , followed by a transition from  $i$  to  $j$ . Hence

11

$$\delta_t(j) = \lambda_t(j) \left[ \max_i \delta_{t-1}(i) A_{i,j} \right] \quad (8.53)$$

12

13   We initialize by setting  $\delta_1(j) = \pi_j \lambda_1(j)$ .

14   We often work in the log domain to avoid numerical issues. Let  $\delta'_t(j) = -\log \delta_t(j)$ ,  $\lambda_t(j) = -\log p(\mathbf{y}_t | \mathbf{z}_t = j)$ ,  $A'(i, j) = -\log p(\mathbf{z}_t = j | \mathbf{z}_{t-1} = i)$ . Then we have

15

$$\delta'_t(j) = \lambda'_t(j) + \left[ \min_i \delta'_{t-1}(i) + A'(i, j) \right] \quad (8.54)$$

16

17   We also need to keep track of the most likely previous (**ancestor**) state, for each possible state  
18   that we end up in:

19

$$a_t(j) = \operatorname{argmax}_i \delta_{t-1}(i) A_{i,j} = \operatorname{argmin}_i \delta'_{t-1}(i) + A'(i, j) \quad (8.55)$$

20

21   That is,  $a_t(j)$  stores the identity of the previous state on the most probable path to  $\mathbf{z}_t = j$ . We will  
22   see why we need this in Section 8.2.7.2.

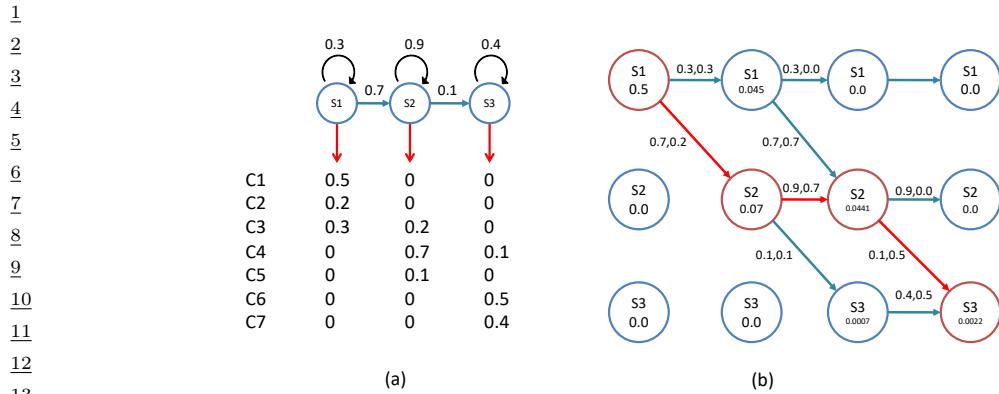
23   **8.2.7.2 Backwards pass**

24   In the backwards pass, we compute the most probable sequence of states using a **traceback** procedure,  
25   as follows:  $\mathbf{z}_t^* = a_{t+1}(\mathbf{z}_{t+1}^*)$ , where we initialize using  $\mathbf{z}_T^* = \operatorname{argmax}_i \delta_T(i)$ . This is just following the  
26   chain of ancestors along the MAP path.

27   If there is a unique MAP estimate, the above procedure will give the same result as picking  
28    $\hat{z}_t = \operatorname{argmax}_j \gamma_t(j)$ , computed by forwards-backwards, as shown in [WF01b]. However, if there are  
29   multiple posterior modes, the latter approach may not find any of them, since it chooses each state  
30   independently, and hence may break ties in a manner that is inconsistent with its neighbors. The  
31   traceback procedure avoids this problem, since once  $\mathbf{z}_t$  picks its most probable state, the previous  
32   nodes condition on this event, and therefore they will break ties consistently.

33   **8.2.7.3 Example**

34   In Figure 8.4(c), we show the Viterbi trace for the casino HMM. We see that, most of the time, the  
35   estimated state corresponds to the true state.



14 *Figure 8.7: Illustration of Viterbi decoding in a simple HMM for speech recognition. (a) A 3-state HMM  
15 for a single phone. We are visualizing the state transition diagram. We assume the observations have been  
16 vector quantized into 7 possible symbols,  $C_1, \dots, C_7$ . Each state  $s_1, s_2, s_3$  has a different distribution over  
17 these symbols. Adapted from Figure 15.20 of [RN02]. (b) Illustration of the Viterbi algorithm applied to this  
18 model, with data sequence  $C_1, C_3, C_4, C_6$ . The columns represent time, and the rows represent states. The  
19 numbers inside the circles represent the  $\delta_t(j)$  value for that state. An arrow from state  $i$  at  $t - 1$  to state  
20  $j$  at  $t$  is annotated with two numbers: the first is the probability of the  $i \rightarrow j$  transition, and the second is  
21 the probability of generating observation  $y_t$  from state  $j$ . The red lines/ circles represent the most probable  
22 sequence of states. Adapted from Figure 24.27 of [RN95].*

23  
24 In Figure 8.7, we give a detailed worked example of the Viterbi algorithm, based on [Rus+95].  
25 Suppose we observe the sequence of discrete observations  $y_{1:4} = (C_1, C_3, C_4, C_6)$ , representing  
26 codebook entries in a vector-quantized version of a speech signal. The model starts in state  $z_1 = S_1$ .  
27 The probability of generating  $x_1 = C_1$  in  $z_1$  is 0.5, so we have  $\delta_1(1) = 0.5$ , and  $\delta_1(i) = 0$  for all other  
28 states. Next we can self-transition to  $S_1$  with probability 0.3, or transition to  $S_2$  with probability 0.7.  
29 If we end up in  $S_1$ , the probability of generating  $x_2 = C_3$  is 0.3; if we end up in  $S_2$ , the probability  
30 of generating  $x_2 = C_3$  is 0.2. Hence we have

$$\delta_2(1) = \delta_1(1)A(1,1)\lambda_2(1) = 0.5 \cdot 0.3 \cdot 0.3 = 0.045 \quad (8.56)$$

$$\delta_2(2) = \delta_1(1)A(1,2)\lambda_2(2) = 0.5 \cdot 0.7 \cdot 0.2 = 0.07 \quad (8.57)$$

31 Thus state 2 is more probable at  $t = 2$ ; see the second column of Figure 8.7(b). The algorithm  
32 continues in this way until we have reached the end of the sequence. Once we have reached the end,  
33 we can follow the red arrows back to recover the MAP path (which is 1,2,2,3).

34 For more details on HMMs for automatic speech recognition (ASR) see e.g., [JM08].

#### 40 8.2.7.4 Time and space complexity

41 The time complexity of Viterbi is clearly  $O(K^2T)$  in general, and the space complexity is  $O(KT)$ ,  
42 both the same as forwards-backwards. If the transition matrix has the form  $A_{i,j} \propto \rho(z_j - z_i)$ , where  
43  $z_i$  is the continuous vector represented by state  $i$  and  $\rho(u)$  is some scalar cost function, such as  
44 Euclidean distance, we can implement Viterbi in  $O(TK)$  time, by using the generalized distance  
45 transform to implement Equation (8.54). See [FHK03; FH12] for details.

46  
47

1

### 8.2.7.5 N-best list

2  
3 There are often multiple paths which have the same likelihood. The Viterbi algorithm returns one of  
4 them, but can be extended to return the top  $N$  paths [SC90; NG01]. This is called the **N-best list**.  
5 Computing such a list can provide a better summary of the posterior uncertainty.

6 In addition, we can perform **discriminative reranking** [CK05] of all the sequences in  $\mathcal{L}_N$ , based  
7 on global features derived from  $(\mathbf{y}_{1:T}, \mathbf{z}_{1:T})$ . This technique is widely used in speech recognition. For  
8 example, consider the sentence “recognize speech”. It is possible that the most probable interpretation  
9 by the system of this acoustic signal is “wreck a nice speech”, or maybe “wreck a nice beach” (see  
10 Figure 34.1). Maybe the correct interpretation is much lower down on the list. However, by using  
11 a re-ranking system, we may be able to improve the score of the correct interpretation based on a  
12 more global context.

13 One problem with the  $N$ -best list is that often the top  $N$  paths are very similar to each other,  
14 rather than representing qualitatively different interpretations of the data. Instead we might want to  
15 generate a more diverse set of paths to more accurately represent posterior uncertainty. One way to  
16 do this is to sample paths from the posterior, as we discuss in Section 8.2.8. Another way is to use a  
17 determinantal point process (Section 31.9.5) which encourages points to be diverse [Bat+12; ZA12].  
18

19

### 8.2.8 Forwards filtering, backwards sampling

20  
21 Rather than computing the single most probable path, it is often useful to sample multiple paths from  
22 the posterior:  $\mathbf{z}_{1:T}^s \sim p(\mathbf{z}_{1:T} | \mathbf{y}_{1:T})$ . We can do this by modifying the forwards filtering backwards  
23 smoothing algorithm from Section 8.2.3, so that we draw samples on the backwards pass, rather than  
24 computing marginals. This is called **forwards filtering backwards sampling** (also sometimes  
25 unfortunately abbreviated to FFBS). In particular, note that we can write the joint from right to left  
26 using

$$p(\mathbf{z}_{1:T} | \mathbf{y}_{1:T}) = p(\mathbf{z}_T | \mathbf{y}_{1:T}) p(\mathbf{z}_{T-1} | \mathbf{z}_T, \mathbf{y}_{1:T}) p(\mathbf{z}_{T-2} | \mathbf{z}_{T-1}, \mathbf{z}_{T}, \mathbf{y}_{1:T}) \cdots p(\mathbf{z}_1 | \mathbf{z}_2, \mathbf{z}_{3:T}, \mathbf{y}_{1:T}) \quad (8.58)$$

$$= p(\mathbf{z}_T | \mathbf{y}_{1:T}) \prod_{t=T-1}^1 p(\mathbf{z}_t | \mathbf{z}_{t+1}, \mathbf{y}_{1:T}) \quad (8.59)$$

32 Thus at step  $t$  we sample  $\mathbf{z}_t^s$  from  $p(\mathbf{z}_t | \mathbf{z}_{t+1}^s, \mathbf{y}_{1:T})$  given in Equation (8.17).

35

## 8.3 Inference based on the Kalman filter

37 In this section, we discuss inference in SSMs where all the distributions are linear Gaussian. This is  
38 called a **linear Gaussian state space model (LG-SSM)** or a **linear dynamical system (LDS)**.  
39 We discuss such models in detail in Section 29.6, but in brief they have the following form:

$$p(\mathbf{z}_t | \mathbf{z}_{t-1}, \mathbf{u}_t) = \mathcal{N}(\mathbf{z}_t | \mathbf{F}_t \mathbf{z}_{t-1} + \mathbf{B}_t \mathbf{u}_t + \mathbf{b}_t, \mathbf{Q}_t) \quad (8.60)$$

$$p(\mathbf{y}_t | \mathbf{z}_t, \mathbf{u}_t) = \mathcal{N}(\mathbf{y}_t | \mathbf{H}_t \mathbf{z}_t + \mathbf{D}_t \mathbf{u}_t + \mathbf{d}_t, \mathbf{R}_t) \quad (8.61)$$

44 where  $\mathbf{z}_t \in \mathbb{R}^{N_z}$  is the hidden state,  $\mathbf{y}_t \in \mathbb{R}^{N_y}$  is the observation, and  $\mathbf{u}_t \in \mathbb{R}^{N_u}$  is the input. (We  
45 have allowed the parameters to be time-varying, for later extensions that we will consider.) We often  
46 assume the means of the process noise and observation noise (i.e., the bias or offset terms) are zero,  
47

<sup>1</sup> so  $\mathbf{b}_t = \mathbf{0}$  and  $\mathbf{d}_t = \mathbf{0}$ . In addition, we often have no inputs, so  $\mathbf{B}_t = \mathbf{D}_t = \mathbf{0}$ . In this case, the model  
<sup>2</sup> simplifies to the following.<sup>3</sup>  
<sup>3</sup>

$$p(\mathbf{z}_t | \mathbf{z}_{t-1}) = \mathcal{N}(\mathbf{z}_t | \mathbf{F}_t \mathbf{z}_{t-1}, \mathbf{Q}_t) \quad (8.62)$$

$$p(\mathbf{y}_t | \mathbf{z}_t) = \mathcal{N}(\mathbf{y}_t | \mathbf{H}_t \mathbf{z}_t, \mathbf{R}_t) \quad (8.63)$$

<sup>7</sup> See Figure 8.1 for the graphical model.

<sup>8</sup> Note that an LG-SSM is just a special case of a Gaussian Bayes net (Section 4.2.3), so the  
<sup>9</sup> entire joint distribution  $p(\mathbf{y}_{1:T}, \mathbf{z}_{1:T} | \mathbf{u}_{1:T})$  is a large multivariate Gaussian with  $N_y N_z T$  dimensions.  
<sup>10</sup> However, it has a special structure that makes it computationally tractable to use, as we show below.  
<sup>11</sup> In particular, we will discuss the **Kalman filter** and **Kalman smoother**, that can perform exact  
<sup>12</sup> filtering and smoothing in  $O(T N_z^3)$  time.  
<sup>13</sup>

### <sup>14</sup> 8.3.1 Examples

<sup>16</sup> Before diving into the theory, we give some motivating examples.  
<sup>17</sup>

#### <sup>18</sup> 8.3.1.1 Tracking and state estimation

<sup>20</sup> A common application of LG-SSMs is for **tracking** objects, such as airplanes or animals, from noisy  
<sup>21</sup> measurements, such as radar or cameras. For example, suppose we want to track an object moving  
<sup>22</sup> in 2d. (We discuss this example in more detail in Section 29.7.1.) The hidden state  $\mathbf{z}_t$  encodes the  
<sup>23</sup> location,  $(x_{t1}, x_{t2})$ , and the velocity,  $(\dot{x}_{t1}, \dot{x}_{t2})$ , of the moving object. The observation  $\mathbf{y}_t$  is a noisy  
<sup>24</sup> version of the location. (The velocity is not observed but can be inferred from the change in location.)  
<sup>25</sup> We assume that we obtain measurements with a sampling period of  $\Delta$ . The new location is the old  
<sup>26</sup> location plus  $\Delta$  times the velocity, plus noise added to all terms:

$$\mathbf{z}_t = \underbrace{\begin{pmatrix} 1 & 0 & \Delta & 0 \\ 0 & 1 & 0 & \Delta \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}}_{\mathbf{F}} \mathbf{z}_{t-1} + \mathbf{q}_t \quad (8.64)$$

<sup>33</sup> where  $\mathbf{q}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_t)$ . The observation extracts the location and adds noise:  
<sup>34</sup>

$$\mathbf{y}_t = \underbrace{\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix}}_{\mathbf{H}} \mathbf{z}_t + \mathbf{r}_t \quad (8.65)$$

<sup>39</sup> where  $\mathbf{r}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{R}_t)$ .

