

Figure 15.3: (a) Representing lasso using a Gaussian scale mixture prior. (b) Graphical model for group lasso with 2 groups, the first has size $G_1 = 2$, the second has size $G_2 = 3$.

Unfortunately, MAP estimation (not to mention full Bayesian inference) with such discrete mixture priors is computationally difficult. Various approximate inference methods have been proposed, including greedy search (see e.g., [SPZ09]) or MCMC (see e.g., [HS09]).

15.2.5 Laplace prior (Bayesian lasso)

A computationally cheap way to achieve sparsity is to perform MAP estimation with a Laplace prior by minimizing the penalized negative log likelihood:

$$\text{PNLL}(\mathbf{w}) = -\log p(\mathcal{D}|\mathbf{w}) - \log p(\mathbf{w}|\lambda) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda\|\mathbf{w}\|_1 \quad (15.84)$$

where $\|\mathbf{w}\|_1 \triangleq \sum_{d=1}^D |w_d|$ is the ℓ_1 norm of \mathbf{w} . This method is called **lasso**, which stands for “least absolute shrinkage and selection operator” [Tib96]. See Section 11.4 of the prequel to this book, [Mur22], for details.

In this section, we discuss posterior inference with this prior; this is known as the **Bayesian lasso** [PC08]. In particular, we assume the following prior:

$$p(\mathbf{w}|\sigma^2) = \prod_j \frac{\lambda}{2\sqrt{\sigma^2}} e^{-\lambda|w_j|/\sqrt{\sigma^2}} \quad (15.85)$$

(Note that conditioning the prior on σ^2 is important to ensure that the full posterior is unimodal.)

To simplify inference, we will represent the Laplace prior as a Gaussian scale mixture, which we discussed in Section 28.2.3.2. In particular, one can show that the Laplace distribution is an infinite

1 weighted sum of Gaussians, where the precision comes from a Gamma distribution:
2

3

$$\text{Laplace}(w|0, \lambda) = \int \mathcal{N}(w|0, \tau^2) \text{Ga}(\tau^2|1, \frac{\lambda^2}{2}) d\tau^2 \quad (15.86)$$

4

5 We can therefore represent the Bayesian lasso model as a hierarchical latent variable model, as shown
6 in Figure 15.3a. The corresponding joint distribution has the following form:
7

8

$$p(\mathbf{y}, \mathbf{w}, \boldsymbol{\tau}, \sigma^2 | \mathbf{X}) = \mathcal{N}(\mathbf{y} | \mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I}_N) \left[\prod_j \mathcal{N}(w_j | 0, \sigma^2 \tau_j^2) \text{Ga}(\tau_j^2 | 1, \lambda^2 / 2) \right] p(\sigma^2) \quad (15.87)$$

9

10 We can also create a GSM to match the **group lasso** prior, which sets multiple coefficients to
11 zero at the same time:
12

13

$$\mathbf{w}_g | \sigma^2, \tau_g^2 \sim \mathcal{N}(\mathbf{0}, \sigma^2 \tau_g^2 \mathbf{I}_{d_g}) \quad (15.88)$$

14

15

$$\tau_g^2 \sim \text{Ga}(\frac{d_g + 1}{2}, \frac{\lambda^2}{2}) \quad (15.89)$$

16

17 where d_g is the size of group g . So we see that there is one variance term per group, each of which
18 comes from a Gamma prior, whose shape parameter depends on the group size, and whose rate
19 parameter is controlled by γ .

20 Figure 15.3b gives an example, where we have 2 groups, one of size 2 and one of size 3. This picture
21 makes it clearer why there should be a grouping effect. For example, suppose $w_{1,1}$ is small; then τ_1^2
22 will be estimated to be small, which will force $w_{1,2}$ to be small, due to shrinkage (c.f., Section 3.7).
23 Conversely, suppose $w_{1,1}$ is large; then τ_1^2 will be estimated to be large, which will allow $w_{1,2}$ to be
24 become large as well.

25 Given these hierarchical models, we can easily derive a Gibbs sampling algorithm (Section 12.3) to
26 sample from the posterior (see e.g., [PC08]). Unfortunately, these posterior samples are not sparse,
27 even though the MAP estimate is sparse. This is because the prior puts infinitesimal probability on
28 the event that each coefficient is zero.
29

30 15.2.6 Horseshoe prior

31

32 The Laplace prior is not suitable for sparse Bayesian models, because posterior samples are not
33 sparse. The spike and slab prior does not have this problem but is often too slow to use (although see
34 [BRG20]). Fortunately, it is possible to devise continuous priors (without discrete latent variables)
35 that are both sparse and computationally efficient. One popular prior of this type is the **horseshoe**
36 prior [CPS10], so-named because of the shape of its density function.

37 In the horseshoe prior, instead of using a Laplace prior for each weight, we use the following
38 Gaussian scale mixture:
39

40

$$w_j \sim \mathcal{N}(0, \lambda_j^2 \tau^2) \quad (15.90)$$

41

42

$$\lambda_j \sim \mathcal{C}_+(0, 1) \quad (15.91)$$

43

44

$$\tau^2 \sim \mathcal{C}_+(0, 1) \quad (15.92)$$

45

where $\mathcal{C}_+(0, 1)$ is the half-Cauchy distribution (Section 2.2.2.4), λ_j is a local shrinkage factor, and τ^2 is a global shrinkage factor. The Cauchy distribution has very fat tails, so λ_j is likely to be either 0 or very far from 0, which emulates the spike and slab prior, but in a continuous way. For more details, see e.g., [Bha+19].

15.2.7 Automatic relevancy determination

An alternative to using posterior inference with a sparsity promoting prior is to use posterior inference with a Gaussian prior, $w_j \sim \mathcal{N}(0, 1/\alpha_j)$, but where we use empirical Bayes to optimize the precisions α_j . That is, we first compute $\hat{\boldsymbol{\alpha}} = \text{argmax}_{\boldsymbol{\alpha}} p(\mathbf{y}|\mathbf{X}, \boldsymbol{\alpha})$, and then compute $\hat{\mathbf{w}} = \text{argmax}_{\mathbf{w}} \mathcal{N}(\mathbf{w}|\mathbf{0}, \hat{\boldsymbol{\alpha}}^{-1})$. Perhaps surprisingly, we will see that this results in a sparse estimate, for reasons we explain in Section 15.2.7.2.

This technique is known as **sparse Bayesian learning** [Tip01] or **automatic relevancy determination (ARD)** [Mac95; Nea96]. It has also been called **NUV** estimation, which stands for “normal prior with unknown variance” [Loe+16]. It was originally developed for neural networks (where sparsity is applied to the first layer weights), but here we apply it to linear models.

15.2.7.1 ARD for linear models

In this section, we explain ARD in more detail, by applying it to linear regression. The likelihood is $p(\mathbf{y}|\mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(\mathbf{y}|\mathbf{w}^\top \mathbf{x}, 1/\beta)$, where $\beta = 1/\sigma^2$. The prior is $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \mathbf{A}^{-1})$, where $\mathbf{A} = \text{diag}(\boldsymbol{\alpha})$. The marginal likelihood can be computed analytically (using Equation (2.90)) as follows:

$$p(\mathbf{y}|\mathbf{X}, \boldsymbol{\alpha}, \beta) = \int \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, (1/\beta)\mathbf{I}_{N_D}) \mathcal{N}(\mathbf{w}|\mathbf{0}, \mathbf{A}^{-1}) d\mathbf{w} \quad (15.93)$$

$$= \mathcal{N}(\mathbf{y}|\mathbf{0}, \beta^{-1}\mathbf{I}_{N_D} + \mathbf{X}\mathbf{A}^{-1}\mathbf{X}^\top) \quad (15.94)$$

$$= \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{C}_{\boldsymbol{\alpha}}) \quad (15.95)$$

where $\mathbf{C}_{\boldsymbol{\alpha}} \triangleq \beta^{-1}\mathbf{I}_N + \mathbf{X}\mathbf{A}^{-1}\mathbf{X}^\top$. This is very similar to the marginal likelihood under the spike-and-slab prior (Section 15.2.4), which is given by

$$p(\mathbf{y}|\mathbf{X}, \mathbf{s}, \sigma_w^2, \sigma_y^2) = \int \mathcal{N}(\mathbf{y}|\mathbf{X}_s \mathbf{w}_s, \sigma_y^2 \mathbf{I}) \mathcal{N}(\mathbf{w}_s|\mathbf{0}_s, \sigma_w^2 \mathbf{I}) d\mathbf{w}_s = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{C}_s) \quad (15.96)$$

where $\mathbf{C}_s = \sigma_y^2 \mathbf{I}_{N_D} + \sigma_w^2 \mathbf{X}_s \mathbf{X}_s^\top$. (Here \mathbf{X}_s refers to the design matrix where we select only the columns of \mathbf{X} where $s_d = 1$.) The difference is that we have replaced the binary $s_j \in \{0, 1\}$ variables with continuous $\alpha_j \in \mathbb{R}^+$, which makes the optimization problem easier.

The objective is the log marginal likelihood, given by

$$\ell(\boldsymbol{\alpha}, \beta) = -\frac{1}{2} \log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\alpha}, \beta) = \log |\mathbf{C}_{\boldsymbol{\alpha}}| + \mathbf{y}^\top \mathbf{C}_{\boldsymbol{\alpha}}^{-1} \mathbf{y} \quad (15.97)$$

There are various algorithms for optimizing $\ell(\boldsymbol{\alpha}, \beta)$, some of which we discuss in Section 15.2.7.3.

ARD can be used as an alternative to ℓ_1 regularization. Although the ARD objective is not convex, it tends to give much sparser results [WW12]. In addition, it can be shown [WRN10] that the ARD objective has many fewer local optima than the ℓ_0 -regularized objective, and hence is much easier to optimize.

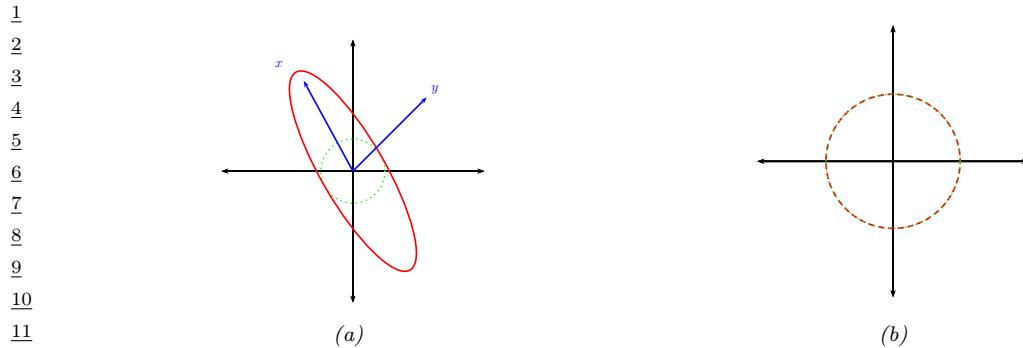


Figure 15.4: Illustration of why ARD results in sparsity. The vector of inputs \mathbf{x} does not point towards the vector of outputs \mathbf{y} , so the feature should be removed. (a) For finite α , the probability density is spread in directions away from \mathbf{y} . (b) When $\alpha = \infty$, the probability density at \mathbf{y} is maximized. Adapted from Figure 8 of [Tip01].

15.2.7.2 Why does ARD result in a sparse solution?

Once we have estimated $\boldsymbol{\alpha}$ and β , we can compute the posterior over the parameters using Bayes rule for Gaussians, to get $p(\mathbf{w}|\mathcal{D}, \hat{\boldsymbol{\alpha}}, \hat{\beta}) = \mathcal{N}(\mathbf{w}|\hat{\mathbf{w}}, \hat{\Sigma})$, where $\hat{\Sigma}^{-1} = \hat{\beta}\mathbf{X}^T\mathbf{X} + \mathbf{A}$ and $\hat{\mathbf{w}} = \hat{\beta}\hat{\Sigma}\mathbf{X}^T\mathbf{y}$. If we have $\hat{\alpha}_d \approx \infty$, then $\hat{\mathbf{w}}_d \approx 0$, so the solution vector will be sparse.

We now give an intuitive argument, based on [Tip01], about when such a sparse solution may be optimal. We shall assume $\beta = 1/\sigma^2$ is fixed for simplicity. Consider a 1d linear regression with 2 training examples, so $\mathbf{X} = \mathbf{x} = (x_1, x_2)$, and $\mathbf{y} = (y_1, y_2)$. We can plot \mathbf{x} and \mathbf{y} as vectors in the plane, as shown in Figure 15.4. Suppose the feature is irrelevant for predicting the response, so \mathbf{x} points in a nearly orthogonal direction to \mathbf{y} . Let us see what happens to the marginal likelihood as we change α . The marginal likelihood is given by $p(\mathbf{y}|\mathbf{x}, \alpha, \beta) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{C}_\alpha)$, where $\mathbf{C}_\alpha = \frac{1}{\beta}\mathbf{I} + \frac{1}{\alpha}\mathbf{x}\mathbf{x}^T$. If α is finite, the posterior will be elongated along the direction of \mathbf{x} , as in Figure 15.4(a). However, if $\alpha = \infty$, we have $\mathbf{C}_\alpha = \frac{1}{\beta}\mathbf{I}$, which is spherical, as in Figure 15.4(b). If $|\mathbf{C}_\alpha|$ is held constant, the latter assigns higher probability density to the observed response vector \mathbf{y} , so this is the preferred solution. In other words, the marginal likelihood “punishes” solutions where α_d is small but $\mathbf{X}_{:,d}$ is irrelevant, since these waste probability mass. It is more parsimonious (from the point of view of Bayesian Occam’s razor) to eliminate redundant dimensions.

Another way to understand the sparsity properties of ARD is as approximate inference in a hierarchical Bayesian model [BT00]. In particular, suppose we put a conjugate prior on each precision, $\alpha_d \sim \text{Ga}(a, b)$, and on the observation precision, $\beta \sim \text{Ga}(c, d)$. Since exact inference with a Student prior is intractable, we can use variational Bayes (Section 10.2.3), with a factored posterior approximation of the form

$$q(\mathbf{w}, \boldsymbol{\alpha}) = q(\mathbf{w})q(\boldsymbol{\alpha}) \approx \mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \prod_d \text{Ga}(\alpha_d | \hat{\alpha}_d, \hat{b}_d) \quad (15.98)$$

ARD approximates $q(\boldsymbol{\alpha})$ by a point estimate. However, in VB, we integrate out $\boldsymbol{\alpha}$; the resulting

1 posterior marginal $q(\mathbf{w})$ on the weights is given by
2

$$\small \begin{array}{l} 3 \\ 4 \quad p(\mathbf{w}|\mathcal{D}) = \int \mathcal{N}(\mathbf{w}|\mathbf{0}, \text{diag}(\boldsymbol{\alpha})^{-1}) \prod_d \text{Ga}(\alpha_d | \hat{a}_d, \hat{b}) d\boldsymbol{\alpha} \\ 5 \end{array} \quad (15.99)$$

6 This is a Gaussian scale mixture, and can be shown to be the same as a multivariate Student
7 distribution (see Section 28.2.3.1), with non-diagonal covariance. Note that the Student has a large
8 spike at 0, which intuitively explains why the posterior mean (which, for a Student distribution, is
9 equal to the posterior mode) is sparse.