<sup>40</sup> Our goal is to use this model to estimate the unknown location (and velocity) of the object given  
<sup>41</sup> the noisy observations. In particular, in the filtering problem, we want to compute  $p(\mathbf{z}_t | \mathbf{y}_{1:t})$  in  
<sup>42</sup> a recursive fashion. Figure 8.8(b) illustrates filtering for the linear Gaussian SSM applied to the  
<sup>43</sup> noisy tracking data in Figure 8.8(a) (shown by the green dots). The filtered estimates are computed  
<sup>44</sup>

<sup>45</sup> 3. Our notation is similar to [Sar13], except he writes  $p(\mathbf{x}_k | \mathbf{x}_{k-1}) = \mathcal{N}(\mathbf{x}_k | \mathbf{A}_{k-1} \mathbf{x}_{k-1}, \mathbf{Q}_{k-1})$  instead of  $p(\mathbf{z}_t | \mathbf{z}_{t-1}) =$   
<sup>46</sup>  $\mathcal{N}(\mathbf{z}_t | \mathbf{F}_t \mathbf{z}_{t-1}, \mathbf{Q}_t)$ , and  $p(\mathbf{y}_k | \mathbf{x}_k) = \mathcal{N}(\mathbf{y}_k | \mathbf{H}_k \mathbf{x}_k, \mathbf{R}_k)$  instead of  $p(\mathbf{y}_t | \mathbf{z}_t) = \mathcal{N}(\mathbf{y}_t | \mathbf{H}_t \mathbf{z}_t, \mathbf{R}_t)$ .

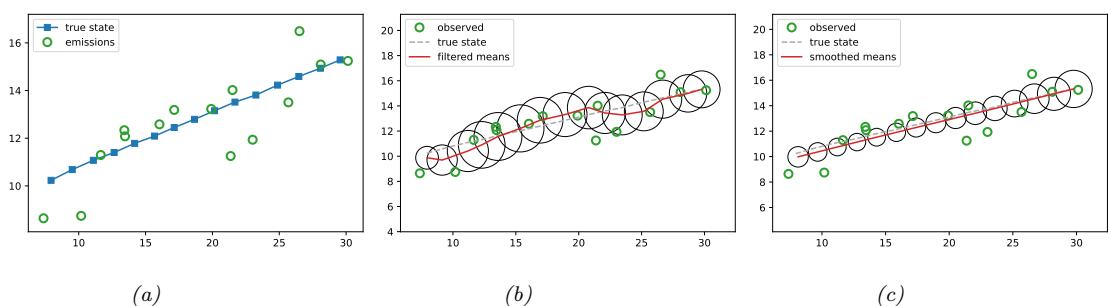


Figure 8.8: Illustration of Kalman filtering and smoothing for a linear dynamical system. (a) Observations (green circles) are generated by an object moving to the right (true location denoted by blue squares). (b) Results of online Kalman filtering. Red cross is the posterior mean, circles are 95% confidence ellipses derived from the posterior covariance. (c) Same as (b), but using offline Kalman smoothing. The MSE in the trajectory for filtering is 3.13, and for smoothing is 1.71. Generated by [kf\\_tracking.ipynb](#).

using the Kalman filter algorithm described in Section 8.3.2. The red line shows the posterior mean estimate of the location, and the black circles show the posterior covariance. We see that the estimated trajectory is less noisy than the raw data, since it incorporates prior knowledge about how the data was generated.

Another task of interest is the smoothing problem where we want to compute  $p(\mathbf{z}_t | \mathbf{y}_{1:T})$  using an offline dataset. Figure 8.8(c) illustrates smoothing for the LG-SSM, implemented using the Kalman smoothing algorithm described in Section 8.3.3. We see that the resulting estimate is smoother, and that the posterior uncertainty is reduced (as visualized by the smaller confidence ellipses).

The disadvantage of the above smoothing method is that we have to wait until all the data has been observed before we start performing inference. **Fixed lag smoothing** is a useful compromise between online and offline estimation; it involves computing  $p(\mathbf{z}_{t-\ell} | \mathbf{y}_{1:t})$ , where  $\ell > 0$  is called the lag. This gives better performance than filtering, but incurs a slight delay. By changing the size of the lag, we can trade off accuracy vs delay.

### 8.3.1.2 Online Bayesian linear regression (recursive least squares)

In Section 29.7.2 we discuss how to use the Kalman filter to recursively compute the exact posterior  $p(\mathbf{w} | \mathcal{D}_{1:t})$  for a linear regression model in an online fashion. This is known as the recursive least squares algorithm. The basic idea is to treat the latent state to be the parameter values,  $\mathbf{z}_t = \mathbf{w}$ , and to define the non-stationary observation model as  $p(\mathbf{y}_t | \mathbf{z}_t) = \mathcal{N}(y_t | \mathbf{x}_t^\top \mathbf{z}_t, \sigma^2)$ , and the dynamics model as  $p(\mathbf{z}_t | \mathbf{z}_{t-1}) = \mathcal{N}(\mathbf{z}_t | \mathbf{z}_{t-1}, \mathbf{0I})$ .

### 8.3.1.3 Time series forecasting

In Section 29.12, we discuss how to use Kalman filtering to perform time series forecasting.

1 **8.3.2 The Kalman filter**

3 The **Kalman Filter (KF)** is an algorithm for exact Bayesian filtering for linear Gaussian state space  
4 models. The resulting algorithm is the Gaussian analog of the HMM filter in Section 8.2.2. The belief  
5 state at time  $t$  is now given by  $p(\mathbf{z}_t | \mathbf{y}_{1:t}) = \mathcal{N}(\mathbf{z}_t | \boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t})$ , where we use the notation  $\boldsymbol{\mu}_{t|t'}$  and  $\boldsymbol{\Sigma}_{t|t'}$   
6 to represent the posterior mean and covariance given  $\mathbf{y}_{1:t'}$ .<sup>4</sup> Since everything is Gaussian, we can  
7 perform the prediction and update steps in closed form, as we explain below (see Section 8.3.2.4 for  
8 the derivation).

10 **8.3.2.1 Predict step**

12 The one-step-ahead prediction for the hidden state, also called the **time update step**, is given by  
13 the following:

$$\underline{15} \quad p(\mathbf{z}_t | \mathbf{y}_{1:t-1}, \mathbf{u}_{1:t}) = \mathcal{N}(\mathbf{z}_t | \boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}) \quad (8.66)$$

$$\underline{17} \quad \boldsymbol{\mu}_{t|t-1} = \mathbf{F}_t \boldsymbol{\mu}_{t-1|t-1} + \mathbf{B}_t \mathbf{u}_t + \mathbf{b}_t \quad (8.67)$$

$$\underline{18} \quad \boldsymbol{\Sigma}_{t|t-1} = \mathbf{F}_t \boldsymbol{\Sigma}_{t-1|t-1} \mathbf{F}_t^\top + \mathbf{Q}_t \quad (8.68)$$

20 **8.3.2.2 Update step**

22 The update step (also called the **measurement step**) can be computed using Bayes rule, as follows:

$$\underline{24} \quad p(\mathbf{z}_t | \mathbf{y}_{1:t}, \mathbf{u}_{1:t}) = \mathcal{N}(\mathbf{z}_t | \boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t}) \quad (8.69)$$

$$\underline{25} \quad \mathbf{m}_t = \mathbf{H}_t \boldsymbol{\mu}_{t|t-1} + \mathbf{D}_t \mathbf{u}_t + \mathbf{d}_t \quad (8.70)$$

$$\underline{27} \quad \mathbf{S}_t = \mathbf{H}_t \boldsymbol{\Sigma}_{t|t-1} \mathbf{H}_t^\top + \mathbf{R}_t \quad (8.71)$$

$$\underline{28} \quad \mathbf{K}_t = \boldsymbol{\Sigma}_{t|t-1} \mathbf{H}_t^\top \mathbf{S}_t^{-1} \quad (8.72)$$

$$\underline{29} \quad \mathbf{e}_t = \mathbf{y}_t - \mathbf{m}_t \quad (8.73)$$

$$\underline{31} \quad \boldsymbol{\mu}_{t|t} = \boldsymbol{\mu}_{t|t-1} + \mathbf{K}_t \mathbf{e}_t \quad (8.74)$$

$$\underline{32} \quad \boldsymbol{\Sigma}_{t|t} = \boldsymbol{\Sigma}_{t|t-1} - \mathbf{K}_t \mathbf{H}_t \boldsymbol{\Sigma}_{t|t-1} \quad (8.75)$$

$$\underline{33} \quad = \boldsymbol{\Sigma}_{t|t-1} - \mathbf{K}_t \mathbf{S}_t \mathbf{K}_t^\top \quad (8.76)$$

35 where  $\mathbf{m}_t$  is the expected observation,  $\mathbf{e}_t$  is the **residual error** or **innovation term**, and  $\mathbf{K}_t$  is the  
36 **Kalman gain matrix**.

37 Note that, by using the matrix inversion lemma, the Kalman gain matrix can also be written as

$$\underline{39} \quad \mathbf{K}_t = \boldsymbol{\Sigma}_{t|t-1} \mathbf{H}_t^\top (\mathbf{H}_t \boldsymbol{\Sigma}_{t|t-1} \mathbf{H}_t^\top + \mathbf{R}_t)^{-1} = (\boldsymbol{\Sigma}_{t|t-1}^{-1} + \mathbf{H}_t^\top \mathbf{R}_t^{-1} \mathbf{H}_t)^{-1} \mathbf{H}_t^\top \mathbf{R}_t^{-1} \quad (8.77)$$

41 This is useful if  $\mathbf{R}_t^{-1}$  is precomputed (e.g., if it is constant over time) and  $N_y \gg N_z$ .

43 4. We represent the mean and covariance of the filtered belief state by  $\boldsymbol{\mu}_{t|t}$  and  $\boldsymbol{\Sigma}_{t|t}$ , but some authors use the notation  
44  $\mathbf{m}_t$  and  $\mathbf{P}_t$  instead. We represent the mean and covariance of the smoothed belief state by  $\boldsymbol{\mu}_{t|T}$  and  $\boldsymbol{\Sigma}_{t|T}$ , but some  
45 authors use the notation  $\mathbf{m}_t^s$  and  $\mathbf{P}_t^s$  instead. Finally, we represent the mean and covariance of the one-step-ahead  
46 posterior predictive distribution,  $p(\mathbf{z}_t | \mathbf{y}_{1:t-1})$ , by  $\boldsymbol{\mu}_{t|t-1}$  and  $\boldsymbol{\Sigma}_{t|t-1}$ , whereas some authors use  $\boldsymbol{\mu}_t^-$  and  $\boldsymbol{\Sigma}_t^-$  instead.

---

### 8.3.2.3 Posterior predictive

The one-step-ahead posterior predictive density for the observations can be computed as follows. (We ignore inputs and bias terms, for notational brevity.) First we compute the one-step-ahead predictive density for latent states:

$$p(\mathbf{z}_t | \mathbf{y}_{1:t-1}) = \int p(\mathbf{z}_t | \mathbf{z}_{t-1}) p(\mathbf{z}_{t-1} | \mathbf{y}_{1:t-1}) d\mathbf{z}_{t-1} = \mathcal{N}(\mathbf{z}_t | \boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}) \quad (8.78)$$

$$= \mathcal{N}(\mathbf{z}_t | \mathbf{F}_t \boldsymbol{\mu}_{t-1|t-1}, \mathbf{F}_t \boldsymbol{\Sigma}_{t-1|t-1} \mathbf{F}_t^\top + \mathbf{Q}_t) = \mathcal{N}(\mathbf{z}_t | \boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}) \quad (8.79)$$

Then we convert this to a prediction about observations by marginalizing out  $\mathbf{z}_t$ :

$$p(\mathbf{y}_t | \mathbf{y}_{1:t-1}) = \int p(\mathbf{y}_t, \mathbf{z}_t | \mathbf{y}_{1:t-1}) d\mathbf{z}_t = \int p(\mathbf{y}_t | \mathbf{z}_t) p(\mathbf{z}_t | \mathbf{y}_{1:t-1}) d\mathbf{z}_t = \mathcal{N}(\mathbf{y}_t | \mathbf{m}_t, \mathbf{S}_t) \quad (8.80)$$

This can also be used to compute the log-likelihood of the observations: The normalization constant of the new posterior can be computed as follows:

$$\log p(\mathbf{y}_{1:T}) = \sum_{t=1}^T \log p(\mathbf{y}_t | \mathbf{y}_{1:t-1}) = \sum_{t=1}^T \log Z_t \quad (8.81)$$

where we define  $p(\mathbf{y}_1 | \mathbf{y}_0) = p(\mathbf{y}_1)$ . This is just a sum of the log probabilities of the one-step-ahead measurement predictions, and is a measure of how “surprised” the model is at each step.

We can generalize the prediction step to predict observations  $K$  steps into the future by first forecasting  $K$  steps in latent space, and then “grounding” the final state into predicted observations. (This is in contrast to an RNN (Section 16.3.4), which requires generating observations at each step, in order to update future hidden states.)

### 8.3.2.4 Derivation

In this section we derive the Kalman filter equations, following [Sar13, p57]. The results are a straightforward application of the rules for manipulating linear Gaussian systems, discussed in Section 2.2.6.

First we derive the prediction step. From Equation (2.81), the joint predictive distribution for states is given by

$$p(\mathbf{z}_{t-1}, \mathbf{z}_t | \mathbf{y}_{1:t-1}) = p(\mathbf{z}_t | \mathbf{z}_{t-1}) p(\mathbf{z}_{t-1} | \mathbf{y}_{1:t-1}) \quad (8.82)$$

$$= \mathcal{N}(\mathbf{z}_t | \mathbf{F}_t \mathbf{z}_{t-1}, \mathbf{Q}_t) \mathcal{N}(\mathbf{z}_{t-1} | \boldsymbol{\mu}_{t-1|t-1}, \boldsymbol{\Sigma}_{t-1|t-1}) \quad (8.83)$$

$$= \mathcal{N}\left(\begin{pmatrix} \mathbf{z}_{t-1} \\ \mathbf{z}_t \end{pmatrix} | \boldsymbol{\mu}', \boldsymbol{\Sigma}'\right) \quad (8.84)$$

where

$$\boldsymbol{\mu}' = \begin{pmatrix} \boldsymbol{\mu}_{t-1|t-1} \\ \mathbf{F}_t \boldsymbol{\mu}_{t-1|t-1} \end{pmatrix}, \quad \boldsymbol{\Sigma}' = \begin{pmatrix} \boldsymbol{\Sigma}_{t-1|t-1} & \boldsymbol{\Sigma}_{t-1|t-1} \mathbf{F}_t^\top \\ \mathbf{F}_t \boldsymbol{\Sigma}_{t-1|t-1} & \mathbf{F}_t \boldsymbol{\Sigma}_{t-1|t-1} \mathbf{F}_t^\top + \mathbf{Q}_t \end{pmatrix} \quad (8.85)$$

Hence the marginal predictive distribution for states is given by

$$p(\mathbf{z}_t | \mathbf{y}_{1:t-1}) = \mathcal{N}(\mathbf{z}_t | \mathbf{F}_t \boldsymbol{\mu}_{t-1|t-1}, \mathbf{F}_t \boldsymbol{\Sigma}_{t-1|t-1} \mathbf{F}_t^\top + \mathbf{Q}_t) = \mathcal{N}(\mathbf{z}_t | \boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}) \quad (8.86)$$

Now we derive the measurement update step. The joint distribution for state and observation is given by

$$p(\mathbf{z}_t, \mathbf{y}_t | \mathbf{y}_{1:t-1}) = p(\mathbf{y}_t | \mathbf{z}_t) p(\mathbf{z}_t | \mathbf{y}_{1:t-1}) \quad (8.87)$$

$$= \mathcal{N}(\mathbf{y}_t | \mathbf{H}_t \mathbf{z}_t, \mathbf{R}_t) \mathcal{N}(\mathbf{z}_t | \boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}) \quad (8.88)$$

$$= \mathcal{N}\left(\begin{pmatrix} \mathbf{z}_t \\ \mathbf{y}_t \end{pmatrix} | \boldsymbol{\mu}'', \boldsymbol{\Sigma}''\right) \quad (8.89)$$

where

$$\boldsymbol{\mu}'' = \begin{pmatrix} \boldsymbol{\mu}_{t|t-1} \\ \mathbf{H}_t \boldsymbol{\mu}_{t|t-1} \end{pmatrix}, \quad \boldsymbol{\Sigma}'' = \begin{pmatrix} \boldsymbol{\Sigma}_{t|t-1} & \boldsymbol{\Sigma}_{t|t-1} \mathbf{H}_t^\top \\ \mathbf{H}_t \boldsymbol{\Sigma}_{t|t-1} & \mathbf{H}_t \boldsymbol{\Sigma}_{t|t-1}^{-1} \mathbf{H}_t^\top + \mathbf{R}_t \end{pmatrix} \quad (8.90)$$