10 Finally, we can also view ARD as a MAP estimation problem with a **non-factorial prior** [WN07].
11 Intuitively, the dependence between the w_j parameters arises, despite the use of a diagonal Gaussian
12 prior, because the prior precision α_j is estimated based after marginalizing out all \mathbf{w} , and hence
13 depends on all the features. Interestingly, [WRN10] prove that MAP estimation with non-factorial
14 priors is strictly better than MAP estimation with any possible factorial prior in the following
15 sense: the non-factorial objective always has fewer local minima than factorial objectives, while still
16 satisfying the property that the global optimum of the non-factorial objective corresponds to the
17 global optimum of the ℓ_0 objective — a property that ℓ_1 regularization, which has no local minima,
18 does not enjoy.
19

20 15.2.7.3 Algorithms for ARD

22 There are various algorithms for optimizing $\ell(\boldsymbol{\alpha}, \beta)$. One approach is to use EM, in which we compute
23 $p(\mathbf{w}|\mathcal{D}, \boldsymbol{\alpha})$ in the E step and then maximize $\boldsymbol{\alpha}$ in the M step. In variational Bayes, we infer both \mathbf{w}
24 and $\boldsymbol{\alpha}$ (see [Dru08] for details). In [WN10], they present a method based on iteratively reweighted ℓ_1
25 estimation.

26 Recently, [HXW17] showed that the nested iterative computations performed these methods can
27 emulated by a recurrent neural network (Section 16.3.4). Furthermore, by training this model, it is
28 possible to achieve much faster convergence than manually designed optimization algorithms.
29

30 15.2.7.4 Relevance vector machines

32 Suppose we create a linear regression model of the form $p(y|\mathbf{x}; \boldsymbol{\theta}) = \mathcal{N}(y|\mathbf{w}^\top \phi(\mathbf{x}), \sigma^2)$, where
33 $\phi(\mathbf{x}) = [\mathcal{K}(\mathbf{x}, \mathbf{x}_1), \dots, \mathcal{K}(\mathbf{x}, \mathbf{x}_N)]$, where $\mathcal{K}()$ is a kernel function (Section 18.2) and $\mathbf{x}_1, \dots, \mathbf{x}_N$ are
34 the N training points. This is called **kernel basis function expansion**, and transforms the input
35 from $\mathbf{x} \in \mathcal{X}$ to $\phi(\mathbf{x}) \in \mathbb{R}^N$. Obviously this model has $O(N)$ parameters, and hence is nonparametric.
36 However, we can use ARD to select a small subset of the exemplars. This technique is called the
37 relevance vector machine (RVM) [Tip01; TF03].
38

39 15.2.8 Multivariate linear regression

40 This section is written by Xinglong Li.
41

42 In this section, we consider the **multivariate linear regression** model, which has the form
43

$$\small \begin{array}{l} 44 \\ 45 \quad \mathbf{Y} = \mathbf{W}\mathbf{X} + \mathbf{E} \\ 46 \end{array} \quad (15.100)$$

45 where $\mathbf{W} \in \mathbb{R}^{N_y \times N_x}$ is the matrix of regression coefficient, $\mathbf{X} \in \mathbb{R}^{N_x \times N_D}$ is the matrix of input
46 features (with each row being an input variable and each column being an observation), $\mathbf{Y} \in \mathbb{R}^{N_y \times N_D}$
47

¹ is the matrix of responses (with each row being an output variable and each *column* being an
² observation), and $\mathbf{E} = [e_1, \dots, e_{N_D}]$ is the matrix of residual errors, where $e_i \stackrel{i.i.d.}{\sim} \mathcal{N}(\mathbf{0}, \Sigma)$. It can be
³ seen from the definition that given Σ , \mathbf{W} and \mathbf{X} , columns of \mathbf{Y} are independently random variables
⁴ following multivariate normal distributions. So the likelihood of the observation is
⁵

$$\frac{7}{8} p(\mathbf{Y}|\mathbf{W}, \mathbf{X}, \Sigma) = \frac{1}{(2\pi)^{N_y \times N_D} |\Sigma|^{N_D/2}} \exp \left(\sum_{i=1}^{N_D} -\frac{1}{2} (\mathbf{y}_i - \mathbf{W}\mathbf{x}_i)^\top \Sigma^{-1} (\mathbf{y}_i - \mathbf{W}\mathbf{x}_i) \right) \quad (15.101)$$

$$\frac{10}{11} = \frac{1}{(2\pi)^{N_y \times N_D} |\Sigma|^{N_D/2}} \exp \left(-\frac{1}{2} \text{trace} ((\mathbf{Y} - \mathbf{W}\mathbf{X})^\top \Sigma^{-1} (\mathbf{Y} - \mathbf{W}\mathbf{X})) \right) \quad (15.102)$$

$$\frac{12}{13} = \mathcal{MN}(\mathbf{Y}|\mathbf{W}\mathbf{X}, \Sigma, \mathbf{I}_{N_D \times N_D}), \quad (15.103)$$

¹⁴ The conjugate prior for this is the **matrix normal inverse Wishart** distribution,

$$\frac{15}{16} \mathbf{W}, \Sigma \sim \text{MNIW}(\mathbf{M}_0, \mathbf{V}_0, \nu_0, \Psi_0) \quad (15.104)$$

¹⁷ where the MNIW is defined by

$$\frac{19}{20} \mathbf{W}|\Sigma \sim \mathcal{MN}(\mathbf{M}_0, \Sigma_0, \mathbf{V}_0) \quad (15.105)$$

$$\frac{21}{22} \Sigma \sim \text{IW}(\nu_0, \Psi_0), \quad (15.106)$$

²² where $\mathbf{V}_0 \in \mathbb{R}_{++}^{N_x \times N_x}$, $\Psi_0 \in \mathbb{R}_{++}^{N_y \times N_y}$ and $\nu_0 > N_x - 1$ is the degree of freedom of the inverse Wishart
²³ distribution.

²⁴ The posterior distribution of $\{\mathbf{W}, \Sigma\}$ still follows a matrix normal inverse Wishart distribution.
²⁵ We follow the derivation in [Fox09, App.F]. Firstly, the density of the joint distribution is

$$\frac{28}{29} p(\mathbf{Y}, \mathbf{W}, \Sigma) \propto |\Sigma|^{-(\nu_0 + N_y + 1 + N_x + N_D)/2} \times \exp \left\{ -\frac{1}{2} \text{trace}(\Omega_0) \right\} \quad (15.107)$$

$$\frac{30}{31} \Omega_0 \triangleq \Psi_0 \Sigma^{-1} + (\mathbf{Y} - \mathbf{W}\mathbf{X})^\top \Sigma^{-1} (\mathbf{Y} - \mathbf{W}\mathbf{X}) + (\mathbf{W} - \mathbf{M}_0)^\top \Sigma^{-1} (\mathbf{W} - \mathbf{M}_0) \mathbf{V}_0 \quad (15.108)$$

³² We firstly aggregate items including \mathbf{W} in the exponent so that it takes the form of a matrix normal
³³ distribution. This is similar to the “completing the square” technique that we used in deriving the
³⁴ conjugate posterior for multivariate normal distributions in Section 3.3.3.3. Specifically,

$$\frac{36}{37} \text{trace} [(\mathbf{Y} - \mathbf{W}\mathbf{X})^\top \Sigma^{-1} (\mathbf{Y} - \mathbf{W}\mathbf{X}) + (\mathbf{W} - \mathbf{M}_0)^\top \Sigma^{-1} (\mathbf{W} - \mathbf{M}_0) \mathbf{V}_0] \quad (15.109)$$

$$\frac{38}{39} = \text{trace} (\Sigma^{-1} [(\mathbf{Y} - \mathbf{W}\mathbf{X})(\mathbf{Y} - \mathbf{W}\mathbf{X})^\top + (\mathbf{W} - \mathbf{M}_0) \mathbf{V}_0 (\mathbf{W} - \mathbf{M}_0)^\top]) \quad (15.110)$$

$$\frac{40}{41} = \text{trace} (\Sigma^{-1} [\mathbf{W}\mathbf{S}_{xx}\mathbf{W}^\top - 2\mathbf{S}_{yx}\mathbf{W}^\top + \mathbf{S}_{yy}]) \quad (15.111)$$

$$\frac{42}{43} = \text{trace} (\Sigma^{-1} [(\mathbf{W} - \mathbf{S}_{yx}\mathbf{S}_{xx}^{-1})\mathbf{S}_{xx}(\mathbf{W} - \mathbf{S}_{yx}\mathbf{S}_{xx}^{-1})^\top + \mathbf{S}_{y|x}]). \quad (15.112)$$

⁴² where

$$\frac{44}{45} \mathbf{S}_{xx} = \mathbf{X}\mathbf{X}^\top + \mathbf{V}_0, \quad \mathbf{S}_{yx} = \mathbf{Y}\mathbf{X}^\top + \mathbf{M}_0 \mathbf{V}_0, \quad (15.113)$$

$$\frac{46}{47} \mathbf{S}_{yy} = \mathbf{Y}\mathbf{Y}^\top + \mathbf{M}_0 \mathbf{V}_0 \mathbf{M}_0^\top, \quad \mathbf{S}_{y|x} = \mathbf{S}_{yy} - \mathbf{S}_{yx}\mathbf{S}_{xx}^{-1}\mathbf{S}_{yx}^\top. \quad (15.114)$$

Therefore, it can be seen from Equation (15.112) that given Σ, \mathbf{W} follows a matrix normal distribution

$$\mathbf{W} | \Sigma, \mathbf{X}, \mathbf{Y} \sim \mathcal{MN}(\mathbf{S}_{yx} \mathbf{S}_{xx}^{-1}, \Sigma, \mathbf{S}_{xx}). \quad (15.115)$$

Marginalizing out \mathbf{W} (which corresponds to removing the terms including \mathbf{W} in the exponent in Equation (15.108)), it can be shown that the posterior distribution of Σ is an inverse Wishart distribution. In fact, by replacing Equation (15.112) to the corresponding terms in Equation (15.108), it can be seen that the only terms left after integrating out \mathbf{W} are $\Sigma^{-1}\Psi$ and $\Sigma^{-1}\mathbf{S}_{y|x}$, which indicates that the scale matrix of the posterior inverse Wishart distribution is $\Psi_0 + \mathbf{S}_{y|x}$.

In conclusion, the joint posterior distribution of $\{\mathbf{W}, \Sigma\}$ given the observation is

$$\mathbf{W}, \Sigma | \mathbf{X}, \mathbf{Y} \sim \text{MNIW}(\mathbf{M}_1, \mathbf{V}_1, \nu_1, \Psi_1) \quad (15.116)$$

$$\mathbf{M}_1 = \mathbf{S}_{yx} \mathbf{S}_{xx}^{-1} \quad (15.117)$$

$$\mathbf{V}_1 = \mathbf{S}_{xx} \quad (15.118)$$

$$\nu_1 = N_{\mathcal{D}} + \nu_0 \quad (15.119)$$

$$\Psi_1 = \Psi_0 + \mathbf{S}_{y|x} \quad (15.120)$$

The MAP estimate of \mathbf{W} and Σ are the mode of the posterior matrix normal inverse Wishart distribution. To derive this, notice that \mathbf{W} only appears in the matrix normal density function in the posterior, so the matrix \mathbf{W} maximizing the posterior density of $\{\mathbf{W}, \Sigma\}$ is the matrix $\hat{\mathbf{W}}$ that maximizes the matrix normal posterior of \mathbf{W} . So the MAP estimate of \mathbf{W} is $\hat{\mathbf{W}} = \mathbf{M}_1 = \mathbf{S}_{yx} \mathbf{S}_{xx}^{-1}$, and this holds for any value of Σ . By plugging $\mathbf{W} = \hat{\mathbf{W}}$ into the joint posterior of $\{\mathbf{W}, \Sigma\}$, and taking derivatives over Σ , it can be seen that the matrix maximizing the density is $(\nu_1 + N_y + N_x + 1)^{-1} \Psi_1$. Since Ψ_1 is positive definite, it is the MAP estimate of Σ .

In conclusion, the MAP estimate of $\{\mathbf{W}, \Sigma\}$ are

$$\hat{\mathbf{W}} = \mathbf{S}_{yx} \mathbf{S}_{xx}^{-1} \quad (15.121)$$

$$\hat{\Sigma} = \frac{1}{\nu_1 + N_y + N_x + 1} (\Psi_0 + \mathbf{S}_{y|x}) \quad (15.122)$$

15.3 Logistic regression

Logistic regression is a very widely used discriminative classification model that maps input vectors $\mathbf{x} \in \mathbb{R}^D$ to a distribution over class labels, $y \in \{1, \dots, C\}$. If $C = 2$, this is known as **binary logistic regression**, and if $C > 2$, it is known as **multinomial logistic regression**, or alternatively, **multiclass logistic regression**.

15.3.1 Binary logistic regression

In the binary case, where $y \in \{0, 1\}$, the model has the following form

$$p(y|\mathbf{x}; \boldsymbol{\theta}) = \text{Ber}(y|\sigma(\mathbf{w}^\top \mathbf{x} + b)) \quad (15.123)$$

where \mathbf{w} are the weights, b is the bias (offset), and σ is the **sigmoid** or **logistic** function, defined by

$$\sigma(a) \triangleq \frac{1}{1 + e^{-a}} \quad (15.124)$$

1 Let $\eta_n = \mathbf{w}^\top \mathbf{x}_n + b$ be the **logits** for example n , and $\mu_n = \sigma(\eta_n) = p(y = 1|\mathbf{x}_n)$ be the mean of
2 the output. Then we can write the log likelihood as the negative cross entropy:
3

4

$$\log p(\mathcal{D}|\boldsymbol{\theta}) = \log \prod_{n=1}^N \mu_n^{y_n} (1 - \mu_n)^{1-y_n} = \sum_{n=1}^N y_n \log \mu_n + (1 - y_n) \log(1 - \mu_n) \quad (15.125)$$

5

6 We can expand this equation into a more explicit form (that is commonly seen in implementations)
7 by performing some simple algebra. First note that
8

9

$$\mu_n = \frac{1}{1 + e^{-\eta_n}} = \frac{e^{\eta_n}}{1 + e^{\eta_n}}, \quad 1 - \mu_n = 1 - \frac{e^{\eta_n}}{1 + e^{\eta_n}} = \frac{1}{1 + e^{\eta_n}} \quad (15.126)$$

10

11 Hence

12

$$\log p(\mathcal{D}|\boldsymbol{\theta}) = \sum_{n=1}^N y_n [\log e^{\eta_n} - \log(1 + e^{\eta_n})] + (1 - y_n) [\log 1 - \log(1 + e^{\eta_n})] \quad (15.127)$$

13

14

$$= \sum_{n=1}^N y_n [\eta_n - \log(1 + e^{\eta_n})] + (1 - y_n) [-\log(1 + e^{\eta_n})] \quad (15.128)$$

15

16

$$= \sum_{n=1}^N y_n \eta_n - \sum_{n=1}^N \log(1 + e^{\eta_n}) \quad (15.129)$$

17

18 Note that the $\log(1 + e^a)$ function is often implemented using `np.log1p(np.exp(a))`.
19

20 15.3.2 Multinomial logistic regression

21 **Multinomial logistic regression** is a discriminative classification model of the following form:
22

23

$$p(y|\mathbf{x}; \boldsymbol{\theta}) = \text{Cat}(y|\text{softmax}(\mathbf{W}\mathbf{x} + \mathbf{b})) \quad (15.130)$$

24

25 where $\mathbf{x} \in \mathbb{R}^D$ is the input vector, $y \in \{1, \dots, C\}$ is the class label, \mathbf{W} is a $C \times D$ weight matrix, \mathbf{b}
26 is C -dimensional bias vector, and $\text{softmax}()$ is the **softmax function**, defined as
27

28

$$\text{softmax}(\mathbf{a}) \triangleq \left[\frac{e^{a_1}}{\sum_{c'=1}^C e^{a_{c'}}}, \dots, \frac{e^{a_C}}{\sum_{c'=1}^C e^{a_{c'}}} \right] \quad (15.131)$$