Finally, we convert this joint into a conditional using Equation (2.39) as follows:

$$p(\mathbf{z}_t | \mathbf{y}_t, \mathbf{y}_{1:t-1}) = \mathcal{N}(\mathbf{z}_t | \boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t}) \quad (8.91)$$

$$\boldsymbol{\mu}_{t|t} = \boldsymbol{\mu}_{t|t-1} + \boldsymbol{\Sigma}_{t|t-1} \mathbf{H}_t^\top (\mathbf{H}_t \boldsymbol{\Sigma}_{t|t-1} \mathbf{H}_t^\top + \mathbf{R}_t)^{-1} [\mathbf{y}_t - \mathbf{H}_t \boldsymbol{\mu}_{t|t-1}] \quad (8.92)$$

$$= \boldsymbol{\mu}_{t|t-1} + \mathbf{K}_t [\mathbf{y}_t - \mathbf{H}_t \boldsymbol{\mu}_{t|t-1}] \quad (8.93)$$

$$\boldsymbol{\Sigma}_{t|t} = \boldsymbol{\Sigma}_{t|t-1} - \boldsymbol{\Sigma}_{t|t-1} \mathbf{H}_t^\top (\mathbf{H}_t \boldsymbol{\Sigma}_{t|t-1} \mathbf{H}_t^\top + \mathbf{R}_t)^{-1} \mathbf{H}_t \boldsymbol{\Sigma}_{t|t-1} \quad (8.94)$$

$$= \boldsymbol{\Sigma}_{t|t-1} - \mathbf{K}_t \mathbf{H}_t \boldsymbol{\Sigma}_{t|t-1} \quad (8.95)$$

where

$$\mathbf{S}_t = \mathbf{H}_t \boldsymbol{\Sigma}_{t|t-1} \mathbf{H}_t^\top + \mathbf{R}_t \quad (8.96)$$

$$\mathbf{K}_t = \boldsymbol{\Sigma}_{t|t-1} \mathbf{H}_t^\top \mathbf{S}_t^{-1} \quad (8.97)$$

### 8.3.2.5 Abstract formulation

We can represent the Kalman filter equations much more compactly by defining various functions that create and manipulate jointly Gaussian systems, as in Section 2.2.6. In particular, suppose  $p(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}})$ , and  $p(\mathbf{y} | \mathbf{z}) = \mathcal{N}(\mathbf{y} | \mathbf{A}\mathbf{z} + \mathbf{b}, \boldsymbol{\Omega})$ . Then the joint is given by  $p(\mathbf{z}, \mathbf{y}) = \mathcal{N}(\mathbf{z}, \mathbf{y} | \tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}})$ , where  $\tilde{\boldsymbol{\mu}} = (\tilde{\boldsymbol{\mu}}, \mathbf{m})$  and  $\tilde{\boldsymbol{\Sigma}} = (\tilde{\boldsymbol{\Sigma}}, \mathbf{C}; \mathbf{C}^\top, \mathbf{S})$ , where the terms  $\mathbf{m}$ ,  $\mathbf{S}$  and  $\mathbf{C}$  are given by **LinGaussPredict** in Algorithm 7. From this, the posterior is given by  $p(\mathbf{z} | \mathbf{y}) = \mathcal{N}(\mathbf{z} | \hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}})$ , where the posterior parameters are given by **GaussCondition** in Algorithm 7.

We can now apply these functions to derive Kalman filtering as follows. In the prediction step, we compute

$$p(\mathbf{z}_{t-1}, \mathbf{z}_t | \mathbf{y}_{1:t-1}) = \mathcal{N}\left(\begin{pmatrix} \boldsymbol{\mu}_{t-1|t-1} \\ \boldsymbol{\mu}_{t|t-1} \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma}_{t-1|t-1} & \boldsymbol{\Sigma}_{t-1,t|t-1} \\ \boldsymbol{\Sigma}_{t,t-1|t-1} & \boldsymbol{\Sigma}_{t|t-1} \end{pmatrix}\right) \quad (8.98)$$

$$(\boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}, \boldsymbol{\Sigma}_{t-1,t|t}) = \text{LinGaussPredict}(\boldsymbol{\mu}_{t-1|t-1}, \boldsymbol{\Sigma}_{t-1|t-1}, \mathbf{F}_t, \mathbf{B}_t \mathbf{u}_t + \mathbf{b}_t, \mathbf{Q}_t) \quad (8.99)$$

from which we get the marginal distribution

$$p(\mathbf{z}_t | \mathbf{y}_{1:t-1}) = \mathcal{N}(\boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}) \quad (8.100)$$

---

1           

2 **Algorithm 7:** Functions for a linear Gaussian system.

---

3 1 def **LinGaussPredict**( $\tilde{\mu}$ ,  $\tilde{\Sigma}$ ,  $\mathbf{A}$ ,  $\mathbf{b}$ ,  $\Omega$ ) :

4 2  $\mathbf{m} = \mathbf{A} \tilde{\mu} + \mathbf{b}$

5 3  $\mathbf{S} = \Omega + \mathbf{A} \tilde{\Sigma} \mathbf{A}^\top$

6 4  $\mathbf{C} = \tilde{\Sigma} \mathbf{A}^\top$

7 5 Return ( $\mathbf{m}$ ,  $\mathbf{S}$ ,  $\mathbf{C}$ )

8           

9 6 def **GaussCondition**( $\tilde{\mu}$ ,  $\tilde{\Sigma}$ ,  $\mathbf{m}$ ,  $\mathbf{S}$ ,  $\mathbf{C}$ ,  $\mathbf{y}$ ) :

10 7  $\mathbf{K} = \mathbf{C} \mathbf{S}^{-1}$

11 8  $\hat{\mu} = \tilde{\mu} + \mathbf{K}(\mathbf{y} - \mathbf{m})$

12 9  $\hat{\Sigma} = \tilde{\Sigma} - \mathbf{K} \mathbf{S} \mathbf{K}^\top$

13 10  $\ell = \log \mathcal{N}(\mathbf{y} | \mathbf{m}, \mathbf{S})$

14 11 Return ( $\hat{\mu}$ ,  $\hat{\Sigma}$ ,  $\ell$ )

15           

16 12 def **LinGaussUpdate**( $\tilde{\mu}$ ,  $\tilde{\Sigma}$ ,  $\mathbf{A}$ ,  $\mathbf{b}$ ,  $\Omega$ ,  $\mathbf{y}$ ) :

17 13  $(\mathbf{m}, \mathbf{S}, \mathbf{C}) = \text{LinGaussPredict}(\tilde{\mu}, \tilde{\Sigma}, \mathbf{A}, \mathbf{b}, \Omega)$

18 14  $(\hat{\mu}, \hat{\Sigma}, \ell) = \text{GaussCondition}(\tilde{\mu}, \tilde{\Sigma}, \mathbf{m}, \mathbf{S}, \mathbf{C}, \mathbf{y})$

19 20 15 Return ( $\hat{\mu}$ ,  $\hat{\Sigma}$ ,  $\ell$ )

---

21           

22           

23           

24 In the update step, we compute the joint distribution

25

$$p(\mathbf{z}_t, \mathbf{y}_t | \mathbf{y}_{1:t-1}) = \mathcal{N}\left(\begin{pmatrix} \boldsymbol{\mu}_{t|t-1} \\ \mathbf{m}_t \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma}_{t|t-1} & \mathbf{C}_t \\ \mathbf{C}_t^\top & \mathbf{S}_t \end{pmatrix}\right) \quad (8.101)$$

26

$$(\mathbf{m}_t, \mathbf{S}_t, \mathbf{C}_t) = \text{LinGaussPredict}(\boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}, \mathbf{H}_t, \mathbf{D}_t \mathbf{u}_t + \mathbf{d}_t, \mathbf{R}_t) \quad (8.102)$$

27 We then condition this on the observations to get the posterior distribution

28

$$p(\mathbf{z}_t | \mathbf{y}_t, \mathbf{y}_{1:t-1}) = p(\mathbf{z}_t | \mathbf{y}_{1:t}) = \mathcal{N}(\boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t}) \quad (8.103)$$

29

$$(\boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t}, \ell_t) = \text{GaussCondition}(\boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}, \mathbf{m}_t, \mathbf{S}_t, \mathbf{C}_t, \mathbf{y}_t) \quad (8.104)$$

30 We can combine the two parts of the update operation by defining the helper function **LinGaussUpdate**  
31 in Algorithm 7. The overall KF algorithm is shown in Algorithm 8.

---

32           

33 **Algorithm 8:** Kalman filter.

---

34

35 1 def **KF**( $\mathbf{F}_{1:T}$ ,  $\mathbf{B}_{1:T}$ ,  $\mathbf{b}_{1:T}$ ,  $\mathbf{Q}_{1:T}$ ,  $\mathbf{H}_{1:T}$ ,  $\mathbf{D}_{1:T}$ ,  $\mathbf{d}_{1:T}$ ,  $\mathbf{R}_{1:T}$ ,  $\mathbf{u}_{1:T}$ ,  $\mathbf{y}_{1:T}$ ,  $\boldsymbol{\mu}_{0|0}$ ,  $\boldsymbol{\Sigma}_{0|0}$ ) :

36 2 **foreach**  $t = 1 : T$  **do**

37 3      $(\boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}, -) = \text{LinGaussPredict}(\boldsymbol{\mu}_{t-1|t-1}, \boldsymbol{\Sigma}_{t-1|t-1}, \mathbf{F}_t, \mathbf{B}_t \mathbf{u}_t + \mathbf{b}_t, \mathbf{Q}_t)$

38 4      $(\boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t}, \ell_t) = \text{LinGaussUpdate}(\boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}, \mathbf{H}_t, \mathbf{D}_t \mathbf{u}_t + \mathbf{d}_t, \mathbf{R}_t)$

39 5 Return  $(\boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t})_{t=1}^T, \sum_{t=1}^T \ell_t$

---

40           

41           

42           

43           

44           

45           

46           

47

1 **8.3.2.6 Numerical issues**

3 In practice, the Kalman filter can encounter numerical issues. One solution is to use the **information**  
4 **filter**, which recursively updates the natural parameters of the Gaussian,  $\Lambda_{t|t} = \Sigma_{t|t}^{-1}$  and  $\eta_{t|t} =$   
5  $\Lambda_t \mu_{t|t}$ , instead of the mean and covariance (see Section 8.3.4). Another solution is the **square root**  
6 **filter**, which works with the Cholesky or QR decomposition of  $\Sigma_{t|t}$ , which is much more numerically  
7 stable than directly updating  $\Sigma_{t|t}$ . These techniques can be combined to create the **square root**  
8 **information filter** (SRIF) [May79]. (According to [Bie06], the SRIF was developed in 1969 for  
9 use in JPL's Mariner 10 mission to Venus.) In [Tol22] they present an approach which uses QR  
10 decompositions instead of matrix inversions, which can also be more stable.  
11

12

13 **8.3.3 The Kalman (RTS) smoother**

14

15 In Section 8.3.2, we described the Kalman filter, which sequentially computes  $p(\mathbf{z}_t | \mathbf{y}_{1:t})$  for each  $t$ .  
16 This is useful for online inference problems, such as tracking. However, in an offline setting, we can  
17 wait until all the data has arrived, and then compute  $p(\mathbf{z}_t | \mathbf{y}_{1:T})$ . By conditioning on past and future  
18 data, our uncertainty will be significantly reduced. This is illustrated in Figure 8.8(c), where we see  
19 that the posterior covariance ellipsoids are smaller for the smoothed trajectory than for the filtered  
20 trajectory.

21 We now explain how to compute the smoothed estimates, using an algorithm called the **RTS**  
22 **Smoother** or **RTSS**, named after its inventors, Rauch, Tung and Striebel [RTS65]. It is also  
23 known as the **Kalman smoothing** algorithm. The algorithm is the linear-Gaussian analog to the  
24 forwards-filtering backwards-smoothing algorithm for HMMs in Section 8.2.3.  
25

26 **8.3.3.1 Algorithm**

27 One can show (see Section 8.3.3.2) that one step of the backwards smoothing pass can be implemented  
28 as follows:

29

$$\underline{31} \quad (\boldsymbol{\mu}_{t+1|t}, \boldsymbol{\Sigma}_{t+1|t}, \boldsymbol{\Sigma}_{t,t+1|t}) = \text{LinGaussPredict}(\boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t}, \mathbf{F}_{t+1}, \mathbf{B}_{t+1}\mathbf{u}_{t+1} + \mathbf{b}_{t+1}, \mathbf{Q}_{t+1}) \quad (8.105)$$

$$\underline{32} \quad (\boldsymbol{\mu}_{t|T}, \boldsymbol{\Sigma}_{t|T}) = \text{GaussSmoothUpdate}(\boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t}, \boldsymbol{\Sigma}_{t,t+1|t}, \boldsymbol{\mu}_{t+1|T}, \boldsymbol{\Sigma}_{t+1|T}) \quad (8.106)$$

33 where the **GaussSmoothUpdate** function is defined in Algorithm 9. This uses the **backwards**  
34 **Kalman gain matrix**:

$$\underline{35} \quad \mathbf{G}_t = \boldsymbol{\Sigma}_{t,t+1|t} \boldsymbol{\Sigma}_{t+1|t}^{-1} = \boldsymbol{\Sigma}_{t|t} \mathbf{F}_{t+1}^T \boldsymbol{\Sigma}_{t+1|t}^{-1} \quad (8.107)$$

36 The algorithm can be initialized from  $\boldsymbol{\mu}_{T|T}$  and  $\boldsymbol{\Sigma}_{T|T}$  from the Kalman filter.  
37

38

39 **8.3.3.2 Derivation**

40

41 In this section, we derive the RTS smoother, following [Sar13, p136]. As in the derivation of the  
42 Kalman filter in Section 8.3.2.4, we make heavy use of the rules for manipulating linear Gaussian  
43 systems, discussed in Section 2.2.6.  
44

45

---

**Algorithm 9:** Updating a Gaussian distribution given updated beliefs on its child.

---

```

1 def GaussSmoothUpdate( $\mu_{t+1|t}$ ,  $\Sigma_{t+1|t}$ ,  $\Sigma_{t,t+1|t}$ ,  $\mu_{t+1|T}$ ,  $\Sigma_{t+1|T}$ )
2    $\mathbf{G}_t = \Sigma_{t,t+1|t} \Sigma_{t+1|t}^{-1}$ 
3    $\mu_{t|T} = \mu_{t|t} + \mathbf{G}_t(\mu_{t+1|T} - \mu_{t+1|t})$ 
4    $\Sigma_{t|T} = \Sigma_{t|t} + \mathbf{G}_t(\Sigma_{t+1|T} - \Sigma_{t+1|t})\mathbf{G}_t^\top$ 
5   Return  $(\mu_{t|T}, \Sigma_{t|T})$ 

```

---

The joint filtered distribution for two consecutive time slices is

$$p(\mathbf{z}_t, \mathbf{z}_{t+1} | \mathbf{y}_{1:t}) = p(\mathbf{z}_{t+1} | \mathbf{z}_t)p(\mathbf{z}_t | \mathbf{y}_{1:t}) = \mathcal{N}(\mathbf{z}_{t+1} | \mathbf{F}_t \mathbf{z}_t, \mathbf{Q}_t) \mathcal{N}(\mathbf{z}_t | \mu_{t|t-1}, \Sigma_{t|t}) \quad (8.108)$$

$$= \mathcal{N}\left(\begin{pmatrix} \mathbf{z}_t \\ \mathbf{z}_{t+1} \end{pmatrix} | \mathbf{m}_1, \mathbf{V}_1\right) \quad (8.109)$$

where

$$\mathbf{m}_1 = \begin{pmatrix} \mu_{t|t-1} \\ \mathbf{F}_t \mu_{t|t-1} \end{pmatrix}, \quad \mathbf{V}_1 = \begin{pmatrix} \Sigma_{t|t} & \Sigma_{t|t} \mathbf{F}_t^\top \\ \mathbf{F}_t \Sigma_{t|t} & \mathbf{F}_t \Sigma_{t|t} \mathbf{F}_t^\top + \mathbf{Q}_t \end{pmatrix} \quad (8.110)$$

By the Markov property for the hidden states we have

$$p(\mathbf{z}_t | \mathbf{z}_{t+1}, \mathbf{y}_{1:T}) = p(\mathbf{z}_t | \mathbf{z}_{t+1}, \mathbf{y}_{1:t}) \quad (8.111)$$

and hence by conditioning the joint distribution  $p(\mathbf{z}_t, \mathbf{z}_{t+1} | \mathbf{y}_{1:t})$  on the future state we get

$$p(\mathbf{z}_t | \mathbf{z}_{t+1}, \mathbf{y}_{1:T}) = \mathcal{N}(\mathbf{z}_t | \mathbf{m}_2(\mathbf{z}_{t+1}), \mathbf{V}_2) \quad (8.112)$$

$$\Sigma_{t+1|t} = \mathbf{F}_t \Sigma_{t|t} \mathbf{F}_t^\top + \mathbf{Q}_{t+1} \quad (8.113)$$

$$\mathbf{G}_t = \Sigma_{t|t} \mathbf{F}_t^\top \Sigma_{t+1|t}^{-1} \quad (8.114)$$

$$\mathbf{m}_2(\mathbf{z}_{t+1}) = \mu_{t|t-1} + \mathbf{G}_t(\mathbf{z}_{t+1} - \mathbf{F}_t \mu_{t|t-1}) \quad (8.115)$$

$$\mathbf{V}_2 = \Sigma_{t|t} - \mathbf{G}_t \Sigma_{t+1|t} \mathbf{G}_t^\top \quad (8.116)$$

The joint smoothed distribution of two consecutive time slices is

$$p(\mathbf{z}_{t+1}, \mathbf{z}_t | \mathbf{y}_{1:T}) = p(\mathbf{z}_{t+1} | \mathbf{y}_{1:T})p(\mathbf{z}_t | \mathbf{z}_{t+1}, \mathbf{y}_{1:T}) \quad (8.117)$$

$$= \mathcal{N}(\mathbf{z}_{t+1} | \mu_{t|T}, \Sigma_{t|T}) \mathcal{N}(\mathbf{z}_t | \mathbf{m}_2(\mathbf{z}_{t+1}), \mathbf{V}_2) \quad (8.118)$$

$$= \mathcal{N}\left(\begin{pmatrix} \mathbf{z}_{t+1} \\ \mathbf{z}_t \end{pmatrix} | \mathbf{m}_3, \mathbf{V}_3\right) \quad (8.119)$$

where

$$\mathbf{m}_3 = \begin{pmatrix} \mu_{t+1|T} \\ \mu_{t|t-1} + \mathbf{G}_t(\mu_{t+1} - \mathbf{F}_t \mu_{t|t-1}) \end{pmatrix}, \quad \mathbf{V}_3 = \begin{pmatrix} \Sigma_{t+1|T} & \Sigma_{t+1|T} \mathbf{G}_T^\top \\ \mathbf{G}_t \Sigma_{t+1|T} & \mathbf{G}_t \Sigma_{t+1|T} \mathbf{G}_t^\top + \mathbf{V}_2 \end{pmatrix} \quad (8.120)$$

1 From this, we can extract the smoothed marginal  
2

$$\underline{4} \quad p(\mathbf{z}_t | \mathbf{y}_{1:T}) = \mathcal{N}(\mathbf{z}_t | \boldsymbol{\mu}_{t|T}, \boldsymbol{\Sigma}_{t|T}) \quad (8.121)$$

$$\underline{5} \quad \boldsymbol{\mu}_{t|T} = \boldsymbol{\mu}_{t|t-1} + \mathbf{G}_t (\boldsymbol{\mu}_{t+1} - \mathbf{F}_t \boldsymbol{\mu}_{t|t-1}) \quad (8.122)$$

$$\underline{6} \quad \boldsymbol{\Sigma}_{t|T} = \mathbf{G}_t \boldsymbol{\Sigma}_{t+1|T} \mathbf{G}_t^\top + \mathbf{V}_2 = \boldsymbol{\Sigma}_{t|t} + \mathbf{G}_t (\boldsymbol{\Sigma}_{t+1|T} - \mathbf{F}_t \boldsymbol{\Sigma}_{t|t} \mathbf{F}_t^\top - \mathbf{Q}_{t+1}) \mathbf{G}_t^\top \quad (8.123)$$

### 9 8.3.3.3 Two-filter smoothing

10 Note that the backwards pass of the Kalman smoother does not need access to the observations,  $\mathbf{y}_{1:T}$ ,  
11 but does need access to the filtered belief states from the forwards pass,  $p(\mathbf{z}_t | \mathbf{y}_{1:t}) = \mathcal{N}(\mathbf{z}_t | \boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t})$ .  
12 There is an alternative version of the algorithm, known as **two-filter smoothing** [FP69; Kit04],  
13 in which we compute the forwards pass as usual, and then separately compute backwards mes-  
14 sages  $p(\mathbf{y}_{t+1:T} | \mathbf{z}_t) \propto \mathcal{N}(\mathbf{z}_t | \boldsymbol{\mu}_{t|t-1}^b, \boldsymbol{\Sigma}_{t|t}^b)$ , similar to the backwards filtering algorithm in HMMs  
15 (Section 8.2.4).

16 However, these backwards messages are conditional likelihoods, not posteriors, which can cause  
17 numerical problems. For example, consider  $t = T$ ; in this case, we need to set the initial covariance  
18 matrix to be  $\boldsymbol{\Sigma}_T^b = \infty \mathbf{I}$ , so that the backwards message has no effect on the filtered posterior (since  
19 there is no evidence beyond step  $T$ ). This problem can be resolved by working in information form.  
20 An alternative approach is to generalize the two-filter smoothing equations to ensure the likelihoods  
21 are normalizable by multiplying them by artificial distributions [BDM10].

22 In general, the RTS smoother is preferred to the two-filter smoother, since it is more numerically  
23 stable, and it is easier to generalize it to the nonlinear case.

24

25

### 26 8.3.3.4 Time and space complexity

27 In general, the Kalman smoothing algorithm takes  $O(N_y^3 + N_z^2 + N_y N_z)$  per step, where there  
28 are  $T$  steps. This can be slow when applied to long sequences. In [SGF21], they describe how to  
29 reduce this to  $O(\log T)$  steps using a **parallel prefix scan** operator that can be run efficiently on  
30 GPUs. In addition, we can reduce the space from  $O(T)$ , to  $O(\log T)$  using the same algorithm as in  
31 Section 8.2.6.

32

33

### 34 8.3.3.5 Forwards filtering backwards sampling

35 To draw posterior samples from the LG-SSM, we can leverage the following result (derived in [Sar13,  
36 p137]):

$$\underline{39} \quad p(\mathbf{z}_t | \mathbf{z}_{t+1}, \mathbf{y}_{1:T}) = \mathcal{N}(\mathbf{z}_t | \tilde{\boldsymbol{\mu}}_t, \tilde{\boldsymbol{\Sigma}}_t) \quad (8.124)$$

$$\underline{40} \quad \tilde{\boldsymbol{\mu}}_t = \boldsymbol{\mu}_{t|t} + \mathbf{G}_t (\mathbf{z}_{t+1} - \mathbf{F}_t \boldsymbol{\mu}_{t|t}) \quad (8.125)$$

$$\underline{41} \quad \tilde{\boldsymbol{\Sigma}}_t = \boldsymbol{\Sigma}_{t|t} - \mathbf{G}_t \boldsymbol{\Sigma}_{t+1|t} \mathbf{G}_t^\top = \boldsymbol{\Sigma}_{t|t} - \boldsymbol{\Sigma}_{t|t} \mathbf{F}_t^\top \boldsymbol{\Sigma}_{t+1|t}^{-1} \boldsymbol{\Sigma}_{t+1|t} \mathbf{G}_t^\top \quad (8.126)$$

$$\underline{42} \quad = \boldsymbol{\Sigma}_{t|t} (\mathbf{I} - \mathbf{F}_t^\top \mathbf{G}_t^\top) \quad (8.127)$$

43 where  $\mathbf{G}_t$  is the backwards Kalman gain defined in Equation (8.107).

44

---

### 8.3.4 Information form filtering and smoothing

*This section was written by Giles Harper-Donnelly.*

In this section, we derive the Kalman filter and smoother algorithms in information form. We will see that this is the “dual” of Kalman filtering / smoothing in moment form. In particular, while computing marginals in moment form is easy, computing conditionals is hard (requires a matrix inverse). Conversely, for information form, computing marginals is hard, but computing conditionals is easy.

#### 8.3.4.1 Filtering: algorithm

The predict step has a similar structure to the update step in moment form. We start with the prior  $p(\mathbf{z}_{t-1}|\mathbf{y}_{1:t-1}, \mathbf{u}_{1:t-1}) = \mathcal{N}_c(\mathbf{z}_{t-1}|\boldsymbol{\eta}_{t-1|t-1}, \boldsymbol{\Lambda}_{t-1|t-1})$  and then compute

$$p(\mathbf{z}_t|\mathbf{y}_{1:t-1}, \mathbf{u}_{1:t}) = \mathcal{N}_c(\mathbf{z}_t|\boldsymbol{\eta}_{t|t-1}, \boldsymbol{\Lambda}_{t|t-1}) \quad (8.128)$$

$$\mathbf{M}_t = \boldsymbol{\Lambda}_{t-1|t-1} + \mathbf{F}_t^\top \mathbf{Q}_t^{-1} \mathbf{F}_t \quad (8.129)$$

$$\mathbf{J}_t = \mathbf{Q}_t^{-1} \mathbf{F}_t \mathbf{M}_t^{-1} \quad (8.130)$$

$$\boldsymbol{\Lambda}_{t|t-1} = \mathbf{Q}_t^{-1} - \mathbf{Q}_t^{-1} \mathbf{F}_t (\boldsymbol{\Lambda}_{t-1|t-1} + \mathbf{F}_t^\top \mathbf{Q}_t^{-1} \mathbf{F}_t)^{-1} \mathbf{F}_t^\top \mathbf{Q}_t^{-1} \quad (8.131)$$

$$= \mathbf{Q}_t^{-1} - \mathbf{J}_t \mathbf{F}_t^\top \mathbf{Q}_t^{-1} \quad (8.132)$$

$$= \mathbf{Q}_t^{-1} - \mathbf{J}_t \mathbf{M}_t \mathbf{J}_t^\top \quad (8.133)$$

$$\boldsymbol{\eta}_{t|t-1} = \mathbf{J}_t \boldsymbol{\eta}_{t-1|t-1} + \boldsymbol{\Lambda}_{t|t-1} (\mathbf{B}_t \mathbf{u}_t + \mathbf{b}_t), \quad (8.134)$$

where  $\mathbf{J}_t$  is analogous to the Kalman gain matrix in moment form Equation (8.72). From the matrix inversion lemma Equation (2.54), we see that Equation (8.131) is the inverse of the predicted covariance  $\boldsymbol{\Sigma}_{t|t-1}$  given in Equation (8.68).

The update step in information form is as follows:

$$p(\mathbf{z}_t|\mathbf{y}_{1:t}, \mathbf{u}_{1:t}) = \mathcal{N}_c(\mathbf{z}_t|\boldsymbol{\eta}_{t|t}, \boldsymbol{\Lambda}_{t|t}) \quad (8.135)$$

$$\boldsymbol{\Lambda}_{t|t} = \boldsymbol{\Lambda}_{t|t-1} + \mathbf{H}_t^\top \mathbf{R}_t^{-1} \mathbf{H}_t \quad (8.136)$$

$$\boldsymbol{\eta}_{t|t} = \boldsymbol{\eta}_{t|t-1} + \mathbf{H}_t^\top \mathbf{R}_t^{-1} (\mathbf{y}_t - \mathbf{D}_t \mathbf{u}_t - \mathbf{d}_t). \quad (8.137)$$

#### 8.3.4.2 Filtering: derivation

For the predict step, we first derive the joint distribution over hidden states at  $t, t-1$ :

$$p(\mathbf{z}_{t-1}, \mathbf{z}_t|\mathbf{y}_{1:t-1}, \mathbf{u}_{1:t}) = p(\mathbf{z}_t|\mathbf{z}_{t-1}, \mathbf{u}_t)p(\mathbf{z}_{t-1}|\mathbf{y}_{1:t-1}, \mathbf{u}_{1:t-1}) \quad (8.138)$$

$$= \mathcal{N}_c(\mathbf{z}_t, |\mathbf{Q}_t^{-1}(\mathbf{F}_t \mathbf{z}_{t-1} + \mathbf{B}_t \mathbf{u}_t + \mathbf{b}_t), \mathbf{Q}_t^{-1}) \quad (8.139)$$

$$\times \mathcal{N}_c(\mathbf{z}_{t-1}, |\boldsymbol{\eta}_{t-1|t-1}, \boldsymbol{\Lambda}_{t-1|t-1}) \quad (8.140)$$

$$= \mathcal{N}_c(\mathbf{z}_{t-1}, \mathbf{z}_t|\boldsymbol{\eta}_{t-1,t|t}, \boldsymbol{\Lambda}_{t-1,t|t}) \quad (8.141)$$

1 where  
2

3

$$\underline{\eta}_{t-1,t|t-1} = \begin{pmatrix} \eta_{t-1|t-1} - \mathbf{F}_t^\top \mathbf{Q}_t^{-1} (\mathbf{B}_t \mathbf{u}_t + \mathbf{b}_t) \\ \mathbf{Q}_t^{-1} (\mathbf{B}_t \mathbf{u}_t + \mathbf{b}_t) \end{pmatrix} \quad (8.142)$$

4

$$\underline{\Lambda}_{t-1,t|t-1} = \begin{pmatrix} \Lambda_{t-1|t-1} + \mathbf{F}_t^\top \mathbf{Q}_t^{-1} \mathbf{F}_t & -\mathbf{F}_t^\top \mathbf{Q}_t^{-1} \\ -\mathbf{Q}_t^{-1} \mathbf{F}_t & \mathbf{Q}_t^{-1} \end{pmatrix} \quad (8.143)$$

5 The information form predicted parameters  $\eta_{t|t-1}, \Lambda_{t|t-1}$  can then be derived using the marginalisa-  
6 tion formulae in Section 2.2.5.4.