29

30 If we define the logits as $\boldsymbol{\eta}_n = \mathbf{W}\mathbf{x}_n + \mathbf{b}$, the probabilities as $\boldsymbol{\mu}_n = \text{softmax}(\boldsymbol{\eta}_n)$, and let \mathbf{y}_n be
31 the one-hot encoding of the label y_n , then the log likelihood can be written as as the negative cross
32 entropy:
33

34

$$\log p(\mathcal{D}|\boldsymbol{\theta}) = \log \prod_{n=1}^N \prod_{c=1}^C \mu_{nc}^{y_{nc}} = \sum_{n=1}^N \sum_{c=1}^C y_{nc} \log \mu_{nc} \quad (15.132)$$

35

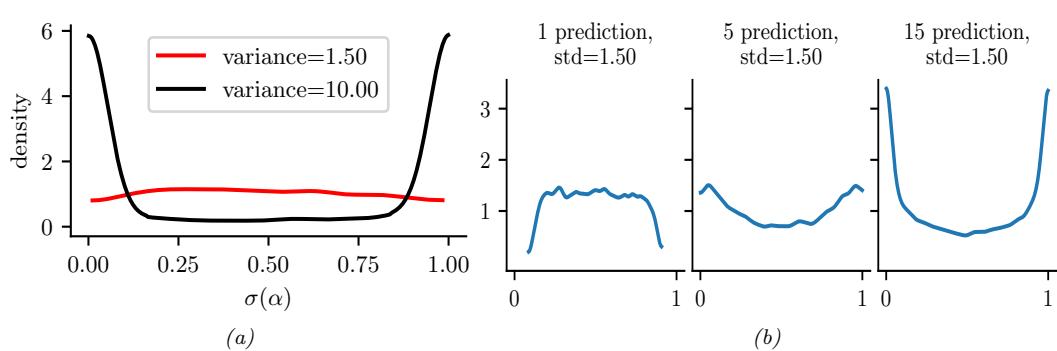


Figure 15.5: (a) Prior on logistic regression output when using $\mathcal{N}(0, \omega)$ prior for the offset term, for $\omega = 10$ or $\omega = 1.5$. Adapted from Figure 11.3 of [McE20]. Generated by `logreg_prior_offset.ipynb`. (b) Distribution over the fraction of 1s we expect to see when using binary logistic regression applied to random binary feature vectors of increasing dimensionality. We use a $\mathcal{N}(0, 1.5)$ prior on the regression coefficients. Adapted from Figure 3 of [Gel+20]. Generated by `logreg_prior.ipynb`.

15.3.3 Priors

As with linear regression, it is standard to use Gaussian priors for the weights in a logistic regression model. It is natural to set the prior mean to 0, to reflect the fact that the output could either increase or decrease in probability depending on the input. But how do we set the prior variance? It is tempting to use a large value, to approximate a uniform distribution, but this is a bad idea. To see why, consider a binary logistic regression model with just an offset term and no features:

$$p(y|\boldsymbol{\theta}) = \text{Ber}(y|\sigma(\alpha)) \quad (15.133)$$

$$p(\alpha) = \mathcal{N}(\alpha|0, \omega) \quad (15.134)$$

If we set the prior to the large value of $\omega = 10$, the implied prior for y is an extreme distribution, with most of its density near 0 or 1, as shown in Figure 15.5a. By contrast, if we use the smaller value of $\omega = 1.5$, we get a flatter distribution, as shown.

If we have input features, the problem gets a little trickier, since the magnitude of the logits will now depend on the number and distribution of the input variables. For example, suppose we generate N random binary vectors \mathbf{x}_n , each of dimension D , where $x_{nd} \sim \text{Ber}(p)$, where $p = 0.8$. We then compute $p(y_n = 1|\mathbf{x}_n) = \sigma(\boldsymbol{\beta}^\top \mathbf{x}_n)$, where $\boldsymbol{\beta} \sim \mathcal{N}(\mathbf{0}, 1.5\mathbf{I})$. We sample S values of $\boldsymbol{\beta}$, and for each one, we sample a vector of labels, $\mathbf{y}_{1:N,s}$ from the above distribution. We then compute the fraction of positive labels, $f_s = \frac{1}{N} \sum_{n=1}^N \mathbb{I}(y_{n,s} = 1)$. We plot the distribution of $\{f_s\}$ as a function of D in Figure 15.5b. We see that the induced prior is initially flat, but eventually becomes skewed towards the extreme values of 0 and 1. To avoid this, we should standardize the inputs, and scale the variance of the prior by $1/\sqrt{D}$. We can also use a heavier tailed distribution, such as a Cauchy or Student [Gel+08; GLM15].

1 2 **15.3.4 Posteriors**

3 Unfortunately, there is no tractable prior that is conjugate to the logistic likelihood. Hence we cannot
4 compute the posterior analytically, unlike with linear regression, even if we use a Gaussian prior.
5 (This mirrors the case with MLE, where we have a closed form solution for linear regression, but not
6 for logistic regression.) Fortunately, there are a range of approximate inference methods we can use,
7 as we discuss in the sections below.

8
9 **15.3.5 Laplace approximation**

10 As we discuss in Section 7.4.3, the Laplace approximation approximates the posterior using a Gaussian.
11 The mean of the Gaussian is equal to the MAP estimate $\hat{\mathbf{w}}$, and the covariance is equal to the inverse
12 Hessian \mathbf{H} computed at the MAP estimate, i.e., $p(\mathbf{w}|\mathcal{D}) \approx \mathcal{N}(\mathbf{w}|\hat{\mathbf{w}}, \mathbf{H})$. We can find the mode using
13 a standard optimization method, and we can then compute the Hessian at the mode analytically or
14 using automatic differentiation.

15 As an example, consider the binary data illustrated in Figure 15.6(a). There are many parameter
16 settings that correspond to lines that perfectly separate the training data; we show 4 example lines.
17 For each decision boundary in Figure 15.6(a), we plot the corresponding parameter vector as point
18 in the log likelihood surface in Figure 15.6(b). These parameters values are $\mathbf{w}_1 = (3, 1)$, $\mathbf{w}_2 = (4, 2)$,
19 $\mathbf{w}_3 = (5, 3)$, and $\mathbf{w}_4 = (7, 3)$. These points all approximately satisfy $\mathbf{w}_i(1)/\mathbf{w}_i(2) \approx \hat{\mathbf{w}}_{\text{mle}}(1)/\hat{\mathbf{w}}_{\text{mle}}(2)$,
20 and hence are close to the orientation of the maximum likelihood decision boundary. The points
21 are ordered by increasing weight norm (3.16, 4.47, 5.83, and 7.62). The unconstrained MLE has
22 $\|\mathbf{w}\| = \infty$, so is infinitely far to the top right.

23 To ensure a unique solution, we use a (spherical) Gaussian prior centered at the origin, $\mathcal{N}(\mathbf{w}|\mathbf{0}, \sigma^2 \mathbf{I})$.
24 The value of σ^2 controls the strength of the prior. If we set $\sigma^2 = \infty$, we force the MAP estimate
25 to be $\mathbf{w} = \mathbf{0}$; this will result in maximally uncertain predictions, since all points \mathbf{x} will produce a
26 predictive distribution of the form $p(y=1|\mathbf{x}) = 0.5$. If we set $\sigma^2 = 0$, the MAP estimate becomes
27 the MLE, resulting in minimally uncertain predictions. (In particular, all positively labeled points
28 will have $p(y=1|\mathbf{x}) = 1.0$, and all negatively labeled points will have $p(y=1|\mathbf{x}) = 0.0$, since the
29 data is separable.) As a compromise (to make a nice illustration), we pick the value $\sigma^2 = 100$.

30 Multiplying this prior by the likelihood results in the unnormalized posterior shown in Figure 15.6(c).
31 The MAP estimate is shown by the blue dot. The Laplace approximation to this posterior is shown
32 in Figure 15.6(d). We see that it gets the mode correct (by construction), but the shape of the
33 posterior is somewhat distorted. (The southwest-northeast orientation captures uncertainty about the
34 magnitude of \mathbf{w} , and the southeast-northwest orientation captures uncertainty about the orientation
35 of the decision boundary.)

36 Next we need to convert the posterior over the parameters into a posterior over predictions, as
37 follows:

$$\frac{40}{41} p(y|\mathbf{x}, \mathcal{D}) = \int p(y|\mathbf{x}, \mathbf{w})p(\mathbf{w}|\mathcal{D})d\mathbf{w} \quad (15.135)$$

42 The simplest way to evaluate this integral is to use a Monte Carlo approximation:

$$\frac{43}{44} p(y=1|\mathbf{x}, \mathcal{D}) \approx \frac{1}{S} \sum_{s=1}^S \sigma(\mathbf{w}_s^\top \mathbf{x}) \quad (15.136)$$

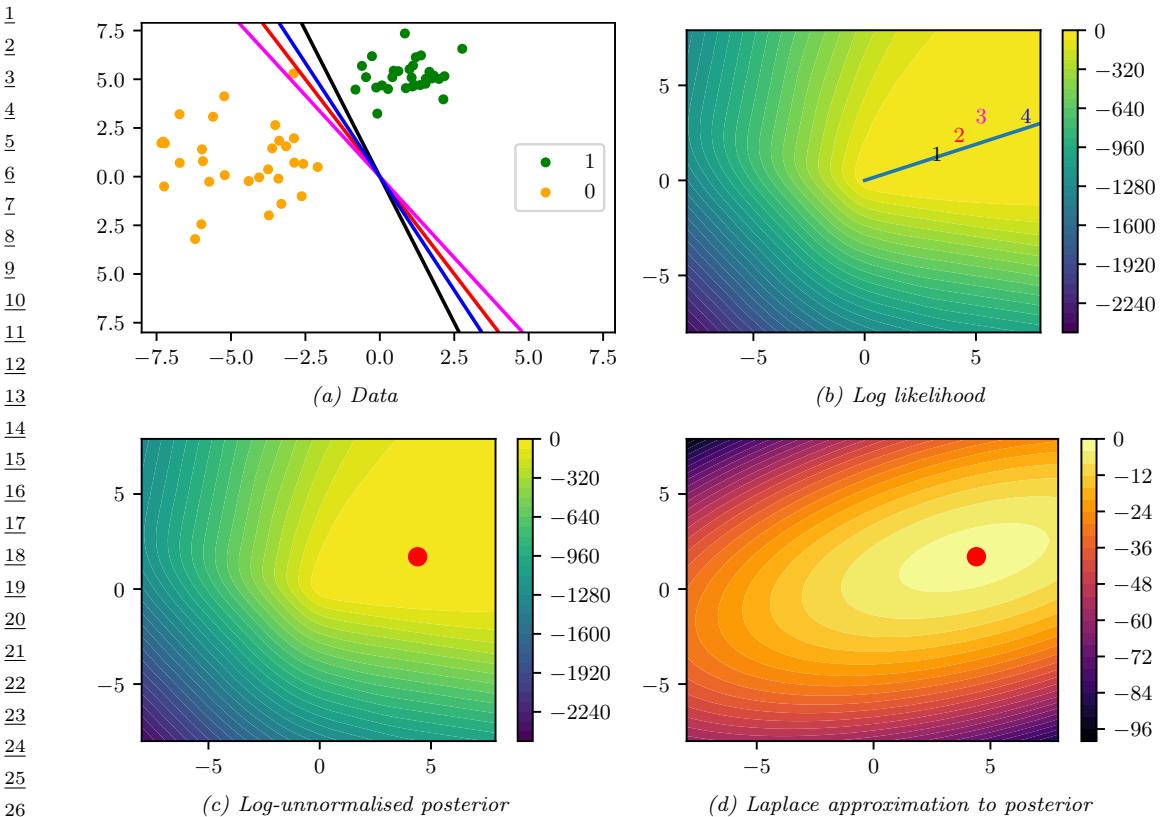


Figure 15.6: (a) Illustration of the data and some decision boundaries, where colors correspond to the points in the (b) panel. (b) Log-likelihood for a logistic regression model. The line is drawn from the origin in the direction of the MLE (which is at infinity). The numbers correspond to 4 points in parameter space, corresponding to the lines in (a). (c) Unnormalised log posterior (assuming vague spherical prior). (d) Laplace approximation to posterior. Adapted from a figure by Mark Girolami. Generated by [logreg_laplace_demo.ipynb](#).

where $\mathbf{w}_s \sim p(\mathbf{w}|\mathcal{D})$.

Alternatively, we can use the deterministic **probit approximation** first suggested in [SL90]. If we define $a = \mathbf{x}^\top \mathbf{w}$, and $p(\mathbf{w}|\mathcal{D}) = \mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$, then we can write this approximation as

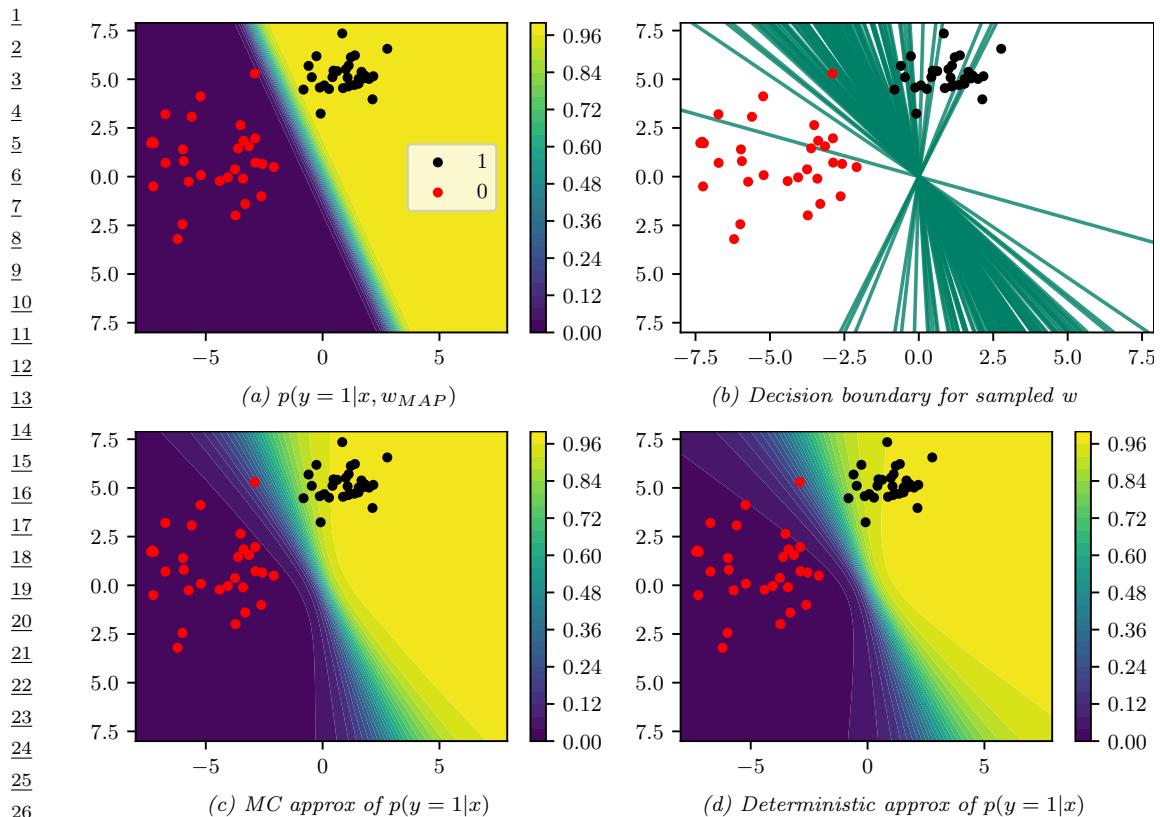
$$p(y=1|\mathbf{x}, \mathcal{D}) \approx \sigma(\kappa(v)m) \quad (15.137)$$

$$\kappa(v) \triangleq (1 + \pi v/8)^{-\frac{1}{2}} \quad (15.138)$$

$$v = \mathbb{V}[a] = \mathbb{V}[\mathbf{x}^\top \mathbf{w}] = \mathbf{x}^\top \boldsymbol{\Sigma} \mathbf{x} \quad (15.139)$$

$$m = \mathbb{E}[a] = \mathbf{x}^\top \boldsymbol{\mu} \quad (15.140)$$

In Figure 15.7, we show contours of the posterior predictive distribution. Figure 15.7(a) shows the plugin approximation using the MAP estimate. We see that there is no uncertainty about the decision boundary, even though we are generating probabilistic predictions over the labels. Figure 15.7(b)



27 *Figure 15.7: Posterior predictive distribution for a logistic regression model in 2d. (a): contours of $p(y =$
28 $1|x, \hat{w}_{MAP})$. (b): samples from the posterior predictive distribution. (c): Averaging over these samples.
29 (d): moderated output (probit approximation). Adapted from a figure by Mark Girolami. Generated by
30 [logreg_laplace_demo.ipynb](#).*

31
32
33 shows what happens when we plug in samples from the Gaussian posterior. Now we see that there is
34 considerable uncertainty about the orientation of the “best” decision boundary. Figure 15.7(c) shows
35 the average of these samples. By averaging over multiple predictions, we see that the uncertainty in
36 the decision boundary “splays out” as we move further from the training data. Figure 15.7(d) shows
37 that the probit approximation gives very similar results to the Monte Carlo approximation.