7 For the update step, we start with the joint distribution over the hidden state and the observation  
8 at  $t$ :

9

$$p(\mathbf{z}_t, \mathbf{y}_t | \mathbf{y}_{1:t-1}, \mathbf{u}_{1:t}) = p(\mathbf{y}_t | \mathbf{z}_t, \mathbf{u}_t) p(\mathbf{z}_t | \mathbf{y}_{1:t-1}, \mathbf{u}_{1:t-1}) \quad (8.144)$$

10

$$= \mathcal{N}_c(\mathbf{y}_t, |\mathbf{R}_t^{-1}(\mathbf{H}_t \mathbf{z}_t + \mathbf{D}_t \mathbf{u}_t + \mathbf{d}_t), \mathbf{R}_t^{-1}) \mathcal{N}_c(\mathbf{z}_t | \eta_{t|t-1}, \Lambda_{t|t-1}) \quad (8.145)$$

11

$$= \mathcal{N}_c(\mathbf{z}_t, \mathbf{y} | \eta_{z,y|t}, \Lambda_{z,y|t}) \quad (8.146)$$

12 where  
13

14

$$\eta_{z,y|t} = \begin{pmatrix} \eta_{t|t-1} - \mathbf{H}_t^\top \mathbf{R}_t^{-1} (\mathbf{D}_t \mathbf{u}_t + \mathbf{d}_t) \\ \mathbf{R}_t^{-1} (\mathbf{D}_t \mathbf{u}_t + \mathbf{d}_t) \end{pmatrix} \quad (8.147)$$

15

$$\Lambda_{z,y|t} = \begin{pmatrix} \Lambda_{t|t-1} + \mathbf{H}_t^\top \mathbf{R}_t^{-1} \mathbf{H}_t & -\mathbf{H}_t^\top \mathbf{R}_t^{-1} \\ -\mathbf{R}_t^{-1} \mathbf{H}_t & \mathbf{R}_t^{-1} \end{pmatrix} \quad (8.148)$$

16 The information form filtered parameters  $\eta_{t|t}, \Lambda_{t|t}$  are then derived using the conditional formulae in  
17 2.2.5.4.

18

### 19 8.3.4.3 Smoothing: algorithm

20 The smoothing equations are as follows:

21

$$p(\mathbf{z}_t | \mathbf{y}_{1:T}) = \mathcal{N}_c(\mathbf{z}_t | \eta_{t|T}, \Lambda_{t|T}) \quad (8.149)$$

22

$$\mathbf{U}_t = \mathbf{Q}_t^{-1} + \Lambda_{t+1|T} - \Lambda_{t+1|t} \quad (8.150)$$

23

$$\mathbf{L}_t = \mathbf{F}_t^\top \mathbf{Q}_t^{-1} \mathbf{U}_t^{-1} \quad (8.151)$$

24

$$\Lambda_{t|T} = \Lambda_{t|t} + \mathbf{F}_t^\top \mathbf{Q}_t^{-1} \mathbf{F}_t - \mathbf{L}_t \mathbf{Q}_t^{-1} \mathbf{F} \quad (8.152)$$

25

$$= \Lambda_{t|t} + \mathbf{F}_t^\top \mathbf{Q}_t^{-1} \mathbf{F}_t - \mathbf{L}_t \mathbf{U}_t \mathbf{L}_t^\top \quad (8.153)$$

26

$$\eta_{t|T} = \eta_{t|t} + \mathbf{L}_t (\eta_{t+1|T} - \eta_{t+1|t}). \quad (8.154)$$

27 The parameters  $\eta_{t|t}$  and  $\Lambda_{t|t}$  are the filtered values from Equations (8.137) and (8.136) respectively.

28 Similarly,  $\eta_{t+1|t}$  and  $\Lambda_{t+1|t}$  are the predicted parameters from Equations (8.134) and (8.131). The  
29 matrix  $\mathbf{L}_t$  is the information form analog to the backward Kalman gain matrix in Equation (8.107).

1

#### 8.3.4.4 Smoothing: derivation

2 From the generic forwards-filtering backwards-smoothing equation, Equation (8.20), we have

3

$$p(\mathbf{z}_t | \mathbf{y}_{1:T}) = p(\mathbf{z}_t | \mathbf{y}_{1:t}) \int \left[ \frac{p(\mathbf{z}_{t+1} | \mathbf{z}_t) p(\mathbf{z}_{t+1} | \mathbf{y}_{1:T})}{p(\mathbf{z}_{t+1} | \mathbf{y}_{1:t})} \right] d\mathbf{z}_{t+1} \quad (8.155)$$

4

$$= \int p(\mathbf{z}_t, \mathbf{z}_{t+1} | \mathbf{y}_{1:t}) \frac{p(\mathbf{z}_{t+1} | \mathbf{y}_{1:T})}{p(\mathbf{z}_{t+1} | \mathbf{y}_{1:t})} d\mathbf{z}_{t+1} \quad (8.156)$$

5

$$= \int \mathcal{N}_c(\mathbf{z}_t, \mathbf{z}_{t+1} | \boldsymbol{\eta}_{t,t+1|t}, \boldsymbol{\Lambda}_{t,t+1|t}) \frac{\mathcal{N}_c(\mathbf{z}_{t+1} | \boldsymbol{\eta}_{t+1|T}, \boldsymbol{\Lambda}_{t+1|T})}{\mathcal{N}_c(\mathbf{z}_{t+1} | \boldsymbol{\eta}_{t+1|t}, \boldsymbol{\Lambda}_{t+1|t})} d\mathbf{z}_{t+1} \quad (8.157)$$

6

$$= \int \mathcal{N}_c(\mathbf{z}_t, \mathbf{z}_{t+1} | \boldsymbol{\eta}_{t,t+1|T}, \boldsymbol{\Lambda}_{t,t+1|T}) d\mathbf{z}_{t+1}. \quad (8.158)$$

7 The parameters of the joint filtering predictive distribution,  $p(\mathbf{z}_t, \mathbf{z}_{t+1} | \mathbf{y}_{1:t})$ , take precisely the same form as those in the filtering derivation described in Section 8.3.4.2:

8

$$\boldsymbol{\eta}_{t,t+1|t} = \begin{pmatrix} \boldsymbol{\eta}_{t|t} \\ \mathbf{0} \end{pmatrix}, \quad \boldsymbol{\Lambda}_{t,t+1|t} = \begin{pmatrix} \boldsymbol{\Lambda}_{t|t} + \mathbf{F}_{t+1}^T \mathbf{Q}_{t+1}^{-1} \mathbf{F}_{t+1} & -\mathbf{F}_{t+1}^T \mathbf{Q}_{t+1}^{-1} \\ -\mathbf{Q}_{t+1}^{-1} \mathbf{F}_{t+1} & \mathbf{Q}_{t+1}^{-1} \end{pmatrix}, \quad (8.159)$$

9 We can now update this potential function by subtracting out the filtered information and adding in the smoothing information, using the rules for manipulating Gaussian potentials described in Section 2.2.7:

10

$$\boldsymbol{\eta}_{t,t+1|T} = \boldsymbol{\eta}_{t,t+1|t} + \begin{pmatrix} \mathbf{0} \\ \boldsymbol{\eta}_{t+1|T} \end{pmatrix} - \begin{pmatrix} \mathbf{0} \\ \boldsymbol{\eta}_{t+1|t} \end{pmatrix} = \begin{pmatrix} \boldsymbol{\eta}_{t|t} \\ \boldsymbol{\eta}_{t+1|T} - \boldsymbol{\eta}_{t+1|t} \end{pmatrix}, \quad (8.160)$$

11 and

12

$$\boldsymbol{\Lambda}_{t,t+1|T} = \boldsymbol{\Lambda}_{t,t+1|t} + \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Lambda}_{t+1|T} \end{pmatrix} - \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Lambda}_{t+1|t} \end{pmatrix} \quad (8.161)$$

13

$$= \begin{pmatrix} \boldsymbol{\Lambda}_{t|t} + \mathbf{F}_{t+1}^T \mathbf{Q}_{t+1}^{-1} \mathbf{F}_{t+1} & -\mathbf{F}_{t+1}^T \mathbf{Q}_{t+1}^{-1} \\ -\mathbf{Q}_{t+1}^{-1} \mathbf{F}_{t+1} & \mathbf{Q}_{t+1}^{-1} + \boldsymbol{\Lambda}_{t+1|T} - \boldsymbol{\Lambda}_{t+1|t} \end{pmatrix} \quad (8.162)$$

14 Applying the information form marginalisation formula Equation (2.46) leads to Equation (8.154) and Equation (8.152).

15

## 8.4 Inference for non-linear and/or non-Gaussian SSMS

16 In general state space models (Section 29.1), the transition and/or emission probability distributions can be arbitrary conditional distributions, and are not restricted to being discrete or Gaussian distributions. This makes inference much harder, because the likelihood is not conjugate to the prior, and the posterior may not have any convenient analytical form.

17

### 8.4.1 Inference based on discretization

18 To illustrate the problem, consider a simple 1d SSM with linear dynamics corrupted by additive Student noise:

19

$$z_t = z_{t-1} + \mathcal{T}_2(0, 1) \quad (8.163)$$

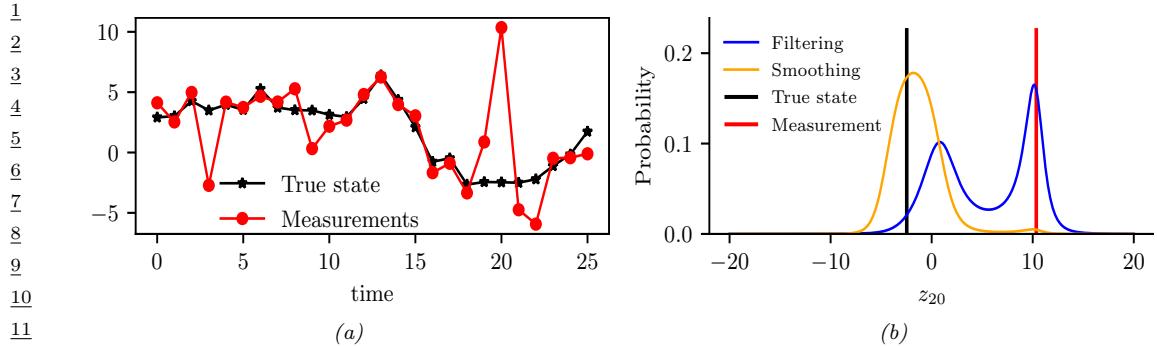


Figure 8.9: (a) Observations and true and estimated state. (b) Marginal distributions for time step  $t = 20$ . Generated by [discretized\\_ssm\\_student.ipynb](#).

The observations are also linear, and are also corrupted by additive Student noise:

$$y_t = z_t + \mathcal{T}_2(0, 1) \quad (8.164)$$

This robust observation model is useful when there are potential outliers in the observed data, such as at time  $t = 20$  in Figure 8.9a. (See also Section 8.6.2 for discussion of robust Kalman filters.)

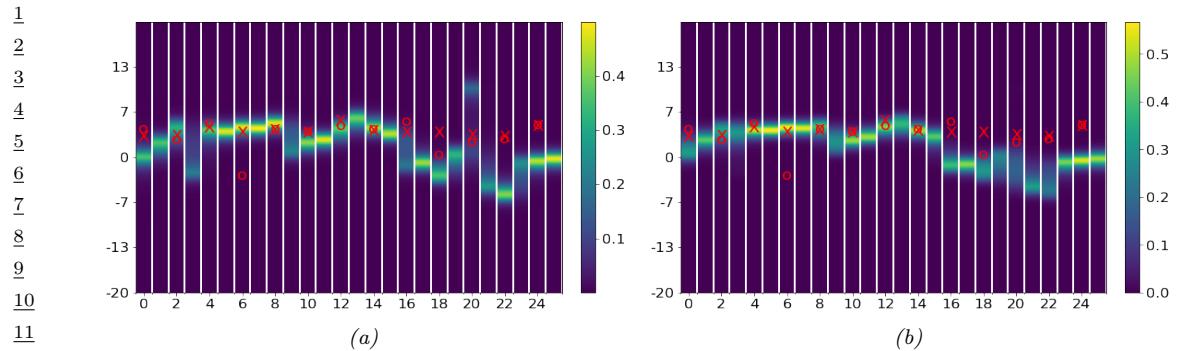
Unfortunately the use of a non-Gaussian likelihood means that the resulting posterior can become multimodal. In models with low-dimensional latent state spaces (e.g., 1-2d), we can always discretize the state space and apply the exact HMM filter and smoother from Section 8.2.4, as proposed in [RG17]. We show the results for filtering and smoother in Figure 8.10a and in Figure 8.10b. We see that at  $t = 20$ , the filtering distribution,  $p(z_t | \mathbf{y}_{1:20})$ , is bimodal, with a mean that is quite far from the true state (see Figure 8.9b for a detailed plot). Such a multi-modal distribution can be approximated by a suitably fine discretization. Unfortunately, tackling the inference problem by discretizing the state space will not scale to higher dimensional problems, due to the curse of dimensionality. In particular, we know that the HMM filter takes  $O(K^2)$  operations per time step, if there are  $K$  states. If we have  $N_z$  dimensions, each discretized into  $B$  bins, then we have  $K = B^{N_z}$ , so the approach becomes intractable beyond 1d.

#### 8.4.2 Inference based on Gaussian approximations

In the sections below, we discuss various deterministic approximate inference schemes which make a Gaussian approximation to the posterior. Most of these algorithms assume the SSM has the following form:

$$\begin{aligned} z_t &= \mathbf{f}(z_{t-1}, u_t) + \mathcal{N}(\mathbf{0}, \mathbf{Q}_t) \\ y_t &= \mathbf{h}(z_t, u_t) + \mathcal{N}(\mathbf{0}, \mathbf{R}_t) \end{aligned} \quad (8.165)$$

where  $\mathbf{z}_t \in \mathbb{R}^{N_z}$  is the hidden state,  $\mathbf{y}_t \in \mathbb{R}^{N_y}$  is the observation,  $\mathbf{u}_t \in \mathbb{R}^{N_u}$  are the optional inputs,  $\mathbf{f} : \mathbb{R}^{N_z+N_u} \rightarrow \mathbb{R}^{N_z}$  is the dynamics model, and  $\mathbf{h} : \mathbb{R}^{N_z+N_u} \rightarrow \mathbb{R}^{N_y}$  is the observation model. The dynamics and/or observations models can be nonlinear, but we assume the noise terms are Gaussian and additive. Our presentation is based on [Sar13, Ch. 5]; for more details, see e.g., [Sar13; Fan+17; Li+17e; Koy+10].)



*Figure 8.10: Discretized posterior of the latent state at each time step. Red cross is the true latent state. Red circle is observation. (a) Filtering. (b) Smoothing. Generated by [discretized\\_ssm\\_student.ipynb](#).*

## 8.5 Inference based on local linearization

In this section, we extend the Kalman filter and smoother to the case where the system dynamics and/or the observation model are nonlinear. (We continue to assume that the noise is additive Gaussian, as in Equation (8.165).) The basic idea is to linearize the dynamics and observation models about the previous state estimate using a first order Taylor series expansion, and then to apply the standard Kalman filter equations from Section 8.3.2. Intuitively we can think of this as approximating a stationary non-linear dynamical system with a non-stationary linear dynamical system. This approach is called the **extended Kalman filter** or **EKF**.

### 8.5.1 Taylor series expansion

Suppose  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and  $\mathbf{y} = \mathbf{g}(\mathbf{x})$ , where  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is a differentiable and invertible function. The pdf for  $\mathbf{y}$  is given by

$$p(\mathbf{y}) = |\det \text{Jac}(\mathbf{g}^{-1})(\mathbf{y})| \mathcal{N}(\mathbf{g}^{-1}(\mathbf{y})|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (8.166)$$

In general this is intractable to compute, so we seek an approximation.

Suppose  $\mathbf{x} = \boldsymbol{\mu} + \boldsymbol{\delta}$ , where  $\boldsymbol{\delta} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$ . Then we can form a first order Taylor series expansion of the function  $\mathbf{g}$  as follows:

$$\mathbf{g}(\mathbf{x}) = \mathbf{g}(\boldsymbol{\mu} + \boldsymbol{\delta}) \approx \mathbf{g}(\boldsymbol{\mu}) + \mathbf{G}(\boldsymbol{\mu})\boldsymbol{\delta} \quad (8.167)$$

where  $\mathbf{G}(\boldsymbol{\mu})$  is the Jacobian of  $\mathbf{g}$  at  $\boldsymbol{\mu}$ :

$$[\mathbf{G}(\boldsymbol{\mu})]_{jj'} = \frac{\partial g_j(\mathbf{x})}{\partial x_{j'}}|_{\mathbf{x}=\boldsymbol{\mu}} \quad (8.168)$$

We now derive the induced Gaussian approximation to  $\mathbf{y} = \mathbf{g}(\mathbf{x})$ . The mean is given by

$$\mathbb{E}[\mathbf{y}] \approx \mathbb{E}[\mathbf{g}(\boldsymbol{\mu}) + \mathbf{G}(\boldsymbol{\mu})\boldsymbol{\delta}] = \mathbf{g}(\boldsymbol{\mu}) + \mathbf{G}(\boldsymbol{\mu})\mathbb{E}[\boldsymbol{\delta}] = \mathbf{g}(\boldsymbol{\mu}) \quad (8.169)$$

1 The covariance is given by  
2

$$\text{Cov} [\mathbf{y}] = \mathbb{E} [(\mathbf{g}(\mathbf{x}) - \mathbb{E}[\mathbf{g}(\mathbf{x})])(\mathbf{g}(\mathbf{x}) - \mathbb{E}[\mathbf{g}(\mathbf{x})])^\top] \quad (8.170)$$

$$\approx \mathbb{E} [(\mathbf{g}(\mathbf{x}) - \mathbf{g}(\boldsymbol{\mu}))(\mathbf{g}(\mathbf{x}) - \mathbf{g}(\boldsymbol{\mu}))^\top] \quad (8.171)$$

$$\approx \mathbb{E} [(\mathbf{g}(\boldsymbol{\mu}) + \mathbf{G}(\boldsymbol{\mu})\boldsymbol{\delta} - \mathbf{g}(\boldsymbol{\mu}))(\mathbf{g}(\boldsymbol{\mu}) + \mathbf{G}(\boldsymbol{\mu})\boldsymbol{\delta} - \mathbf{g}(\boldsymbol{\mu}))^\top] \quad (8.172)$$

$$= \mathbb{E} [(\mathbf{G}(\boldsymbol{\mu})\boldsymbol{\delta})(\mathbf{G}(\boldsymbol{\mu})\boldsymbol{\delta})^\top] \quad (8.173)$$

$$= \mathbf{G}(\boldsymbol{\mu})\mathbb{E} [\boldsymbol{\delta}\boldsymbol{\delta}^\top]\mathbf{G}(\boldsymbol{\mu})^\top \quad (8.174)$$

$$= \mathbf{G}(\boldsymbol{\mu})\Sigma\mathbf{G}(\boldsymbol{\mu})^\top \quad (8.175)$$

13  
14 **Algorithm 10:** Linearized approximation to a joint Gaussian distribution.  
15

---

16 1 def **LinearizedPredict**( $\boldsymbol{\mu}$ ,  $\Sigma$ ,  $\mathbf{g}$ ,  $\Omega$ ) :  
17 2  $\mathbf{m} = \mathbf{g}(\boldsymbol{\mu})$   
18 3  $\mathbf{G} = \text{Jac}(\mathbf{g})(\boldsymbol{\mu})$   
19 4  $\mathbf{S} = \mathbf{G}\Sigma\mathbf{G}^\top + \Omega$   
20 5  $\mathbf{C} = \Sigma\mathbf{G}^\top$   
21 6 Return ( $\boldsymbol{\mu}$ ,  $\mathbf{S}$ ,  $\mathbf{C}$ )

---

22 When deriving the EKF, we need to compute the joint distribution  $p(\mathbf{x}, \mathbf{y})$  where  
24

$$\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma), \mathbf{y} = \mathbf{g}(\mathbf{x}) + \mathbf{q}, \mathbf{q} \sim \mathcal{N}(\mathbf{0}, \Omega) \quad (8.176)$$

27 where  $\mathbf{q}$  is independent of  $\mathbf{x}$ . We can compute this by defining the augmented function  $\tilde{\mathbf{g}}(\mathbf{x}) = [\mathbf{x}, \mathbf{g}(\mathbf{x})]$   
28 and following the procedure above. The resulting linear approximation to the joint is  
29

$$\begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} \boldsymbol{\mu} \\ \mathbf{m} \end{pmatrix}, \begin{pmatrix} \Sigma & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{S} \end{pmatrix} \right) \quad (8.177)$$

32 where the parameters are computed using Algorithm 10. We can then condition this joint Gaussian  
34 on the observed value  $\mathbf{y}$  to get the posterior.

35 It is also possible to derive an approximation for the case of non-additive Gaussian noise, where  
36  $\mathbf{y} = \mathbf{g}(\mathbf{x}, \mathbf{q})$ . See [Sar13, Sec 5.1] for details.

37

### 38 8.5.2 The extended Kalman filter (EKF) 39

40 We now derive the extended Kalman filter for performing approximate inference in the model given by  
41 Equation (8.165). We first linearize the dynamics model around  $\boldsymbol{\mu}_{t-1|t-1}$  to get an approximation to  
42 the one-step-ahead predictive distribution  $p(\mathbf{z}_t | \mathbf{y}_{1:t-1}, \mathbf{u}_{1:t}) = \mathcal{N}(\mathbf{z}_t | \boldsymbol{\mu}_{t|t-1}, \Sigma_{t|t-1})$ . We then linearize  
43 the observation model around  $\boldsymbol{\mu}_{t|t-1}$ , and then performing a Gaussian update. (In Section 8.5.3, we  
44 consider linearizing around a different point that gives better accuracy.)

45 In more detail, we can write one step of the EKF algorithm as follows, where we use the notation  
46 from Section 8.3.2.5:

47

$$(\boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}, -) = \text{LinearizedPredict}(\boldsymbol{\mu}_{t-1|t-1}, \boldsymbol{\Sigma}_{t-1|t-1}, \mathbf{f}(\cdot, \mathbf{u}_t), \mathbf{Q}_t) \quad (8.178)$$

$$(\mathbf{m}_t, \mathbf{S}_t, \mathbf{C}_t) = \text{LinearizedPredict}(\boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}, \mathbf{h}(\cdot, \mathbf{u}_t), \mathbf{R}_t) \quad (8.179)$$

$$(\boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t}, \ell_t) = \text{GaussCondition}(\boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}, \mathbf{m}_t, \mathbf{S}_t, \mathbf{C}_t, \mathbf{y}_t) \quad (8.180)$$

The EKF is widely used because it is simple and relatively efficient. However, there are two cases when the EKF works poorly. The first is when the prior covariance is large. In this case, the prior distribution is broad, so we end up sending a lot of probability mass through different parts of the function that are far from  $\boldsymbol{\mu}_{t-1|t-1}$ , where the function has been linearized. The other setting where the EKF works poorly is when the function is highly nonlinear near the current mean (see Figure 8.12a).

A more accurate approach is to use a second-order Taylor series approximation, known as the **second order EKF**. The resulting updates can still be computed in closed form (see [Sar13, Sec 5.2] for details). Other more accurate extensions are discussed in Section 8.5.3 and Section 8.7.2.

### 8.5.3 Iterated EKF

**Algorithm 11:** Iterated extended Kalman filter.

```

1 def IEKF( $\mathbf{f}, \mathbf{Q}, \mathbf{h}, \mathbf{R}, \mathbf{y}_{1:T}, J$ ) :
2   foreach  $t = 1 : T$  do
3     Predict step:
4      $(\boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}, -) = \text{LinearizedPredict}(\boldsymbol{\mu}_{t-1|t-1}, \boldsymbol{\Sigma}_{t-1|t-1}, \mathbf{f}(\cdot, \mathbf{u}_t), \mathbf{Q}_t)$ 
5     Update step:
6      $\boldsymbol{\mu}_{t|t} = \boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t} = \boldsymbol{\Sigma}_{t|t-1}$ 
7     foreach  $j = 1 : J$  do
8        $(\mathbf{m}_t, \mathbf{S}_t, \mathbf{C}_t) = \text{LinearizedPredict}(\boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t}, \mathbf{h}(\cdot, \mathbf{u}_t), \mathbf{R}_t)$ 
9        $(\boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t}, \ell_t) = \text{GaussCondition}(\boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}, \mathbf{m}_t, \mathbf{S}_t, \mathbf{C}_t, \mathbf{y}_t)$ 
31  Return  $(\boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t})_{t=1}^T$ 
```

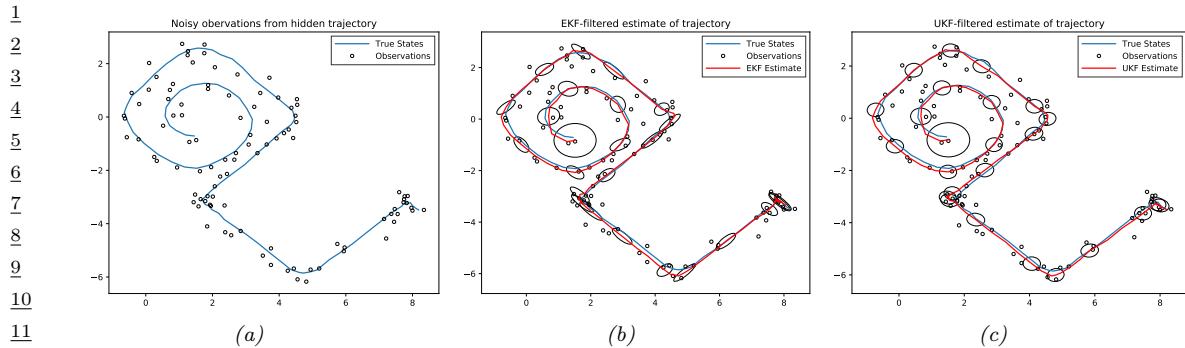
We can improve performance of the EKF by repeatedly re-linearizing the measurement model around the current posterior,  $\boldsymbol{\mu}_{t|t}$ , instead of  $\boldsymbol{\mu}_{t|t-1}$ ; this is called the **iterated EKF** [BC93]. See Algorithm 11 for the pseudocode. Unfortunately the IEKF can diverge. A robust IEKF method, that uses line search to perform damped (partial) updates of  $\boldsymbol{\mu}_{t|t}^j$ , is presented in [SHA15].

### 8.5.4 The extended Kalman smoother (EKS)

We can extend the EKF to the offline smoothing case, resulting in the **extended Kalman smoother**, also called the **extended RTS smoother**. We just need to linearize around the filtered mean, and then apply the standard RTS update equations:

$$(\boldsymbol{\mu}_{t+1|t}, \boldsymbol{\Sigma}_{t+1|t}, \boldsymbol{\Sigma}_{t,t+1|t}) = \text{LinearizedPredict}(\boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t}, \mathbf{f}(\cdot, \mathbf{u}_{t+1}), \mathbf{Q}_{t+1}) \quad (8.181)$$

$$(\boldsymbol{\mu}_{t|T}, \boldsymbol{\Sigma}_{t|T}) = \text{GaussSmoothUpdate}(\boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t}, \boldsymbol{\Sigma}_{t,t+1|t}, \boldsymbol{\mu}_{t+1|T}, \boldsymbol{\Sigma}_{t+1|T}) \quad (8.182)$$



*Figure 8.11: Illustration of filtering applied to a 2d nonlinear dynamical system. (a) True underlying state and observed data. (b) Extended Kalman filter estimate. Generated by `ekf_spiral.ipynb`. (c) Unscented Kalman filter estimate. Generated by `ukf_spiral.ipynb`.*

See [Sar13, Sec 9.1] for more details.

For improved accuracy, you can use the **iterated EKS** which relinearizes the model at the previous MAP estimate (see Section 8.9.4).

### 8.5.5 Examples

In this section, we give some examples of EKF/EKS.

#### 8.5.5.1 Tracking a point spiraling in 2d

In Section 8.3.1.1, we considered an example of state estimation and tracking of an object moving in 2d under a linear dynamics model with a linear observation model. However, motion and observation models are often nonlinear. For example, consider an object that is moving along a curved trajectory, such as this:

$$\mathbf{f}(\mathbf{z}) = (z_1 + \Delta \sin(z_2), z_2 + \Delta \cos(z_1)) \quad (8.183)$$

where  $\Delta$  is the step size. For simplicity, we assume full visibility of the state vector (modulo observation noise), so  $\mathbf{h}(\mathbf{z}) = \mathbf{z}$ .

Despite the simplicity of this model, exact inference is intractable. However, we can easily apply the EKF. The results are shown in Figure 8.11b.

#### 8.5.5.2 Neural network training

In Section 17.5.2, we show how to use the EKF to perform online parameter inference for an MLP regression model.

## 8.6 Other variants of the Kalman filter

In this section, we briefly mention some other variants of Kalman filtering. For a more extensive review, see [Li+17e].

---

### 8.6.1 Ensemble Kalman filter

The **ensemble Kalman filter (EnKF)** is a technique developed in the geoscience (meteorology) community to perform approximate online inference in large nonlinear systems. In particular, it is mostly used for problems where the hidden state represents an unknown physical quantity (e.g., temperature and pressure) at each point on a spatial grid, and the measurements are sparse and spatially localized. Combining this information over space and time is called **data assimilation**.

The canonical reference is [Eve09], but a more accessible tutorial (using the same Bayesian signal processing approach we adopt in this chapter) is in [Rot+17].

The key idea is to represent the belief state  $p(\mathbf{z}_t | \mathbf{y}_{1:t})$  by a finite number of samples  $\mathbf{Z}_{t|t} = \{\mathbf{z}_{t|t}^s : s = 1 : N_s\}$ , where each  $\mathbf{z}_{t|t}^s \in \mathbb{R}^{N_z}$ . In contrast to particle filtering (Section 13.2), the samples are updated in a manner that closely resembles the Kalman filter, so there is no importance sampling or resampling step. The downside is that the posterior does not converge to the true Bayesian posterior even as  $N_s \rightarrow \infty$  [LGMT11], except in the linear-Gaussian case. However, sometimes the performance of EnKF can be better for small number of samples (although this depends of course on the PF proposal distribution).