39 15.3.6 MCMC inference

40
41 Markov chain Monte Carlo, or MCMC, is often considered the “gold standard” for approximate
42 inference, since it makes no explicit assumptions about the form of the posterior. Instead, it just
43 approximates it non-parametrically using a set of S samples:

$$44 \\ 45 q(\boldsymbol{\theta}|\mathcal{D}) \approx \frac{1}{S} \sum_{s=1}^S \delta(\boldsymbol{\theta} - \boldsymbol{\theta}^s) \quad (15.141)$$

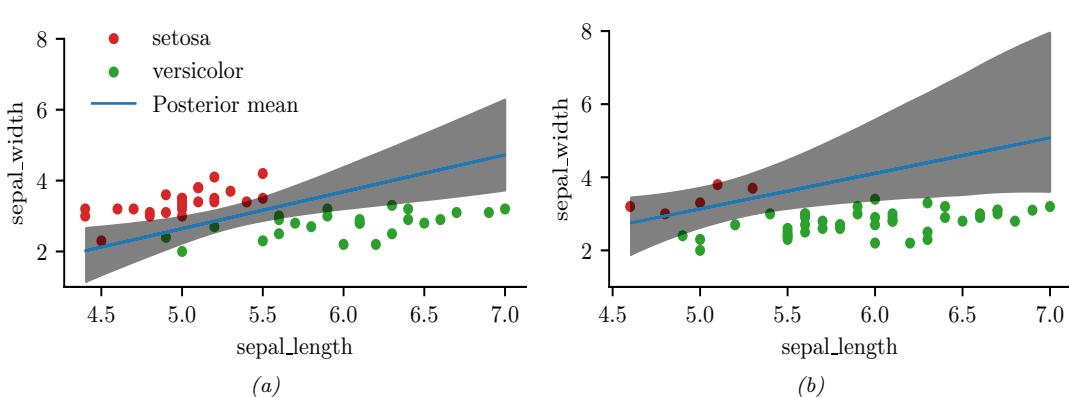


Figure 15.8: Illustration of the posterior over the decision boundary for classifying iris flowers (setosa vs versicolor) using 2 input features. (a) 25 examples per class. Adapted from Figure 4.5 of [Mar18]. (b) 5 examples of class 0, 45 examples of class 1. Adapted from Figure 4.8 of [Mar18]. Generated by [logreg_iris_bayes_2d.ipynb](#).

where $\boldsymbol{\theta}^s \sim p(\boldsymbol{\theta}|\mathcal{D})$ are samples from the posterior.

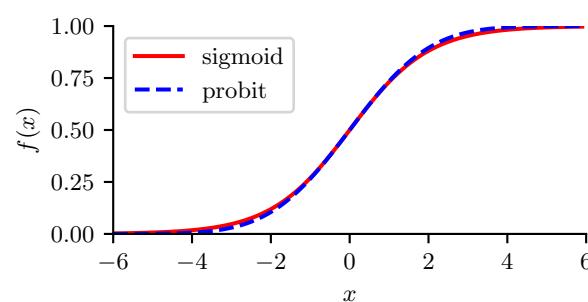
To efficiently compute these samples, we can use the method of Hamiltonian Monte Carlo or HMC, which we describe in Section 12.5. This relies on our ability to compute the gradient of the log joint, $\nabla_{\boldsymbol{\theta}} \log p(\mathcal{D}, \boldsymbol{\theta})$, which we can compute using automatic differentiation.

Let us apply HMC to a 2-dimensional, 2-class version of the iris classification problem, where we just use two input features, sepal length and sepal width, and two classes, Virginica and Non-Virginica. The decision boundary is the set of points (x_1^*, x_2^*) such that $\sigma(b + w_1 x_1^* + w_2 x_2^*) = 0.5$. Such points must lie on the following line:

$$x_2^* = -\frac{b}{w_2} + \left(-\frac{w_1}{w_2} x_1^* \right) \quad (15.142)$$

We can therefore compute an MC approximation to the posterior over decision boundaries by sampling the parameters from the posterior, $(w_1, w_2, b) \sim p(\boldsymbol{\theta}|\mathcal{D})$, and plugging them into the above equation, to get $p(x_1^*, x_2^* | \mathcal{D})$. The results of this method (using a vague Gaussian prior for the parameters) are shown in Figure 15.8a. The solid line is the posterior mean, and the shaded interval is a 95% credible interval. As before, we see that the uncertainty about the location of the boundary is higher as we move away from the training data.

In Figure 15.8b, we show what happens to the decision boundary when we have unbalanced classes. We notice two things. First, the posterior uncertainty increases, because we have less data from the blue class. Second, we see that the posterior mean of the decision boundary shifts towards the class with less data. This follows from linear discriminant analysis, where one can show that changing the class prior changes the location of the decision boundary, so that more of the input space gets mapped to the class which is higher a priori. (See [Mur22, Sec 9.2] for details.)



¹¹ *Figure 15.9: The logistic (sigmoid) function $\sigma(x)$ in solid red, with the Gaussian cdf function $\Phi(\lambda x)$ in dotted
¹² blue superimposed. Here $\lambda = \sqrt{\pi}/8$, which was chosen so that the derivatives of the two curves match at
¹³ $x = 0$. Adapted from Figure 4.9 of [Bis06]. Generated by [probit_plot.ipynb](#).*

¹⁴

¹⁵

¹⁶ 15.3.7 Variational inference

¹⁷

¹⁸ As we discuss in Section 10.1, variational inference converts approximate inference into an optimization
¹⁹ problem. It does this by choosing an approximate distribution $q(\mathbf{w}; \psi)$ and optimising the variational
²⁰ parameters ψ to maximize the evidence lower bound (ELBO). This has the effect of making
²¹ $q(\mathbf{w}; \psi) \approx p(\mathbf{w} | \mathcal{D})$ in the sense that the KL divergence is small. There are several ways to tackle
²² this: use a stochastic estimate of the ELBO (see Section 10.3.3), use the conditionally conjugate VI
²³ method of [Supplementary](#) Section 10.3.1.2, or use a “local” VI method that creates a quadratic lower
²⁴ bound to the logistic function (see [Supplementary](#) Section 15.1).

²⁵

²⁶ 15.3.8 Assumed density filtering

²⁷

²⁸ In Section 8.10.3, we discuss how to use assumed density filtering (ADF) to recursively approximate
²⁹ the posterior $p(\mathbf{w} | \mathcal{D}_{1:t})$ for a logistic regression model in an online fashion.

³⁰

³¹ 15.4 Probit regression

³²

³³ In this section, we discuss **probit regression**, which is similar to binary logistic regression except
³⁴ it uses $\mu_n = \Phi(a_n)$ instead of $\mu_n = \sigma(a_n)$ as the mean function, where Φ is the cdf of the standard
³⁵ normal, and $a_n = \mathbf{w}^\top \mathbf{x}_n$. The corresponding link function is therefore $a_n = \ell(\mu_n) = \Phi^{-1}(\mu_n)$; the
³⁶ inverse of the Gaussian cdf is known as the **probit function**.

³⁷ The Gaussian cdf Φ is very similar to the logistic function, as shown in Figure 15.9. Thus probit
³⁸ regression and “regular” logistic regression behave very similarly. However, probit regression has some
³⁹ advantages. In particular, it has a simple interpretation as a latent variable model (see Section 15.4.1),
⁴⁰ which arises from the field of **choice theory** as studied in economics (see e.g., [Koo03]). This also
⁴¹ simplifies the task of Bayesian parameter inference.

⁴²

⁴³ 15.4.1 Latent variable interpretation

⁴⁴

⁴⁵ We can interpret $a_n = \mathbf{w}^\top \mathbf{x}_n$ as a factor that is proportional to how likely a person is respond
⁴⁶ positively (generate $y_n = 1$) given input \mathbf{x}_n . However, typically there are other unobserved factors that
⁴⁷

influence someone's response. Let us model these hidden factors by Gaussian noise, $\epsilon_n \sim \mathcal{N}(0, 1)$. Let the combined preference for positive outcomes be represented by the latent variable $z_n = \mathbf{w}^\top \mathbf{x}_n + \epsilon_n$. We assume that the person will pick the positive label iff this latent factor is positive rather than negative, i.e.,

$$y_n = \mathbb{I}(z_n \geq 0) \quad (15.143)$$

When we marginalize out z_n , we recover the probit model:

$$p(y_n = 1 | \mathbf{x}_n, \mathbf{w}) = \int \mathbb{I}(z_n \geq 0) \mathcal{N}(z_n | \mathbf{w}^\top \mathbf{x}_n, 1) dz_n \quad (15.144)$$

$$= p(\mathbf{w}^\top \mathbf{x}_n + \epsilon_n \geq 0) = p(\epsilon_n \geq -\mathbf{w}^\top \mathbf{x}_n) \quad (15.145)$$

$$= 1 - \Phi(-\mathbf{w}^\top \mathbf{x}_n) = \Phi(\mathbf{w}^\top \mathbf{x}_n) \quad (15.146)$$

Thus we can think of probit regression as a threshold function applied to noisy input.

We can interpret logistic regression in the same way. However, in that case the noise term ϵ_n comes from a **logistic distribution**, defined as follows:

$$f(y|\mu, s) \triangleq \frac{e^{-\frac{y-\mu}{s}}}{s(1 + e^{-\frac{y-\mu}{s}})^2} \quad (15.147)$$

The cdf of this distribution is given by

$$F(y|\mu, s) = \frac{1}{1 + e^{-\frac{y-\mu}{s}}} \quad (15.148)$$

It is clear that if we use logistic noise with $\mu = 0$ and $s = 1$ we recover logistic regression. However, it is computationally easier to deal with Gaussian noise, as we show below.

15.4.2 Maximum likelihood estimation

In this section, we discuss some methods for fitting probit regression using MLE.

15.4.2.1 MLE using SGD

We can find the MLE for probit regression using standard gradient methods. Let $\mu_n = \mathbf{w}^\top \mathbf{x}_n$, and let $\tilde{y}_n \in \{-1, +1\}$. Then the gradient of the log-likelihood for a single example n is given by

$$\mathbf{g}_n \triangleq \frac{d}{d\mathbf{w}} \log p(\tilde{y}_n | \mathbf{w}^\top \mathbf{x}_n) = \frac{d\mu_n}{d\mathbf{w}} \frac{d}{d\mu_n} \log p(\tilde{y}_n | \mathbf{w}^\top \mathbf{x}_n) = \mathbf{x}_n \frac{\tilde{y}_n \phi(\mu_n)}{\Phi(\tilde{y}_n \mu_n)} \quad (15.149)$$

where ϕ is the standard normal pdf, and Φ is its cdf. Similarly, the Hessian for a single case is given by

$$\mathbf{H}_n = \frac{d}{d\mathbf{w}^2} \log p(\tilde{y}_n | \mathbf{w}^\top \mathbf{x}_n) = -\mathbf{x}_n \left(\frac{\phi(\mu_n)^2}{\Phi(\tilde{y}_n \mu_n)^2} + \frac{\tilde{y}_n \mu_n \phi(\mu_n)}{\Phi(\tilde{y}_n \mu_n)} \right) \mathbf{x}_n^\top \quad (15.150)$$

This can be passed to any gradient-based optimizer.

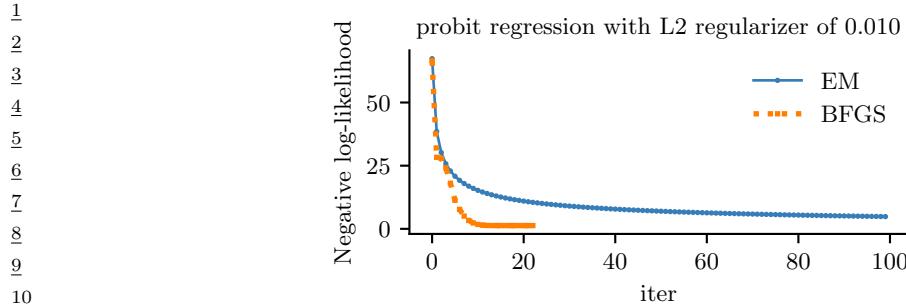


Figure 15.10: Fitting a probit regression model in 2d using a quasi-Newton method or EM. Generated by [probit_reg_demo.ipynb](#).

15.4.2.2 MLE using EM

We can use the latent variable interpretation of probit regression to derive an elegant EM algorithm for fitting the model. The complete data log likelihood has the following form, assuming a $\mathcal{N}(\mathbf{0}, \mathbf{V}_0)$ prior on \mathbf{w} :

$$\ell(\mathbf{z}, \mathbf{w} | \mathbf{V}_0) = \log p(\mathbf{y} | \mathbf{z}) + \log \mathcal{N}(\mathbf{z} | \mathbf{X}\mathbf{w}, \mathbf{I}) + \log \mathcal{N}(\mathbf{w} | \mathbf{0}, \mathbf{V}_0) \quad (15.151)$$

$$= \sum_n \log p(y_n | z_n) - \frac{1}{2} (\mathbf{z} - \mathbf{X}\mathbf{w})^\top (\mathbf{z} - \mathbf{X}\mathbf{w}) - \frac{1}{2} \mathbf{w}^\top \mathbf{V}_0^{-1} \mathbf{w} \quad (15.152)$$

The posterior in the E step is a **truncated Gaussian**:

$$p(z_n | y_n, \mathbf{x}_n, \mathbf{w}) = \begin{cases} \mathcal{N}(z_n | \mathbf{w}^\top \mathbf{x}_n, 1) \mathbb{I}(z_n > 0) & \text{if } y_n = 1 \\ \mathcal{N}(z_n | \mathbf{w}^\top \mathbf{x}_n, 1) \mathbb{I}(z_n < 0) & \text{if } y_n = 0 \end{cases} \quad (15.153)$$

In Equation (15.152), we see that \mathbf{w} only depends linearly on \mathbf{z} , so we just need to compute $\mathbb{E}[z_n | y_n, \mathbf{x}_n, \mathbf{w}]$, so we just need to compute the posterior mean. One can show that this is given by

$$\mathbb{E}[z_n | \mathbf{w}, \mathbf{x}_n] = \begin{cases} \mu_n + \frac{\phi(\mu_n)}{1 - \Phi(-\mu_n)} = \mu_n + \frac{\phi(\mu_n)}{\Phi(\mu_n)} & \text{if } y_n = 1 \\ \mu_n - \frac{\phi(\mu_n)}{\Phi(-\mu_n)} = \mu_n - \frac{\phi(\mu_n)}{1 - \Phi(\mu_n)} & \text{if } y_n = 0 \end{cases} \quad (15.154)$$

where $\mu_n = \mathbf{w}^\top \mathbf{x}_n$.

In the M step, we estimate \mathbf{w} using ridge regression, where $\boldsymbol{\mu} = \mathbb{E}[\mathbf{z}]$ is the output we are trying to predict. Specifically, we have

$$\hat{\mathbf{w}} = (\mathbf{V}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\mu} \quad (15.155)$$

The EM algorithm is simple, but can be much slower than direct gradient methods, as illustrated in Figure 15.10. This is because the posterior entropy in the E step is quite high, since we only observe that z is positive or negative, but are given no information from the likelihood about its magnitude. Using a stronger regularizer can help speed convergence, because it constrains the range of plausible z values. In addition, one can use various speedup tricks, such as data augmentation [DM01].

15.4.3 Bayesian inference

It is possible to use the latent variable formulation of probit regression in Section 15.4.2.2 to derive a simple Gibbs sampling algorithm for approximating the posterior $p(\mathbf{w}|\mathcal{D})$ (see e.g., [AC93; HH06]).

The key idea is to use an auxiliary latent variable, which, when conditioned on, makes the whole model a conjugate linear-Gaussian model. The full conditional for the latent variables is given by

$$p(z_i|y_i, \mathbf{x}_i, \mathbf{w}) = \begin{cases} \mathcal{N}(z_i|\mathbf{w}^T \mathbf{x}_i, 1)\mathbb{I}(z_i > 0) & \text{if } y_i = 1 \\ \mathcal{N}(z_i|\mathbf{w}^T \mathbf{x}_i, 1)\mathbb{I}(z_i < 0) & \text{if } y_i = 0 \end{cases} \quad (15.156)$$

Thus the posterior is a truncated Gaussian. We can sample from a truncated Gaussian, $\mathcal{N}(z|\mu, \sigma)\mathbb{I}(a \leq z \leq b)$ in two steps: first sample $u \sim U(\Phi((a - \mu)/\sigma), \Phi((b - \mu)/\sigma))$, then set $z = \mu + \sigma\Phi^{-1}(u)$ [Rob95a].

The full conditional for the parameters is given by

$$p(\mathbf{w}|\mathcal{D}, \mathbf{z}, \boldsymbol{\lambda}) = \mathcal{N}(\mathbf{w}_N, \mathbf{V}_N) \quad (15.157)$$

$$\mathbf{V}_N = (\mathbf{V}_0^{-1} + \mathbf{X}^T \mathbf{X})^{-1} \quad (15.158)$$

$$\mathbf{w}_N = \mathbf{V}_N(\mathbf{V}_0^{-1} \mathbf{m}_0 + \mathbf{X}^T \mathbf{z}) \quad (15.159)$$

For further details, see e.g., [AC93; FSF10]. It is also possible to use variational Bayes, which tends to be much faster (see e.g., [GR06a; FDZ19]).

15.4.4 Ordinal probit regression

One advantage of the latent variable interpretation of probit regression is that it is easy to extend to the case where the response variable is ordered in some way, such as the outputs low, medium and high. This is called **ordinal regression**. The basic idea is as follows. If there are C output values, we introduce $C + 1$ thresholds γ_j and set

$$y_n = j \quad \text{if} \quad \gamma_{j-1} < z_n \leq \gamma_j \quad (15.160)$$

where $\gamma_0 \leq \dots \leq \gamma_C$. For identifiability reasons, we set $\gamma_0 = -\infty$, $\gamma_1 = 0$ and $\gamma_C = \infty$. For example, if $C = 2$, this reduces to the standard binary probit model, whereby $z_n < 0$ produces $y_n = 0$ and $z_n \geq 0$ produces $y_n = 1$. If $C = 3$, we partition the real line into 3 intervals: $(-\infty, 0]$, $(0, \gamma_2]$, (γ_2, ∞) . We can vary the parameter γ_2 to ensure the right relative amount of probability mass falls in each interval, so as to match the empirical frequencies of each class label. See e.g., [AC93] for further details.

Finding the MLEs for this model is a bit trickier than for binary probit regression, since we need to optimize for \mathbf{w} and $\boldsymbol{\gamma}$, and the latter must obey an ordering constraint. See e.g., [KL09] for an approach based on EM. It is also possible to derive a simple Gibbs sampling algorithm for this model (see e.g., [Hof09, p216]).

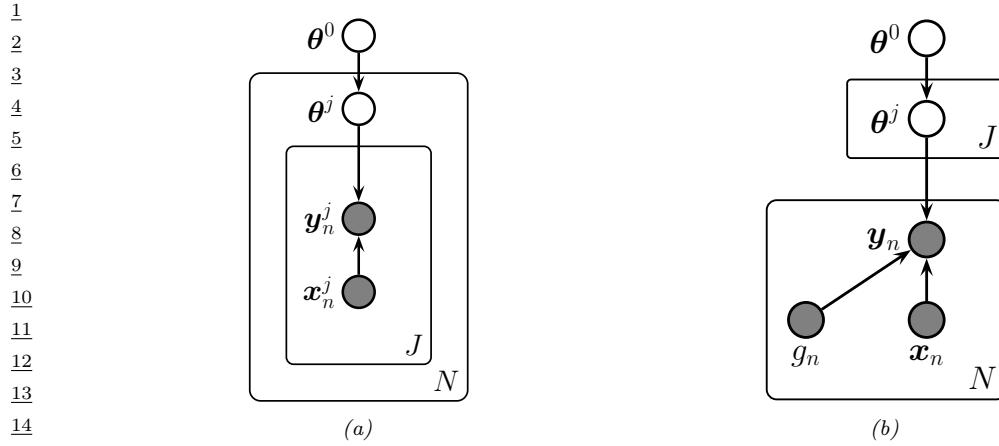


Figure 15.11: Hierarchical Bayesian discriminative models with J groups. (a) Nested formulation. (b) Non-nested formulation, with group indicator $g_n \in \{1, \dots, J\}$.

15.4.5 Multinomial probit models

Now consider the case where the response variable can take on C unordered categorical values, $y_n \in \{1, \dots, C\}$. The **multinomial probit** model is defined as follows:

$$z_{nc} = \mathbf{w}_c^\top \mathbf{x}_{nc} + \epsilon_{nc} \quad (15.161)$$

$$\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{R}) \quad (15.162)$$

$$y_n = \arg \max_c z_{nc} \quad (15.163)$$

See e.g., [DE04; GR06b; Sco09; FSF10] for more details on the model and its connection to multinomial logistic regression.

If instead of setting $y_n = \arg \max_c z_{ic}$ we use $y_{nc} = \mathbb{I}(z_{nc} > 0)$, we get a model known as **multivariate probit**, which is one way to model C correlated binary outcomes (see e.g., [TMD12]).

15.5 Multi-level (hierarchical) GLMs

Suppose we have a set of J related datasets, each of which contains a series of N_j datapoints $\mathcal{D}_j = \{(\mathbf{x}_n^j, \mathbf{y}_n^j) : n = 1 : N_j\}$. There are 3 main ways to fit models in such a setting: we could fit J separate models, $p(\mathbf{y}|\mathbf{x}; \mathcal{D}_j)$, which might result in overfitting if some \mathcal{D}_j are small; we could pool all the data to get $\mathcal{D} = \bigcup_{j=1}^J \mathcal{D}_j$ and fit a single model, $p(\mathbf{y}|\mathbf{x}; \mathcal{D})$, which might result in underfitting; or we can use a **hierarchical Bayesian model**, also called a **multilevel model** or **partially pooled model**, in which we assume each group has its own parameters, θ^j , but that these have something in common, as modeled by a shared global prior $p(\theta^0)$. (Note that each group could be a single individual.) The overall model has the form

$$p(\theta^{0:J}, \mathcal{D}) = p(\theta^0) \prod_{j=1}^J \left[p(\theta^j | \theta^0) \prod_{n=1}^{N_j} p(\mathbf{y}_n^j | \mathbf{x}_n^j, \theta^j) \right] \quad (15.164)$$

1 See Figure 15.11a, which represents the model using nested plate notation.
 2

3 It is often more convenient to represent the model as in Figure 15.11b, which eliminates the
 4 nested plates (and hence the double indexing of variables) by associating a group indicator variable
 5 $g_n \in \{1, \dots, J\}$, which specifies which set of parameters to use for each data point. Thus the model
 6 now has the form

$$7 \quad p(\boldsymbol{\theta}^{0:J}, \mathcal{D}) = p(\boldsymbol{\theta}^0) \left[\prod_{j=1}^J p(\boldsymbol{\theta}^j | \boldsymbol{\theta}^0) \right] \left[\prod_{n=1}^N p(\mathbf{y}_n | \mathbf{x}_n, g_n, \boldsymbol{\theta}) \right] \quad (15.165)$$

11 where

$$12 \quad p(\mathbf{y}_n | \mathbf{x}_n, g_n, \boldsymbol{\theta}) = \prod_{j=1}^J p(\mathbf{y}_n | \mathbf{x}_n, \boldsymbol{\theta}^j)^{\mathbb{I}(g_n=j)} \quad (15.166)$$

16 If the likelihood function is a GLM, this hierarchical model is called a **hierarchical GLM** [LN96].
 17 This class of models is very widely used in applied statistics. For much more details, see e.g., [GH07;
 18 GHV20; Gel+22].

20 15.5.1 Generalized linear mixed models (GLMMs)

22 Suppose that the prior on the per-group parameters is Gaussian, so $p(\boldsymbol{\theta}^j | \boldsymbol{\theta}^0) = \mathcal{N}(\boldsymbol{\theta}^j | \boldsymbol{\theta}^0, \boldsymbol{\Sigma}^j)$. If we
 23 have a GLM likelihood, the model becomes

$$25 \quad p(\mathbf{y}_n | \mathbf{x}_n, g_n = j, \boldsymbol{\theta}) = p(\mathbf{y}_n | \ell(\eta_n)) \quad (15.167)$$

$$26 \quad \eta_n = \mathbf{x}_n^\top \boldsymbol{\theta}^j = \mathbf{x}_n^\top (\boldsymbol{\theta}^0 + \boldsymbol{\epsilon}^j) = \mathbf{x}_n^\top \boldsymbol{\theta}^0 + \mathbf{x}_n^\top \boldsymbol{\epsilon}^j \quad (15.168)$$

28 where ℓ is the link function, and $\boldsymbol{\epsilon}^j \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$. This is known as a **generalized linear mixed**
 29 **model (GLMM)** or **mixed effects model**. The shared (common) parameters $\boldsymbol{\theta}^0$ are called **fixed**
 30 **effects**, and the group-specific offsets $\boldsymbol{\epsilon}^j$ are called **random effects**.² We can see that the random
 31 effects model group-specific deviations or idiosyncrasies away from the shared fixed parameters.
 32 Furthermore, we see that the random effects are correlated, which allows us to model dependencies
 33 between the observations that would not be captured by a standard GLM.

35 15.5.2 Model fitting

37 We can fit GLMMs, and hierarchical models more generally, using standard Bayesian inference
 38 methods. We can use a variety of algorithms, such as HMC (see e.g., [BG13]), variational Bayes (see
 39 e.g., [HOW11; TN13]), expectation propagation (see e.g., [KW18]), etc. In this section, we use HMC,
 40 since it is simple and efficient.³

41 2. Note that there are multiple definitions of the terms “fixed effects” and random effects”, as explained in this blog
 42 post by Andrew Gelman: https://statmodeling.stat.columbia.edu/2005/01/25/why_i_dont_use/.

43 3. There are many standard software packages for HMC analysis of hierarchical GLMs, such as **Bambi** (<https://github.com/bambinos/bambi>), which is a Python wrapper on top of PyMC, Blackjax and numpyro samplers;
 44 **RStanARM** (<https://cran.r-project.org/web/packages/rstanarm/index.html>), which is an R wrapper on top of
 45 Stanl and **BRMS** ([https://cran.r-project.org/web/packages\(brms/index.html](https://cran.r-project.org/web/packages(brms/index.html)), which is another R wrapper on
 46 top of Stan, but which also needs a C++ compiler.

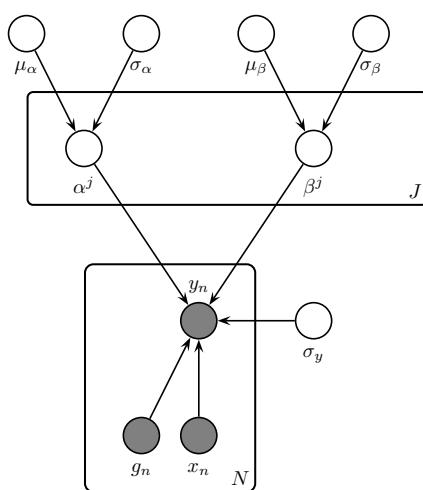


Figure 15.12: A hierarchical Bayesian linear regression model for the radon problem.

15.5.3 Example: radon regression

In this section, we give an example of a hierarchical Bayesian linear regression model. We apply it to a simplified version of the `radon` example from [Gel+14a, Sec 9.4].

Radon is known to be the highest cause of lung cancer in non-smokers, so reducing it where possible is desirable. To help with this, we fit a regression model, that predicts the (log) radon level as a function of the location of the house, as represented by a categorical feature indicating its county, and a binary feature representing whether the house has a basement or not. We use a dataset consisting of $J = 85$ counties in Minnesota; each county has between 2 and 80 measurements.

We assume the following likelihood:

$$p(y_n | x_n, g_n = j, \theta) = \mathcal{N}(y_n | \alpha_j + \beta_j x_n, \sigma_y^2) \quad (15.169)$$

where $g_n \in \{1, \dots, J\}$ is the county for house i , and $x_n \in \{0, 1\}$ indicates if the floor is at level 0 (i.e., in the basement) or level 1 (i.e., above ground). Intuitively we expect the radon levels to be lower in houses without basements, since they are more insulated from the earth which is the source of the radon.

Since some counties have very few data points, we use a hierarchical prior in which we assume $\alpha_j \sim \mathcal{N}(\mu_\alpha, \sigma_\alpha^2)$ and $\beta_j \sim \mathcal{N}(\mu_\beta, \sigma_\beta^2)$. We use weak priors for the parameters: $\mu_\alpha \sim \mathcal{N}(0, 1)$, $\mu_\beta \sim \mathcal{N}(0, 1)$, $\sigma_\alpha \sim \mathcal{C}_+(1)$, $\sigma_\beta \sim \mathcal{C}_+(1)$, $\sigma_y \sim \mathcal{C}_+(1)$. See Figure 15.12 for the graphical model.

15.5.3.1 Posterior inference

Figure 15.13 shows the posterior marginals for μ_α , μ_β , α_j and β_j . We see that μ_β is close to -0.6 with high probability, which confirms our suspicion that having $x = 1$ (i.e., no basement) decreases the amount of radon in the house. We also see that the distribution of the α_j parameters is quite variable, due to different base rates across the counties.