The posterior mean and covariance can be derived from the ensemble of samples as follows:

$$\tilde{\boldsymbol{\mu}}_{t|t} = \frac{1}{N_s} \sum_{s=1}^{N_s} \mathbf{z}_{t|t}^s = \frac{1}{N_s} \mathbf{Z}_{t|t} \mathbf{1} \quad (8.184)$$

$$\tilde{\boldsymbol{\Sigma}}_{t|t} = \frac{1}{N_s - 1} \sum_{s=1}^{N_s} (\mathbf{z}_{t|t}^s - \tilde{\boldsymbol{\mu}}_{t|t})(\mathbf{z}_{t|t}^s - \tilde{\boldsymbol{\mu}}_{t|t})^\top = \frac{1}{N_s - 1} \tilde{\mathbf{Z}}_{t|t} \tilde{\mathbf{Z}}_{t|t}^\top \quad (8.185)$$

where  $\tilde{\mathbf{Z}}_{t|t} = \mathbf{Z}_{t|t} - \tilde{\boldsymbol{\mu}}_{t|t} \mathbf{1}^\top$ .

We update the samples as follows. For the time update, we first draw  $N_s$  system noise variables  $\mathbf{q}_t^s \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_t)$ , and then we pass these, and the previous state estimate, through the dynamics model to get the one-step-ahead state predictions,  $\mathbf{z}_{t|t-1}^s = \mathbf{f}(\mathbf{z}_{t-1|t-1}^s, \mathbf{q}_t^s)$ , from which we get  $\mathbf{Z}_{t|t-1} = \{\mathbf{z}_{t|t-1}^s\}$ , which has size  $N_z \times N_s$ . Next we draw  $N_s$  observation noise variables  $\mathbf{r}_t^s \sim \mathcal{N}(\mathbf{0}, \mathbf{R}_t)$ , and use them to compute the one-step-ahead observation predictions,  $\mathbf{y}_{t|t-1}^s = \mathbf{h}(\mathbf{z}_{t|t-1}^s, \mathbf{r}_t^s)$  and  $\mathbf{Y}_{t|t-1} = \{\mathbf{y}_{t|t-1}^s\}$ , which has size  $N_y \times N_s$ . Finally we compute the measurement update using

$$\mathbf{Z}_{t|t} = \mathbf{Z}_{t|t-1} + \tilde{\mathbf{K}}_t (\mathbf{y}_t \mathbf{1}^\top - \mathbf{Y}_{t|t-1}) \quad (8.186)$$

which is the analog of Equation (8.74).

We now discuss how to compute  $\tilde{\mathbf{K}}_t$ , which is the analog of the Kalman gain matrix in Equation (8.72). First note that we can write the exact Kalman gain matrix (in the linear-Gaussian case) as  $\mathbf{K}_t = \Sigma_{t|t-1} \mathbf{H}^\top \mathbf{S}_t^{-1} = \mathbf{C}_t \mathbf{S}_t^{-1}$ , where  $\mathbf{S}_t$  is the covariance of the measurements, and  $\mathbf{C}_t$  is the cross-covariance between the state and output predictions. In the EnKF, we approximate  $\mathbf{S}_t$  and  $\mathbf{C}_t$  empirically as follows. First we compute the deviations from predictions:

$$\tilde{\mathbf{Z}}_{t|t-1} = \mathbf{Z}_{t|t-1} (\mathbf{I} - \frac{1}{N} \mathbf{1} \mathbf{1}^\top), \quad \tilde{\mathbf{Y}}_{t|t-1} = \mathbf{Y}_{t|t-1} (\mathbf{I} - \frac{1}{N} \mathbf{1} \mathbf{1}^\top) \quad (8.187)$$

Then we compute the sample covariance matrices

$$\tilde{\mathbf{S}}_t = \frac{1}{N_s - 1} \tilde{\mathbf{Y}}_{t|t-1} \tilde{\mathbf{Y}}_{t|t-1}^\top, \quad \tilde{\mathbf{C}}_t = \frac{1}{N_s - 1} \tilde{\mathbf{Z}}_{t|t-1} \tilde{\mathbf{Y}}_{t|t-1}^\top \quad (8.188)$$

1 Finally we compute  
2

$$\tilde{\mathbf{K}}_t = \tilde{\mathbf{C}}_t \tilde{\mathbf{S}}_t^{-1} \quad (8.189)$$

5 We now compare the computational complexity to the KF algorithm. Recall that  $N_z$  is the number  
6 of latent dimensions  $N_y$  is the number of observed dimensions, and  $N_s$  is the number of samples. We  
7 will assume  $N_z > N_s > N_y$ , as occurs in most geospatial problems. The EnKF time update takes  
8  $O(N_z^2 N_s)$  operations, and the measurement update takes  $O(N_z N_y N_s)$ . By contrast, in the KF, the  
9 time update takes  $O(N_z^3)$  operations, and the measurement update takes  $O(N_z^2 N_y)$ . So we see that  
10 the EnKF is faster for high dimensional state spaces, because it uses a low-rank approximation to  
11 the posterior covariance.

12 Unfortunately, if  $N_s$  is too small, the EnKF can become overconfident, and the filter can diverge.  
13 Various heuristics (e.g., covariance inflation) have been proposed to fix this. However, most of these  
14 methods are ad-hoc. A variety of more well-principled solutions have also been proposed, see e.g.,  
15 [FK13b; Rei13].

16

### 17 8.6.2 Robust Kalman filters

18

19 In practice we often have noise that is non-Gaussian. A common example is when we have clutter, or  
20 outliers, in the observation model, or sudden changes in the process model. In this case, we might  
21 use the Laplace distribution [Ara+09] or the Student- $t$  distribution [Ara10; RÖG13; Ara+17] as  
22 noise models.

23 [Hua+17b] proposes a variational Bayes (Section 10.2.3) approach, that allows the dynamical  
24 prior and the observation model to both be (linear) Student distributions, but where the posterior is  
25 approximated at each step using a Gaussian, conditional on the noise scale matrix, which is modeled  
26 using an inverse Wishart distribution. An extension of this, to handle mixture distributions, can be  
27 found in [Hua+19].

28

### 29 8.6.3 Dual EKF

30

31 In this section, we briefly discuss one approach to estimating the parameters of an SSM. In an offline  
32 setting, we can use EM, SGD or Bayesian inference to compute an approximation to  $p(\boldsymbol{\theta}|\mathbf{y}_{1:T})$  (see  
33 Section 29.8). In the online setting, we want to compute  $p(\boldsymbol{\theta}_t|\mathbf{y}_{1:t})$ . We can do this by adding the  
34 parameters to the state space, possibly with an **artificial dynamics**,  $p(\boldsymbol{\theta}_t|\boldsymbol{\theta}_{t-1}) = \mathcal{N}(\boldsymbol{\theta}_t|\boldsymbol{\theta}_{t-1}, \epsilon \mathbf{I})$ ,  
35 and then performing joint inference of states and parameters. The latent variables at each step  
36 now contain the latent states,  $\mathbf{z}_t$ , and the latent parameters,  $\boldsymbol{\theta}_t$ . One approach to performing  
37 approximating inference in such a model is to use the **dual EKF**, in which one EKF performs state  
38 estimation and the other EKF performs parameter estimation [WN01].

39

## 40 8.7 Inference based on the unscented transform

41

42 In this section, we replace the local linearization of the model with a different approximation known  
43 as the **unscented transform**. When applied to Bayesian filtering, we get the **unscented Kalman**  
44 **filter (UKF)**, since it is a version of the EKF that “doesn’t stink” [JU97; JUDW00]. The key idea  
45 is this: instead of computing a linear approximation to the function and then passing a Gaussian  
46 through it, we instead pass a deterministically chosen set of points, known as **sigma points**, through  
47

the function, and fit a Gaussian to the resulting transformed points; this is known as the unscented transform (see Section 8.7). Using the unscented transform for the transition and observation models gives the the overall method, which is also called a **sigma point filter** [VDMW03].

The main advantage of the UKF over the EKF is that it can be more accurate, and more stable. (The EKF can sometimes lead to large errors and divergence of the filter [IX00; VDMW03].) In addition, the UKF does not need to compute Jacobians of the observation and dynamics model, so it can be applied to non-differentiable models, or ones with hard constraints. However, the UKF can be slower, since it requires  $N_z$  evaluations of the dynamics and observation models. In addition, it has 3 hyper-parameters that need to be set.

### 8.7.1 The unscented transform

---

**Algorithm 12:** Computing sigma points using unscented transform.

---

```

1 def SigmaPoints( $\mu, \Sigma; \alpha, \beta, \kappa$ ) :
2    $n = \text{dimensionality of } \mu$ 
3    $\lambda = \alpha^2(n + \kappa) - n$ 
4   Compute a set of  $2n + 1$  sigma points:
5      $\mathcal{X}_0 = \mu, \mathcal{X}_i = \mu + \sqrt{n + \lambda} [\sqrt{\Sigma}]_{:,i}, \mathcal{X}_{i+n} = \mu - \sqrt{n + \lambda} [\sqrt{\Sigma}]_{:,i}$ 
6   Compute a set of  $2n + 1$  weights for the mean and covariance:
7      $w_0^m = \frac{\lambda}{n+\lambda}, w_0^c = \frac{\lambda}{n+\lambda} + (1 - \alpha^2 + \beta), w_i^m = w_i^c = \frac{1}{2(n+\lambda)}$ 
8   Return  $(\mathcal{X}_{0:2n}, w_{0:2n}^m, w_{0:2n}^c)$ 

```

---



---

**Algorithm 13:** Unscented approximation to a joint Gaussian distribution.

---

```

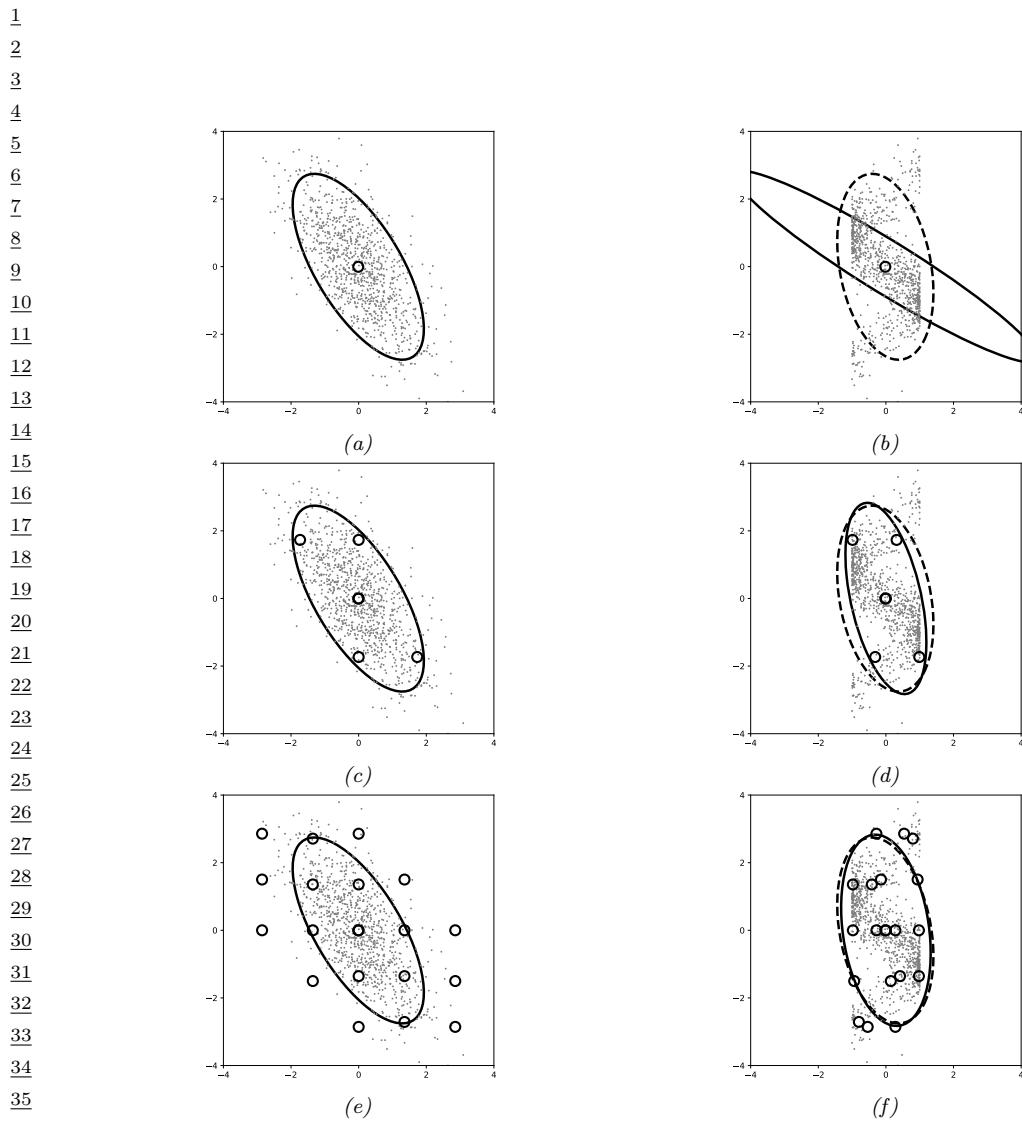
1 def UnscentedPredict( $\mu, \Sigma, g, \Omega; \alpha, \beta, \kappa$ ) :
2    $(\mathcal{X}_{0:2n}, w_{0:2n}^m, w_{0:2n}^c) = \text{SigmaPoints}(\mu, \Sigma; \alpha, \beta, \kappa)$ 
3    $\mathcal{Y}_i = g(\mathcal{X}_i), i = 0 : 2n$ 
4    $\mathbf{m} = \sum_{i=0}^{2n} w_i^m \mathcal{Y}_i$ 
5    $\mathbf{S} = \sum_{i=0}^{2n} w_i^c (\mathcal{Y}_i - \mu_U)(\mathcal{Y}_i - \mu_U)^T + \Omega$ 
6    $\mathbf{C} = \sum_{i=0}^{2n} w_i^c (\mathcal{X}_i - \mu)(\mathcal{Y}_i - \mu_U)^T$ 
7   Return  $(\mathbf{m}, \mathbf{S}, \mathbf{C})$ 

```

---

Suppose we have two random variables  $\mathbf{x} \sim \mathcal{N}(\mu, \Sigma)$  and  $\mathbf{y} = g(\mathbf{x})$ , where  $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ . The unscented transform forms a Gaussian approximation to  $p(\mathbf{y})$  using the following process. First we compute a set of  $2n + 1$  sigma points,  $\mathcal{X}_i$ , and corresponding weights,  $w_i^m$  and  $w_i^c$ , using Algorithm 12, for  $i = 0 : 2n$ . (The notation  $\mathbf{M}_{:,i}$  means the  $i$ 'th column of matrix  $\mathbf{M}$ ,  $\sqrt{\Sigma}$  is the matrix square root, so  $\sqrt{\Sigma}\sqrt{\Sigma}^T = \Sigma$ .) Next we propagate the sigma points through the nonlinear function to get the following  $2n + 1$  outputs:

$$\mathcal{Y}_i = g(\mathcal{X}_i), i = 0 : 2n \quad (8.190)$$



37 *Figure 8.12: Illustration of different ways to approximate the distribution induced by a nonlinear transformation*  
38  $f : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ . (a) Data from the source distribution,  $\mathcal{D} = \{\mathbf{x}_i \sim p(\mathbf{x})\}$ , with Gaussian approximation  
39 superimposed. (b) The dots show a Monte Carlo approximation to  $p(f(\mathbf{x}))$  derived from  $\mathcal{D}' = \{f(\mathbf{x}_i)\}$ .  
40 The dotted ellipse is a Gaussian approximation to this target distribution, computed from the empirical  
41 moments. The solid ellipse is Taylor Transform. (c) Unscented sigma points. (d) Unscented transform. (e)  
42 Gauss-Hermite points (order 5). (f) GH transform. Adapted from Figures 5.3–5.4 of [Sar13]. Generated by  
43 gaussian\_transforms.ipynb.