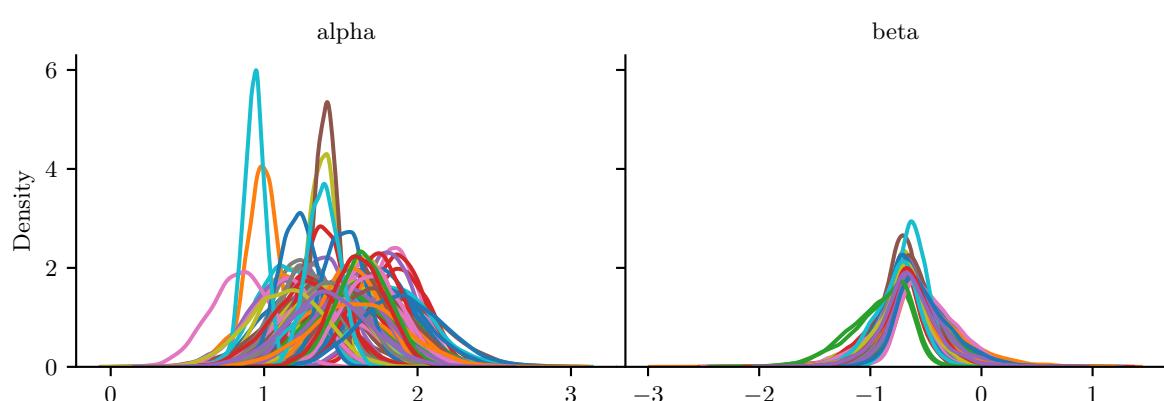


Figure 15.13: Posterior marginals for α_c and β_c for each county in the radon model. Generated by [linreg_hierarchical_non_centered.ipynb](#).

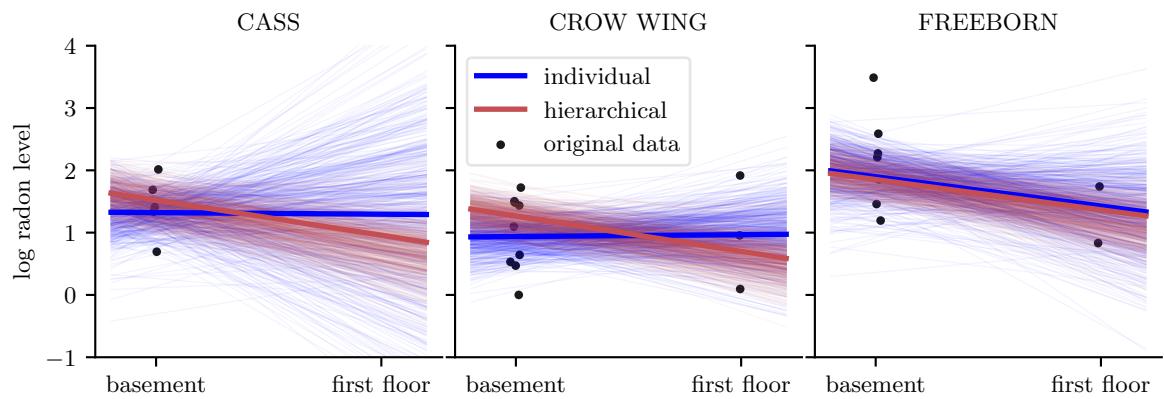


Figure 15.14: Predictions from the radon model for 3 different counties in Minnesota. Black dots are observed datapoints. Red represents results of hierarchical (shared) prior, blue represents results of non-hierarchical prior. Thick lines are the result of using the posterior mean, thin lines are the result of using posterior samples. Generated by [linreg_hierarchical_non_centered.ipynb](#).

Figure 15.14 shows predictions from the hierarchical and non-hierarchical model for 3 different counties. We see that the predictions from the hierarchical model are more consistent across counties, and work well even if there are no examples of certain feature combinations for a given county (e.g., there are no houses without basements in the sample from Cass county). If we sample data from the posterior predictive distribution, and compare it to the real data, we find that the RMSE is 0.13 for the non-hierarchical model and 0.08 for the hierarchical model, indicating that the latter fits better.

15.5.3.2 Non-centered parameterization

One problem that frequently arises in hierarchical models is that the parameters be very correlated. This can cause computational problems when performing inference.

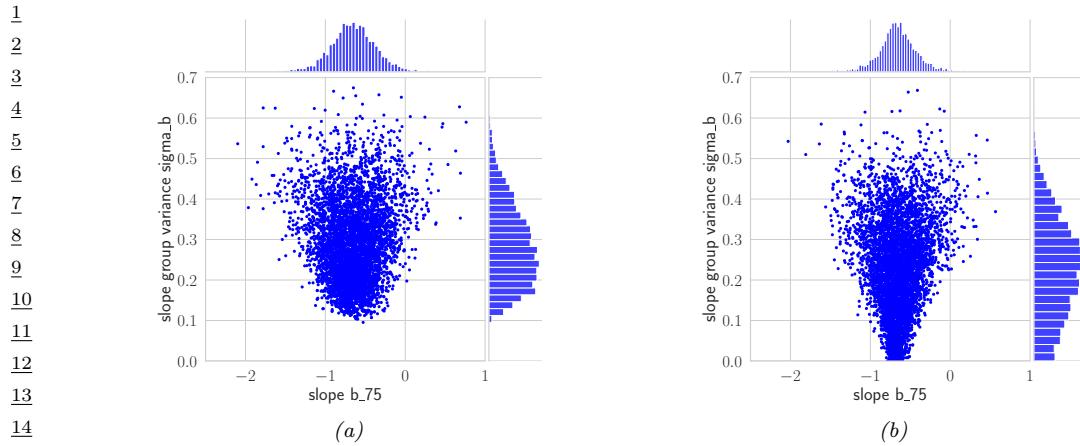


Figure 15.15: (a) Bivariate posterior $p(\beta_j, \sigma_\beta | \mathcal{D})$ for the hierarchical radon model for county $j = 75$ using centered parameterization. (b) Similar to (a) except we plot $p(\tilde{\beta}_j, \sigma_\beta | \mathcal{D})$ for the non-centered parameterization.
Generated by [linreg_hierarchical_non_centered.ipynb](#).

Figure 15.15a gives an example where we plot $p(\beta_j, \sigma_\beta | \mathcal{D})$ for some specific county j . If we believe that σ_β is large, then β_j is “allowed” to vary a lot, and we get the broad distribution at the top of the figure. However, if we believe that σ_β is small, then β_j is constrained to be close to the global prior mean of μ_β , so we get the narrow distribution at the bottom of the figure. This is often called **Neal’s funnel**, after a paper by Radford Neal [Nea03]. It is difficult for many algorithms (especially sampling algorithms) to explore parts of parameter space at the bottom of the funnel. This is evident from the marginal posterior for σ_β shown (as a histogram) on the right hand side of the plot: we see that it excludes the interval $[0, 0.1]$, thus ruling out models in which we shrink β_j all the way to 0. In cases where a covariate has no useful predictive role, we would like to be able to induce sparsity, so we need to overcome this problem.

A simple solution to this is to use a **non-centered parameterization** [PR03]. That is, we replace $\beta_j \sim \mathcal{N}(\mu_\beta, \sigma_\beta^2)$ with $\beta_j = \mu_\beta + \tilde{\beta}_j \sigma_\beta$, where $\tilde{\beta}_j \sim \mathcal{N}(0, 1)$ represents the *offset* from the global mean, μ_β . The correlation between $\tilde{\beta}_j$ and σ_β is much less, as shown in Figure 15.15b. See Section 12.6.5 for more details.

16 Deep neural networks

16.1 Introduction

The term “**deep neural network**” or **DNN**, in its modern usage, refers to any kind of differentiable function that can be expressed as a **computation graph**, where the nodes are primitive operations (like matrix multiplication), and edges represent numeric data in the form of vectors, matrices, or tensors. In its simplest form, this graph can be constructed as a linear series of nodes or “**layers**”. The term “deep” refers to models with many such layers.

In Section 16.2 we discuss some of the basic building blocks (node types) that are used in the field. In Section 16.3 we give examples of common architectures which are constructed from these building blocks. In Section 6.2 we show how we can efficiently compute the gradient of functions defined on such graphs. If the function computes the scalar loss of the model’s predictions given a training set, we can pass this gradient to an optimization routine, such as those discussed in Chapter 6, in order to fit the model. Fitting such models to data is called “**deep learning**”.

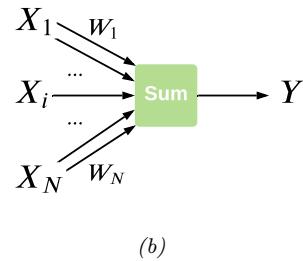
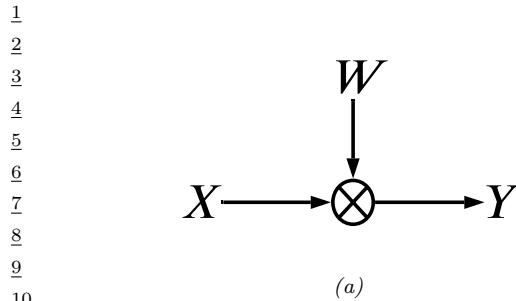
We can combine DNNs with probabilistic models in two different ways. The first is to use them to define nonlinear functions which are used inside conditional distributions. For example, we may construct a classifier using $p(y|\mathbf{x}, \boldsymbol{\theta}) = \text{Cat}(y|\text{softmax}(f(\mathbf{x}; \boldsymbol{\theta})))$, where $f(\mathbf{x}; \boldsymbol{\theta})$ is a neural network that maps inputs \mathbf{x} and parameters $\boldsymbol{\theta}$ to output logits. Or we may construct a joint probability distribution over multiple variables using a directed graphical model (Chapter 4) where each CPD $p(\mathbf{x}_i|\text{pa}(\mathbf{x}_i))$ is a DNN. This lets us construct expressive probability models.

The other way we can combine DNNs and probabilistic models is to use DNNs to approximate the posterior distribution, i.e., we learn a function f to compute $q(\mathbf{z}|f(\mathcal{D}; \boldsymbol{\phi}))$, where \mathbf{z} are the hidden variables (latents and/or parameters), \mathcal{D} are the observed variables (data), f is an **inference network**, and $\boldsymbol{\phi}$ are its parameters; for details, see Section 10.3.6. Note that in this latter, setting the joint model $p(\mathbf{z}, \mathcal{D})$ may be a “traditional” model without any “neural” components. For example, it could be a complex simulator. Thus the DNN is just used for computational purposes, not statistical / modeling purposes.

More details on DNNs can be found in such books as [Zha+20a; Cho21; Gér19; GBC16], as well as a multitude of online courses. For a more theoretical treatment, see e.g., [Ber+21; Cal20; Aro+21; RY21].

16.2 Building blocks of differentiable circuits

In this section we discuss some common building blocks used in constructing neural networks. We denote the input to a block as \mathbf{x} and the output as \mathbf{y} .



11 *Figure 16.1: An articial “neuron”, the most basic building block of a DNN. (a) The output y is a weighted
12 combination of the inputs \mathbf{x} , where the weights vector is denoted by \mathbf{w} . (b) Alternative depiction of the
13 neuron’s behavior. The bias term b can be emulated by defining $w_N = b$ and $X_N = 1$.*

14

15

16 16.2.1 Linear layers

17

18 The most basic building block of a DNN is a single “**neuron**”, which corresponds to a real-valued
19 signal y computed by multiplying a vector-valued input signal \mathbf{x} by a weight vector \mathbf{w} , and then
20 adding a bias term b . That is,

21
$$y = f(\mathbf{x}; \theta) = \mathbf{w}^\top \mathbf{x} + b \tag{16.1}$$

22

23 where $\theta = (\mathbf{w}, b)$ are the parameters for the function f . This is depicted in Figure 16.1. (The bias
24 term is omitted for clarity.)

25 It is common to group a set of neurons together into a **layer**. We can then represent the activations
26 of a layer with D units as a vector $\mathbf{z} \in \mathbb{R}^D$. We can transform an input vector of activations \mathbf{x} into
27 an output vector \mathbf{y} by multiplying by a weight matrix \mathbf{W} , an adding an offset vector or bias term \mathbf{b}
28 to get

29
$$\mathbf{y} = f(\mathbf{x}; \theta) = \mathbf{W}\mathbf{x} + \mathbf{b} \tag{16.2}$$

30

31 where $\theta = (\mathbf{W}, \mathbf{b})$ are the parameters for the function f . This is called a **linear layer**, or **fully**
32 **connected layer**.

33 It is common to prepend the bias vector onto the first column of the weight matrix, and to append
34 a 1 to the vector \mathbf{x} , so that we can write this more compactly as $\mathbf{x} = \tilde{\mathbf{W}}^\top \tilde{\mathbf{x}}$, where $\tilde{\mathbf{W}} = [\mathbf{W}, \mathbf{b}]$ and
35 $\tilde{\mathbf{x}} = [\mathbf{x}, 1]$. This allows us to ignore the bias term from our notation if we want to.

36

37 16.2.2 Non-linearities

38 A stack of linear layers is equivalent to a single linear layer where we multiply together all the
39 weight matrices. To get more expressive power we can transform each layer by passing it elementwise
40 (pointwise) through a nonlinear function called an **activation function**. This is denoted by
41

42
$$\mathbf{y} = \varphi(\mathbf{x}) = [\varphi(x_1), \dots, \varphi(x_D)] \tag{16.3}$$

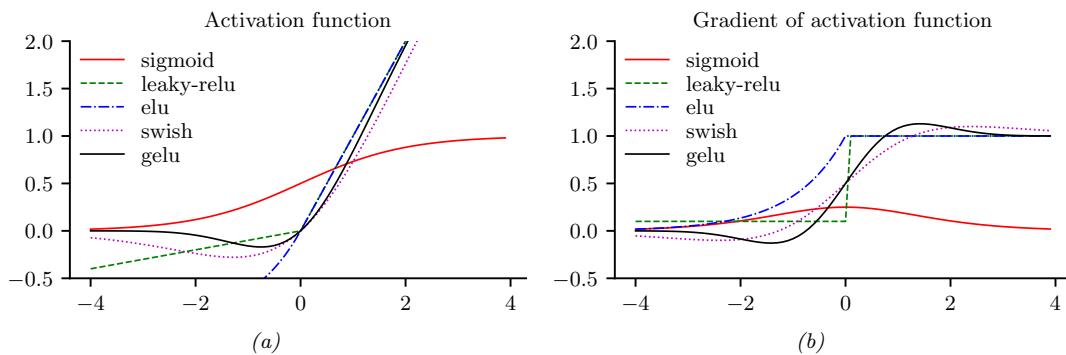
43

44 See Table 16.1 for a list of some common activation functions, and Figure 16.2 for a visualization.
45 For more details, see e.g., [Mur22, Sec 13.2.3].

46

Name	Definition	Range	Reference
Sigmoid	$\sigma(a) = \frac{1}{1+e^{-a}}$	$[0, 1]$	
Hyperbolic tangent	$\tanh(a) = 2\sigma(2a) - 1$	$[-1, 1]$	
Softplus	$\sigma_+(a) = \log(1 + e^a)$	$[0, \infty]$	[GBB11]
Rectified linear unit	$\text{ReLU}(a) = \max(a, 0)$	$[0, \infty]$	[GBB11; KSH12a]
Leaky ReLU	$\max(a, 0) + \alpha \min(a, 0)$	$[-\infty, \infty]$	[MHN13]
Exponential linear unit	$\max(a, 0) + \min(a(e^a - 1), 0)$	$[-\infty, \infty]$	[CUH16]
Swish	$a\sigma(a)$	$[-\infty, \infty]$	[RZL17]
GELU	$a\Phi(a)$	$[-\infty, \infty]$	[HG16]

Table 16.1: List of some popular activation functions for neural networks.

Figure 16.2: (a) Some popular activation functions. “ReLU” stands for “restricted linear unit”. “GELU” stands for “Gaussian error linear unit”. (b) Plot of their gradients. Generated by [activation_fun_deriv.ipynb](#).

16.2.3 Convolutional layers

When dealing with image data, we can apply the same weight matrix to each local patch of the image, in order to reduce the number of parameters. If we “slide” this weight matrix over the image and add up the results, we get a technique known as **convolution**; in this case the weight matrix is often called a “**kernel**” or “**filter**”.

More precisely, let $\mathbf{X} \in \mathbb{R}^{H \times W}$ be the input image, and $\mathbf{W} \in \mathbb{R}^{h \times w}$ be the kernel. The output is denoted by $\mathbf{Z} = \mathbf{X} \circledast \mathbf{W}$, where (ignoring boundary conditions) we have the following:¹

$$Z_{i,j} = \sum_{u=0}^{h-1} \sum_{v=0}^{w-1} x_{i+u, j+v} w_{u,v} \quad (16.4)$$

Essentially we compare a local patch of \mathbf{x} , of size $h \times w$ and centered at (i, j) , to the filter \mathbf{w} ; the output just measures how similar the input patch is to the filter. We can define convolution in 1d or 3d in an analogous manner. Note that the spatial size of the outputs may be smaller than inputs,

¹ 1. Note that, technically speaking, we are using **cross correlation** rather than convolution. However, these terms are used interchangeably in deep learning.

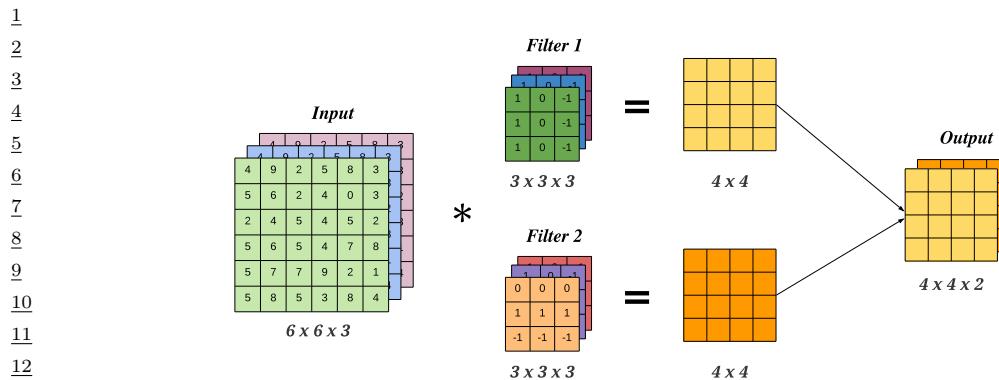


Figure 16.3: A 2d convolutional layer with 3 input channels and 2 output channels. The kernel has size 3×3 and we use stride 1 with 0 padding, so the the 6×6 input gets mapped to the 4×4 output.

due to boundary effects, although this can be solved by using **padding**. See [Mur22, Sec 14.2.1] for more details.

We can repeat this process for multiple layers of inputs, and by using multiple filters, we can generate multiple layers of output. In general, if we have C input channels, and we want to map it to D output (feature) channels, then we define D kernels, each of size $h \times w \times C$, where h, w are the height and width of the kernel. The d 'th output feature map is obtained by convolving all C input feature maps with the d 'th kernel, and then adding up the results elementwise:

$$z_{i,j,d} = \sum_{u=0}^{h-1} \sum_{v=0}^{w-1} \sum_{c=0}^{C-1} x_{i+u,j+v,c} w_{u,v,c,d} \quad (16.5)$$

This is called a **convolutional layer**, and is illustrated in Figure 16.3.

The advantage of a convolutional layer compared to using a linear layer is that the weights of the kernel are shared across locations in the input. Thus if a pattern in the input shifts locations, the corresponding output activation will also shift. This is called **shift equivariance**. In some cases, we want the output to be the same, no matter where the input pattern occurs; this is called **shift invariance**, and can be obtained by using a **pooling layer**, which computes the maximum or average value in each local patch of the input. (Note that pooling layers have no free (learnable) parameters.) Other forms of invariance can also be captured by neural networks (see e.g., [CW16; FWW21]).

16.2.4 Residual (skip) connections

If we stack a large number of nonlinear layers together, the signal may get squashed to zero or may blow up to infinity, depending on the magnitude of the weights, and the nature of the nonlinearities. Similar problems can plague gradients that are passed backwards through the network (see Section 6.2). To reduce the effect of this we can add **skip connections**, also called **residual connections**, which allow the signal to skip one or more layers, which prevents it from being modified. For example,

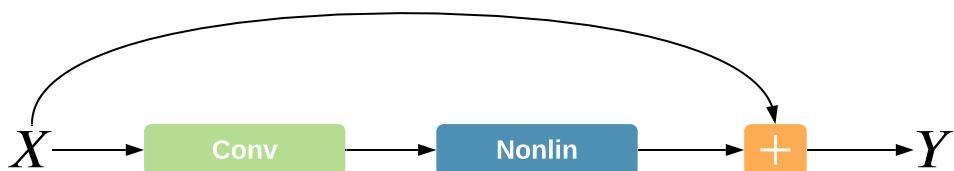


Figure 16.4: A residual connection around a convolutional layer.

Figure 16.4 illustrates a network that computes

$$\mathbf{y} = f(\mathbf{x}; \mathbf{W}) = \varphi(\text{conv}(\mathbf{x}; \mathbf{W})) + \mathbf{x} \quad (16.6)$$

Now the convolutional layer only needs to learn an offset or residual to add (or subtract) to the input to match the desired output, rather than predicting the output directly. Such residuals are often small in size, and hence are easier to learn using neurons with weights that are bounded (e.g., close to 1).

16.2.5 Normalization layers

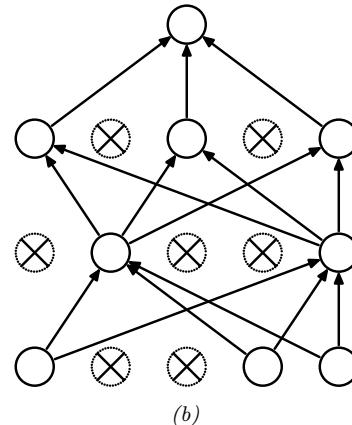
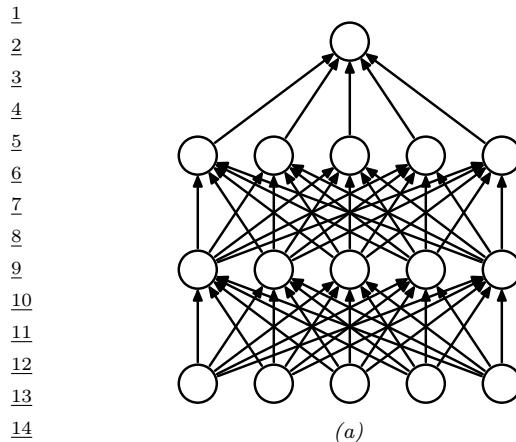
To learn an input-output mapping, it is often best if the inputs are standardized, meaning that they have zero mean and unit standard deviation. This ensures that the required magnitude of the weights is small, and comparable across dimensions. To ensure that the internal activations have this property, it is common to add **normalization layers**.

The most common approach is to use **batch normalization (BN)** [IS15]. However this relies on having access to a batch of $B > 1$ input examples. Various alternatives have been proposed to overcome the need of having an input batch, such as **layer normalization** [BKH16], **instance normalization** [UVL16], **group normalization** [WH18], **filter response normalization** [SK20], etc. More details can be found in [Mur22, Sec 14.2.4].

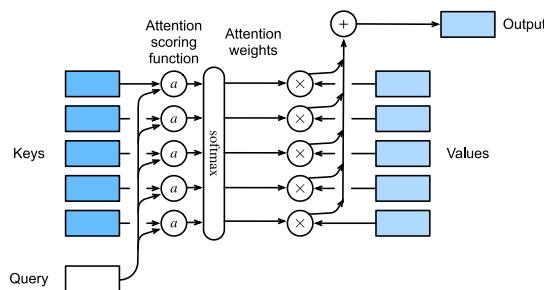
16.2.6 Dropout layers

Neural networks often have millions of parameters, and thus can sometimes overfit, especially when trained on small datasets. There are many ways to ameliorate this effect, such as applying regularizers to the weights, or adopting a fully Bayesian approach (see Chapter 17). Another common heuristic is known as **dropout** [Sri+14a], in which edges are randomly omitted each time the network is used, as illustrated in Figure 16.5. More precisely, if w_{lij} is the weight of the edge from node i in layer $l - 1$ to node j in layer $l + 1$, then we replace it with $\theta_{lij} = w_{lij}\epsilon_{li}$, where $\epsilon_{li} \sim \text{Ber}(1 - p)$, where p is the drop probability, and $1 - p$ is the keep probability. Thus if we sample $\epsilon_{li} = 0$, then all of the weights going out of unit i in layer $l - 1$ into any j in layer l will be set to 0.

During training, the gradients will be zero for the weights connected to a neuron which has been switched “off”. However, since we resample ϵ_{lij} every time the network is used, different combinations of weights will be updated on each step. The result is an **ensemble** of networks, each with slightly different sparse graph structures.



16 *Figure 16.5: Illustration of dropout.* (a) A standard neural net with 2 hidden layers. (b) An example of a
17 thinned net produced by applying dropout with $p = 0.5$. Units that have been dropped out are marked with an
18 x . From Figure 1 of [Sri+14a]. Used with kind permission of Geoff Hinton.



29 *Figure 16.6: Attention layer.* (a) Mapping a single query \mathbf{q} to a single output, given a set of keys and values.
30 From Figure 10.3.1 of [Zha+20a]. Used with kind permission of Aston Zhang.

31
32
33
34 At test time, we usually turn the dropout noise off, so the model acts deterministically. To ensure
35 the weights have the same expectation at test time as they did during training (so the input activation
36 to the neurons is the same, on average), at test time we should use $\mathbb{E}[\theta_{lij}] = w_{lij}\mathbb{E}[\epsilon_{li}]$. For Bernoulli
37 noise, we have $\mathbb{E}[\epsilon] = 1 - p$, so we should multiply the weights by the keep probability, $1 - p$, before
38 making predictions. We can, however, use dropout at test time if we wish. This is called **Monte
39 Carlo dropout** (see Section 17.3.1).
40

41

42 16.2.7 Attention layers

43
44 In non-parametric kernel based prediction methods, such as Gaussian processes (Chapter 18), we
45 compare the input $\mathbf{x} \in \mathbb{R}^{d_k}$ to each of the training examples $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ using a kernel to get
46 a vector of similarity scores, $\boldsymbol{\alpha} = [\mathcal{K}(\mathbf{x}, \mathbf{x}_i)]_{i=1}^n$. We then use this to retrieve a weighted combination
47

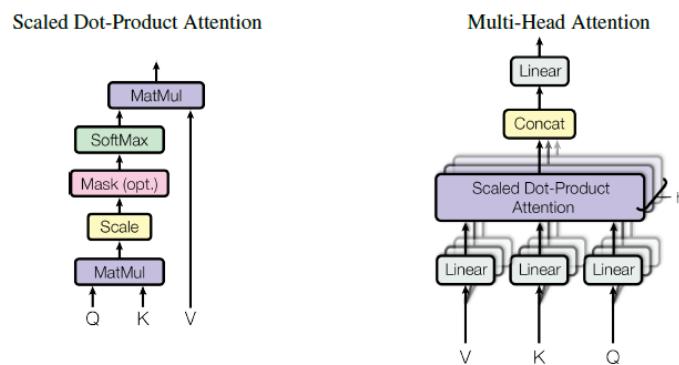


Figure 16.7: (a) Scaled dot-product attention in matrix form. (b) Multi-head attention. From Figure 2 of [Vas+17b]. Used with kind permission of Ashish Vaswani.

of the corresponding m target values $\mathbf{y}_i \in \mathbb{R}^{d_v}$ as follows:

$$\hat{\mathbf{y}} = \sum_{i=1}^n \alpha_i \mathbf{y}_i \quad (16.7)$$

See Section 18.3.7 for details.

We can make a differentiable and parametric version of this as follows (see [Tsa+19] for details). First we replace the stored examples matrix \mathbf{X} with a learned embedding, to create a set of stored **keys**, $\mathbf{K} = \mathbf{W}^K \mathbf{X} \in \mathbb{R}^{n \times d_k}$. Similarly we replace the stored output matrix \mathbf{Y} with a learned embedding, to create a set of stored **values**, $\mathbf{V} = \mathbf{W}^V \mathbf{Y} \in \mathbb{R}^{n \times d_v}$. Finally we embed the input to create a **query**, $\mathbf{q} = \mathbf{W}^Q \mathbf{x} \in \mathbb{R}^{d_k}$. The parameters to be learned are the three embedding matrices.

To ensure the output is a differentiable function of the input, we replace the fixed kernel function with a soft **attention layer**. More precisely, we define

$$\text{Attn}(\mathbf{q}, (\mathbf{k}_1, \mathbf{v}_1), \dots, (\mathbf{k}_n, \mathbf{v}_n)) = \text{Attn}(\mathbf{q}, (\mathbf{k}_{1:n}, \mathbf{v}_{1:n})) = \sum_{i=1}^n \alpha_i(\mathbf{q}, \mathbf{k}_{1:n}) \mathbf{v}_i \quad (16.8)$$

where $\alpha_i(\mathbf{q}, \mathbf{k}_{1:n})$ is the i 'th **attention weight**; these weights satisfy $0 \leq \alpha_i(\mathbf{q}, \mathbf{k}_{1:n}) \leq 1$ for each i and $\sum_i \alpha_i(\mathbf{q}, \mathbf{k}_{1:n}) = 1$.

The attention weights can be computed from an **attention score** function $a(\mathbf{q}, \mathbf{k}_i) \in \mathbb{R}$, that computes the similarity of query \mathbf{q} to key \mathbf{k}_i . For example, we can use (scaled) **dot product attention**, which has the form

$$a(\mathbf{q}, \mathbf{k}) = \mathbf{q}^\top \mathbf{k} / \sqrt{d_k} \quad (16.9)$$

(The scaling by $\sqrt{d_k}$ is to reduce the dependence of the output on the dimensionality of the vectors.) Given the scores, we can compute the attention weights using the softmax function:

$$\alpha_i(\mathbf{q}, \mathbf{k}_{1:n}) = \text{softmax}_i([a(\mathbf{q}, \mathbf{k}_1), \dots, a(\mathbf{q}, \mathbf{k}_n)]) = \frac{\exp(a(\mathbf{q}, \mathbf{k}_i))}{\sum_{j=1}^n \exp(a(\mathbf{q}, \mathbf{k}_j))} \quad (16.10)$$

1 See Figure 16.6 for an illustration.
2

3 In some cases, we want to restrict attention to a subset of the dictionary, corresponding to valid
4 entries. For example, we might want to pad sequences to a fixed length (for efficient minibatching),
5 in which case we should “mask out” the padded locations. This is called **masked attention**. We
6 can implement this efficiently by setting the attention score for the masked entries to a large negative
7 number, such as -10^6 , so that the corresponding softmax weights will be 0.

8 In practice, we usually deal with minibatches of n vectors at a time. Let the corresponding matrices
9 of queries, keys and values be denoted by $\mathbf{Q} \in \mathbb{R}^{n \times d_k}$, $\mathbf{K} \in \mathbb{R}^{n \times d_k}$, $\mathbf{V} \in \mathbb{R}^{n \times d_v}$. Let

10

$$\mathbf{z}_j = \sum_{i=1}^n \alpha_i(\mathbf{q}_j, \mathbf{K}) \mathbf{v}_i \quad (16.11)$$

11

12 be the j 'th output corresponding to the j 'th query. We can compute all outputs $\mathbf{Z} \in \mathbb{R}^{n \times d_v}$ in
13 parallel using

14

$$\mathbf{Z} = \text{Attn}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}\right)\mathbf{V} \quad (16.12)$$

15

16 where the softmax function softmax is applied row-wise. See Figure 16.7(left) for an illustration.