Finally we estimate the mean and covariance of the resulting set of points:

$$\mathbb{E}[\mathbf{g}(\mathbf{x})] \approx \mathbf{m} = \sum_{i=0}^{2n} w_i^m \mathcal{Y}_i \quad (8.191)$$

$$\text{Cov}[\mathbf{g}(\mathbf{x})] \approx \mathbf{S} = \sum_{i=0}^{2n} w_i^c (\mathcal{Y}_i - \boldsymbol{\mu}_U)(\mathcal{Y}_i - \boldsymbol{\mu}_U)^\top \quad (8.192)$$

Now suppose we want to approximate the joint distribution  $p(\mathbf{x}, \mathbf{y})$ , where  $\mathbf{y} = \mathbf{g}(\mathbf{x}) + \mathbf{e}$ , and  $\mathbf{e} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Omega})$ . By defining the augmented function  $\tilde{\mathbf{g}}(\mathbf{x}) = (\mathbf{x}, \mathbf{g}(\mathbf{x}))$ , and applying the above procedure (and adding extra noise), we get

$$\begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} \sim \mathcal{N}\left(\begin{pmatrix} \boldsymbol{\mu} \\ \mathbf{m} \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{S} \end{pmatrix}\right) \quad (8.193)$$

where the parameters are computed using Algorithm 13.

The sigma points and their weights depend on three hyper-parameters,  $\alpha$ ,  $\beta$  and  $\kappa$ , which determine the spread of the sigma points around the mean. A typical recommended setting for these is  $\alpha = 10^{-3}$ ,  $\kappa = 1$ ,  $\beta = 2$  [Bit16].

In Figure 8.12(a-b), we show the linearized Taylor transform discussed in Section 8.5.1 applied to a nonlinear function. In Figure 8.12(c-d), we show the corresponding unscented transform, which we can see is more accurate. In fact, the unscented transform is a third-order method in the sense that the mean of  $\mathbf{y}$  is exact for polynomials up to order 3. However the covariance is only exact for linear functions (first order polynomials). We can compute even more accurate approximations using other forms of numerical integration. For example, we can use **Gauss-Hermite integration** of order  $p$ , which will be exact for polynomials of order up to  $2p - 1$ . See [Sar13, Sec 6.3] for details, and Figure 8.12(e-f) for an illustration. However, this comes at a price: the number of sigma points is now  $p^n$ , whereas it is only  $2n + 1$  for the unscented transform.

### 8.7.2 The unscented Kalman filter (UKF)

The UKF applies the unscented transform twice to in Equation (8.165), once to approximate passing through the system model  $\mathbf{f}$ , and once to approximate passing through the measurement model  $\mathbf{h}$ . By analogy to Section 8.3.2.5, we can derive the UKF algorithm as follows:

$$(\boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}, -) = \text{UnscentedPredict}(\boldsymbol{\mu}_{t-1|t-1}, \boldsymbol{\Sigma}_{t-1|t-1}, \mathbf{f}(\cdot, \mathbf{u}_t), \mathbf{Q}_t) \quad (8.194)$$

$$(\mathbf{m}_t, \mathbf{S}_t, \mathbf{C}_t) = \text{UnscentedPredict}(\boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}, \mathbf{h}(\cdot, \mathbf{u}_t), \mathbf{R}_t) \quad (8.195)$$

$$(\boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t}, \ell_t) = \text{GaussCondition}(\boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}, \mathbf{m}_t, \mathbf{S}_t, \mathbf{C}_t, \mathbf{y}_t) \quad (8.196)$$

See [Sar13, p86] for more details.

### 8.7.3 The unscented Kalman smoother (UKS)

The **unscented Kalman smoother**, also called the **unscented RTS smoother** [Sar08], is a simple modification of the usual Kalman smoothing method, where we approximate the nonlinearity

1 by the unscented transform. By analogy to Section 8.3.3, the update is as follows:  
2

$$\underline{3} \quad (\boldsymbol{\mu}_{t+1|t}, \boldsymbol{\Sigma}_{t+1|t}, \boldsymbol{\Sigma}_{t,t+1|t}) = \text{UnscentedPredict}(\boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t}, \mathbf{f}(\cdot, \mathbf{u}_{t+1}), \mathbf{Q}_{t+1}) \quad (8.197)$$

$$\underline{4} \quad (\boldsymbol{\mu}_{t|T}, \boldsymbol{\Sigma}_{t|T}) = \text{GaussSmoothUpdate}(\boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t}, \boldsymbol{\Sigma}_{t,t+1|t}, \boldsymbol{\mu}_{t+1|T}, \boldsymbol{\Sigma}_{t+1|T}) \quad (8.198)$$

5 See [Sar13, p148] for more details.  
6

### 7 8.7.4 Examples

8 In this section we give some examples of the UKF.

#### 9 8.7.4.1 Tracking a point spiraling in 2d

10 Let us revisit the 2d nonlinear tracking problem from Section 8.5.5.1. In Figure 8.11c, we see that  
11 the UKF algorithm (with  $\alpha = 1$ ,  $\beta = 0$ ,  $\kappa = 2$ ) works well on this problem.  
12

#### 13 8.7.4.2 COVID-19 risk score estimation

14 During the COVID-19 pandemic in 2020, many governments created contact tracing apps, in which  
15 the distance between people was estimated (anonymously) by measuring bluetooth signal strength  
16 between mobile phones. This can be done using the UKF, as explained in [Lov+20]. Indeed, this  
17 approach formed the basis of the UK COVID-19 risk score estimation app [BCH20].<sup>5</sup>  
18

## 19 8.8 Inference based on moment matching

20 In this section, we discuss a simple unified framework, known as the **general Gaussian filter** or  
21 **GGF** [IX00; Wu+06], which includes EKF, UKF, and various other algorithms. Our presentation is  
22 based on [Sar13, Ch. 6].  
23

### 24 8.8.1 Gaussian moment matching

25 To explain the approach, we temporarily drop the time indices, and the conditioning on past  
26 information, and consider a single time step of inference. Furthermore, we will use the shorthand  
27

$$\underline{28} \quad \int_{\mathbf{x}} g(\mathbf{x}) = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} g(x_1, \dots, x_N) dx_1 \cdots dx_N \quad (8.199)$$

29 Let  $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and  $p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{g}(\mathbf{x}), \boldsymbol{\Omega})$  for some function  $\mathbf{g}$ . Let  $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y}|\mathbf{x})$   
30 be the exact joint distribution. The best Gaussian approximation to the joint can be obtained by  
31 **moment matching**, i.e.,  
32

$$\underline{33} \quad q(\mathbf{x}, \mathbf{y}) = \mathcal{N}\left(\begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} \mid \begin{pmatrix} \boldsymbol{\mu} \\ \boldsymbol{\mu}_M \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma} & \mathbf{C}_M \\ \mathbf{C}_M^\top & \mathbf{S}_M \end{pmatrix}\right) \quad (8.200)$$

34 35 5. Once the distance has been estimated, it needs to be combined with other signals, such as contact duration and  
36 infectiousness level of the index case, to estimate the risk of transmission. For an ML based approach to this problem,  
37 see [MKS21].  
38

1 where

2

$$\mu_M = \int_{\mathbf{x}} \mathbf{g}(\mathbf{x}) \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (8.201)$$

3

$$\mathbf{S}_M = \int_{\mathbf{x}} (\mathbf{g}(\mathbf{x}) - \mu_M)(\mathbf{g}(\mathbf{x}) - \mu_M)^T \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) + \boldsymbol{\Omega} \quad (8.202)$$

4

$$\mathbf{C}_M = \int_{\mathbf{x}} (\mathbf{x} - \boldsymbol{\mu})(\mathbf{g}(\mathbf{x}) - \mu_M)^T \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (8.203)$$

5 We can either compute these integrals by linearizing  $\mathbf{g}$  and using closed-form expressions, or by using  
6 numerical integration, as we discuss in Section 8.8.3.

7 Once we have computed the joint  $q(\mathbf{x}, \mathbf{y})$ , we can compute the marginal  $q(\mathbf{y})$  and the posterior  
8  $q(\mathbf{x} | \mathbf{y})$  using the usual equations for Gaussian distributions.

### 9 8.8.2 The general Gaussian filter

10 Let us now apply this method in the filtering context. We assume the prior has the form  
11  $p(\mathbf{z}_{t-1} | \mathbf{y}_{1:t-1}) = \mathcal{N}(\boldsymbol{\mu}_{t-1|t-1}, \boldsymbol{\Sigma}_{t-1|t-1})$ . In the prediction step, we compute a Gaussian approxima-  
12 tion to  $p(\mathbf{z}_{t-1}, \mathbf{z}_t | \mathbf{y}_{1:t-1})$ , from which we can compute the marginal  $p(\mathbf{z}_t | \mathbf{y}_{1:t-1}) \approx \mathcal{N}(\boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1})$ .  
13 In the update step, we compute a Gaussian approximation to  $p(\mathbf{z}_t, \mathbf{y}_t | \mathbf{y}_{1:t-1})$ , from which we can  
14 compute the posterior  $p(\mathbf{z}_t | \mathbf{y}_t, \mathbf{y}_{1:t-1}) \approx \mathcal{N}(\boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t})$ .

### 15 8.8.3 Implementation using numerical integration

16 To implement the above integrals in practice, we usually need some approximations. One approach  
17 is to linearize  $\mathbf{g}$  around the predicted mean, as in Section 8.5.1. This gives rise to the EKF. However,  
18 this approach can lead to large errors and sometimes divergence of the filter [IX00; VDMW03].

19 A more stable method is to use **numerical integration**<sup>6</sup>, which can be written as

20

$$\int_{\mathbf{z}} f(\mathbf{z}) \approx \sum_{k=1}^K w^k f(\mathbf{z}^k) \quad (8.204)$$

21 for a suitable set of evaluation points  $\mathbf{z}^k$  and weights  $w^k$ .

22 One approach is to use the unscented transform, as in Algorithm 13. We can denote this by

23

$$(\boldsymbol{\mu}_M, \mathbf{S}_M, \mathbf{C}_M) = \text{UnscentedPredict}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{g}, \boldsymbol{\Omega}) \quad (8.205)$$

24 Using this inside the GGF is equivalent to the UKF in Section 8.7.2. Alternatively, we can use  
25 **spherical cubature integration**, which gives rise to the **cubature Kalman filter** or **CKF** [AH09].  
26 This turns out (see [Sar13, p110]) to be a special case of the UKF, with  $2n + 1$  sigma points, and  
27 hyper-parameters values of  $\alpha = 1$  and  $\beta = 0$  (with  $\kappa$  left free).

28 A more accurate approximation uses **Gauss-Hermite integration**, which allows the user to  
29 select more sigma points, as we discussed in Line 7. This gives rise to the **Gauss-Hermite Kalman**  
30 **filter** or **GHKF** [IX00], also known as the **quadrature Kalman filter** or **QKF** [AHE07].

31 6. One-dimensional integrals are called **quadratures**, and multi-dimensional integrals are called **cubatures**.

<sup>1</sup> We can also approximate the integrals with Monte Carlo. Note, however, that this is not the same  
<sup>2</sup> as particle filtering (Section 13.2), which approximates the conditional  $p(\mathbf{z}_t|\mathbf{y}_{1:t})$  rather than the  
<sup>3</sup> joint  $p(\mathbf{z}_t, \mathbf{y}_t|\mathbf{y}_{1:t-1})$  (see Section 8.10.1 for discussion of this difference).  
<sup>4</sup>

## <sup>5</sup> 8.9 Inference based on statistical linearization

<sup>6</sup> In this section, we discuss an approach to approximate inference which is based on generalized  
<sup>7</sup> Gaussian filtering (Section 8.8), but which chooses the parameters of the linear approximation to  
<sup>8</sup> minimize the mean squared error. This approach is therefore called **statistical linear regression**  
<sup>9</sup> or **SLR** [LBS01; AHE07]. SLR is more accurate than Taylor linearization, which is only optimal at  
<sup>10</sup> a single point, since it takes the entire distribution over the hidden quantities into account. We can  
<sup>11</sup> show that the UKF is a special case of SLR. However, SLR is more general, and can be extended to  
<sup>12</sup> handle non-Gaussian noise (see Section 8.9.5).  
<sup>13</sup>

### <sup>14</sup> 8.9.1 Statistical linear regression

<sup>15</sup> In this section, we describe the SLR method, following [GFSS17]. We drop the time index for brevity.  
<sup>16</sup> Suppose we have the prior  $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and likelihood  $p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{g}(\mathbf{x}), \boldsymbol{\Omega})$ . We want to  
<sup>17</sup> approximate the joint,  $p(\mathbf{x}, \mathbf{y})$ , by a Gaussian, from which we can compute the marginal  $p(\mathbf{y})$  and  
<sup>18</sup> the posterior  $p(\mathbf{x}|\mathbf{y})$  in the usual way.

<sup>19</sup> We use the following enabling approximation:

$$\mathbf{g}(\mathbf{x}) \approx \mathbf{Ax} + \mathbf{b} + \mathbf{e} \quad (8.206)$$

<sup>20</sup> where  $\mathbf{e} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Omega})$  is the additional noise term introduced by the approximation. This will induce  
<sup>21</sup> the joint distribution  $q(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta})$ , where  $\boldsymbol{\theta} = (\mathbf{A}, \mathbf{b}, \boldsymbol{\Omega})$ . We choose the parameters of the approximate  
<sup>22</sup> model so that the induced joint is as close as possible to the true joint:  
<sup>23</sup>

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} D_{\text{KL}}(p(\mathbf{x}, \mathbf{y}) \parallel q(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta})) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \int_{\mathbf{x}, \mathbf{y}} p(\mathbf{x}, \mathbf{y}) \log q(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta}) \quad (8.207)$$

<sup>24</sup> This can be solved by making the moments of  $q$  match the moments of  $p$  (see Section 5.1.3.4).  
<sup>25</sup>

<sup>26</sup> One can also show that this objective is equivalent to minimizing the mean squared error of the  
<sup>27</sup> approximation  
<sup>28</sup>

$$(\hat{\mathbf{A}}, \hat{\mathbf{b}}) = \underset{\mathbf{A}, \mathbf{b}}{\operatorname{argmin}} \mathbb{E}[(\mathbf{g}(\mathbf{x}) - \mathbf{Ax} - \mathbf{b})^\top (\mathbf{g}(\mathbf{x}) - \mathbf{Ax} - \mathbf{b})] \quad (8.208)$$

<sup>29</sup> where  $\hat{\boldsymbol{\Omega}}$  is the covariance of the resulting errors:  
<sup>30</sup>

$$\hat{\boldsymbol{\Omega}} = \mathbb{E}[(\mathbf{g}(\mathbf{x}) - \hat{\mathbf{Ax}} - \hat{\mathbf{b}})(\mathbf{g}(\mathbf{x}) - \hat{\mathbf{Ax}} - \hat{\mathbf{b}})^\top] \quad (8.209)$$

<sup>31</sup> One can show that the resulting solution is given by  
<sup>32</sup>

$$\hat{\mathbf{A}} = \mathbf{C}^\top \boldsymbol{\Sigma}^{-1} \quad (8.210a)$$

$$\hat{\mathbf{b}} = \mathbf{m} - \hat{\mathbf{A}}\boldsymbol{\mu} \quad (8.210b)$$

$$\hat{\boldsymbol{\Omega}} = \mathbf{S} - \hat{\mathbf{A}}\boldsymbol{\Sigma}\hat{\mathbf{A}}^\top \quad (8.210c)$$