17 To increase the flexibility of the model, we often use a **multi-head attention** layer, as illustrated
18 in Figure 16.7(right). Let the i 'th head be

19

$$\mathbf{h}_i = \text{Attn}(\mathbf{Q}\mathbf{W}_i^Q, \mathbf{K}\mathbf{W}_i^K, \mathbf{V}\mathbf{W}_i^V) \quad (16.13)$$

20

21 where $\mathbf{W}_i^Q \in \mathbb{R}^{d \times d_k}$, $\mathbf{W}_i^K \in \mathbb{R}^{d \times d_k}$ and $\mathbf{W}_i^V \in \mathbb{R}^{d \times d_v}$ are linear projection matrices. We define the
22 output of the MHA layer to be

23

$$\mathbf{Z} = \text{MHA}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Concat}(\mathbf{h}_1, \dots, \mathbf{h}_h)\mathbf{W}^O \quad (16.14)$$

24

25 where h is the number of heads, and $\mathbf{W}^O \in \mathbb{R}^{hd_v \times d}$. Having multiple heads can increase performance
26 of the layer, in the event that some of the weight matrices are poorly initialized; after training, we
27 can often remove all but one of the heads [MLN19].

28 When the output of one attention layer is used as input to another, the method is called **self-
29 attention**. This is the basis of the transformer model, which we discuss in Section 16.3.5.

30 16.2.8 Recurrent layers

31 We can make the model be **stateful** by augmenting the input \mathbf{x} with the current state \mathbf{s}_t , and then
32 computing the output and the new state using some kind of function:

33

$$(\mathbf{y}, \mathbf{s}_{t+1}) = f(\mathbf{x}, \mathbf{s}_t) \quad (16.15)$$

34

35 This is called a **recurrent layer**, as shown in Figure 16.8. This forms the basis of **recurrent neural
36 networks**, discussed in Section 16.3.4. In a vanilla RNN, the function f is a simple MLP, but it
37 may also use attention (Section 16.2.7).

38

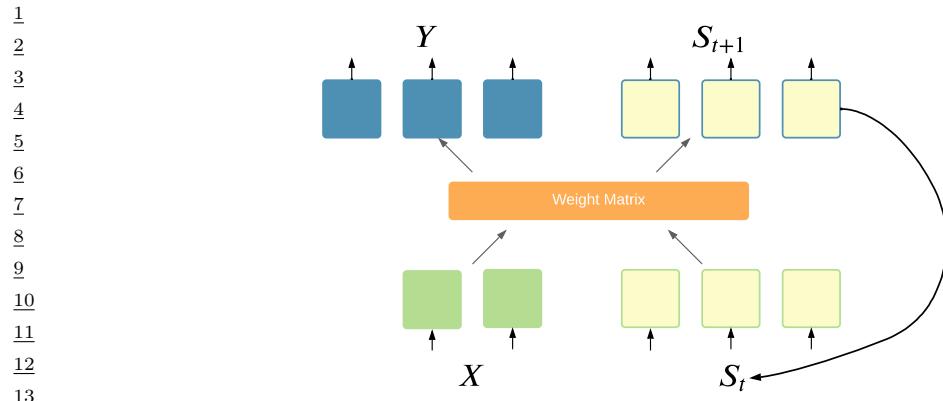


Figure 16.8: Recurrent layer.

16.2.9 Multiplicative layers

In this section, we discuss **multiplicative layers**, which are useful for combining different information sources. Our presentation follows [Jay+20].

Suppose we have inputs $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{z} \in \mathbb{R}^m$. In a linear layer (and, by extension, convolutional layers), it is common to concatenate the inputs to get $f(\mathbf{x}, \mathbf{z}) = \mathbf{W}[\mathbf{x}; \mathbf{z}] + \mathbf{b}$, where $\mathbf{W} \in \mathbb{R}^{k \times (m+n)}$ and $\mathbf{b} \in \mathbb{R}^k$. We can increase the expressive power of the model by using **multiplicative interactions**, such as the following **bilinear form**:

$$f(\mathbf{x}, \mathbf{z}) = \mathbf{z}^\top \mathbb{W} \mathbf{x} + \mathbf{U} \mathbf{z} + \mathbf{V} \mathbf{x} + \mathbf{b} \quad (16.16)$$

where $\mathbb{W} \in \mathbb{R}^{m \times n \times k}$ is a weight tensor, defined such that

$$(\mathbf{z}^\top \mathbb{W} \mathbf{x})_k = \sum_{ij} z_i \mathbb{W}_{ijk} x_j \quad (16.17)$$

That is, the k 'th entry of the output is the weighted inner product of \mathbf{z} and \mathbf{x} , where the weight matrix is the k 'th “slice” of \mathbb{W} . The other parameters have size $\mathbf{U} \in \mathbb{R}^{k \times m}$, $\mathbf{V} \in \mathbb{R}^{k \times n}$, and $\mathbf{b} \in \mathbb{R}^k$.

This formulation includes many interesting special cases. In particular, a **hypernetwork** [HDL17] can be viewed in this way. A hypernetwork is a neural network that generates parameters for another neural network. In particular, we replace $f(\mathbf{x}; \boldsymbol{\theta})$ with $f(\mathbf{x}; g(\mathbf{z}; \boldsymbol{\phi}))$. If f and g are affine, this is equivalent to a multiplicative layer. To see this, let $\mathbf{W}' = \mathbf{z}^\top \mathbb{W} + \mathbf{V}$ and $\mathbf{b}' = \mathbf{U} \mathbf{z} + \mathbf{b}$. If we define $g(\mathbf{z}; \boldsymbol{\Phi}) = [\mathbf{W}', \mathbf{b}']$, and $f(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{W}' \mathbf{x} + \mathbf{b}'$, we recover Equation (16.16).

We can also view the gating layers used in RNNs (Section 16.3.4) as a form of multiplicative interaction. In particular, if we the hypernetwork computes the diagonal matrix $\mathbf{W}' = \sigma(\mathbf{z}^\top \mathbb{W} + \mathbf{V}) = \text{diag}(a_1, \dots, a_n)$, then we can define $f(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) = \mathbf{a}(\mathbf{z}) \odot \mathbf{x}$, which is the standard gating mechanism. Attention mechanisms (Section 16.2.7) are also a form of multiplicative interaction, although they involve three-way interactions, between query, key and value.

Another variant arises if the hypernetwork just computes a scalar weight for each channel of a convolutional layer, plus a bias term:

$$f(\mathbf{x}, \mathbf{z}) = \mathbf{a}(\mathbf{z}) \odot \mathbf{x} + \mathbf{b}(\mathbf{z}) \quad (16.18)$$

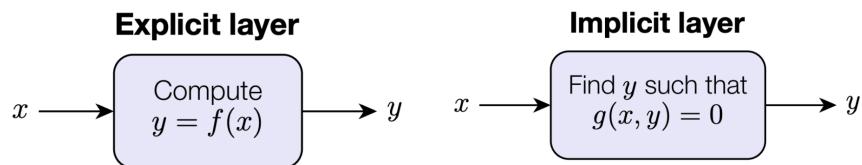


Figure 16.9: Explicit vs implicit layers.

This is called **FiLM**, which stands for “Feature-wise Linear Modulation” [Per+18]. For a detailed tutorial on the FiLM layer and its many applications, see <https://distill.pub/2018/feature-wise-transformations>.

16.2.10 Implicit layers

So far we have focused on **explicit layers**, which specify how to transform the input to the output using $y = f(x)$. We can also define **implicit layers**, which specify the output indirectly, in terms of a constraint function:

$$\mathbf{y} \in \underset{\mathbf{y}}{\operatorname{argmin}} f(\mathbf{x}, \mathbf{y}) \text{ such that } g(\mathbf{x}, \mathbf{y}) = 0 \quad (16.19)$$

The details on how to find a solution to this constrained optimization problem can vary depending on the problem. For example, we may need to run an inner optimization routine, or call a differential equation solver. The main advantage of this approach is that the inner computations do not need to be stored explicitly, which saves a lot of memory. Furthermore, once the solution has been found, we can propagate gradients through the whole layer, by leveraging the implicit function theorem. This lets us use higher level primitives inside an end-to-end framework. For more details, see [GHC21] and <http://implicit-layers-tutorial.org/>.

16.3 Canonical examples of neural networks

In this section, we give several “canonical” examples of neural network architectures that are widely used for different tasks.

16.3.1 Multi-layer perceptrons (MLP)

A **multi-layer perceptron (MLP)**, also called a **feedforward neural network (FFNN)**, is one of the simplest kinds of neural networks. It consists of a series of L linear layers, combined with elementwise nonlinearities:

$$f(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{W}_L \varphi_L (\mathbf{W}_{L-1} \varphi_{L-1} (\cdots \varphi_1 (\mathbf{W}_1 \mathbf{x}) \cdots)) \quad (16.20)$$

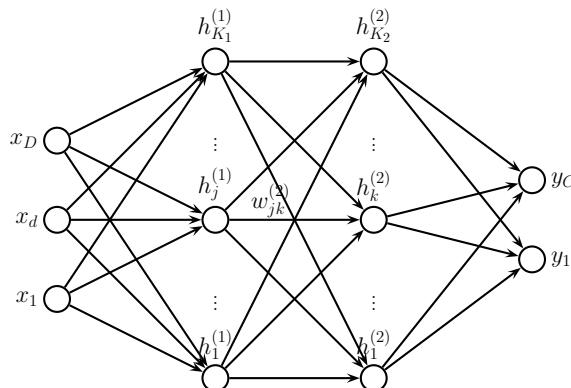


Figure 16.10: A feedforward neural network with D inputs, K_1 hidden units in layer 1, K_2 hidden units in layer 2, and C outputs. $w_{jk}^{(l)}$ is the weight of the connection from node j in layer $l - 1$ to node k in layer l .

For example, Figure 16.10 shows an MLP with 1 input layer of D units, 2 hidden layers of K_1 and K_2 units, and 1 output layer with C units. The k 'th hidden unit in layer l is given by

$$h_k^{(l)} = \varphi_l \left(b_k^{(l)} + \sum_{j=1}^{K_{l-1}} w_{jk}^{(l)} h_j^{(l-1)} \right) \quad (16.21)$$

where φ_l is the nonlinear activation function at layer l .

For a classification problem, the final nonlinearity is usually the softmax function. However, it is also common for the final layer to have linear activations, in which case the outputs are interpreted as logits; the loss function used during training then converts to (log) probabilities internally.

We can also use MLPs for regression. Figure 16.11 shows how we can make a model for **heteroskedastic** nonlinear regression. (The term “heteroskedastic” just means that the predicted output variance is input-dependent, rather than a constant.) This function has two outputs which compute $f_\mu(\mathbf{x}) = \mathbb{E}[y|\mathbf{x}, \boldsymbol{\theta}]$ and $f_\sigma(\mathbf{x}) = \sqrt{\mathbb{V}[y|\mathbf{x}, \boldsymbol{\theta}]}$. We can share most of the layers (and hence parameters) between these two functions by using a common “backbone” and two output “heads”, as shown in Figure 16.11. For the μ head, we use a linear activation, $\varphi(a) = a$. For the σ head, we use a softplus activation, $\varphi(a) = \sigma_+(a) = \log(1 + e^a)$. If we use linear heads and a nonlinear backbone, the overall model is given by

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(y|w_\mu^\top f(\mathbf{x}; \mathbf{w}_{\text{shared}}), \sigma_+(w_\sigma^\top f(\mathbf{x}; \mathbf{w}_{\text{shared}}))) \quad (16.22)$$

16.3.2 Convolutional neural networks (CNN)

A vanilla **convolutional neural network** or **CNN** consists of a series of convolutional layers, pooling layers, linear layers, and nonlinearities. See Figure 16.12 for an example. More sophisticated architectures, such as the **ResNet** model [He+16a; He+16b], add skip (residual) connections, normalization layers, etc. The **ConvNeXt** model of [Liu+22] is considered the current (as of February 2022) state of the art CNN architecture for a wide variety of vision tasks. See e.g., [Mur22, Ch.14] for more details on CNNs.

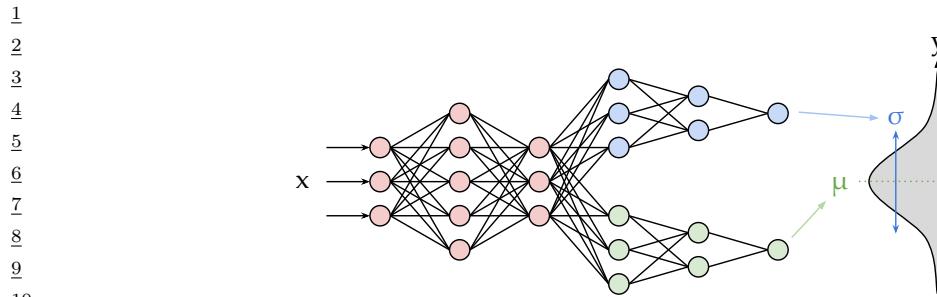


Figure 16.11: Illustration of an MLP with a shared “backbone” and two output “heads”, one for predicting the mean and one for predicting the variance. From <https://brendanhasz.github.io/2019/07/23/bayesian-density-net.html>. Used with kind permission of Brendan Hasz.

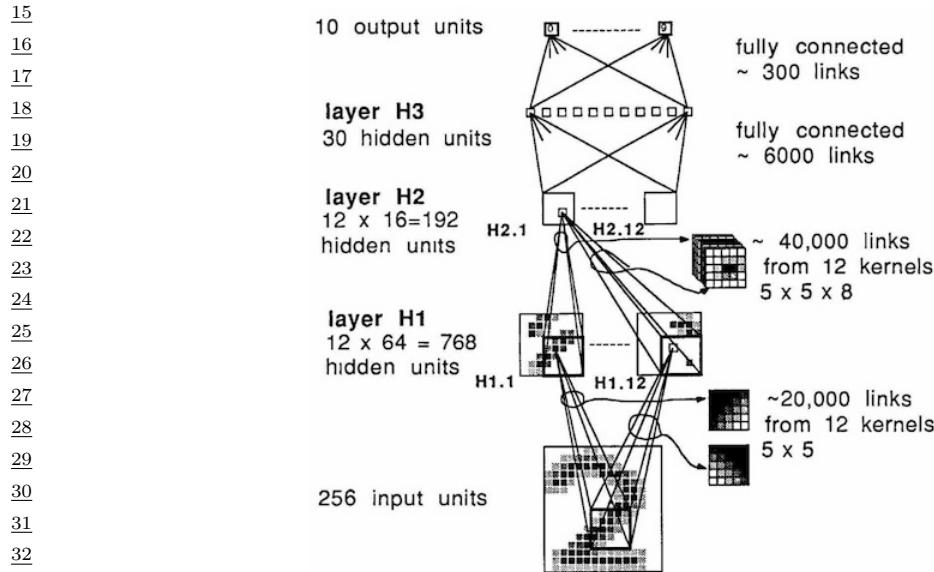


Figure 16.12: One of the first CNNs ever created, for classifying MNIST images. From Figure 3 of [LeC+89]. For a “modern” implementation, see lecun1989.ipynb.

16.3.3 Autoencoders

An **autoencoder** is a neural network that maps inputs \mathbf{x} to a low-dimensional latent space using an **encoder**, $\mathbf{z} = f_e(\mathbf{x})$, and then attempts to reconstruct the inputs using a **decoder**, $\hat{\mathbf{x}} = f_d(\mathbf{z})$. The model is trained to minimize

$$\mathcal{L}(\boldsymbol{\theta}) = \|\mathbf{r}(\mathbf{x}) - \mathbf{x}\|_2^2 \quad (16.23)$$

where $\mathbf{r}(\mathbf{x}) = f_d(f_e(\mathbf{x}))$. (We can also replace squared error with more general conditional log likelihoods.) See Figure 16.13 for an illustration of a 3 layer AE.

47

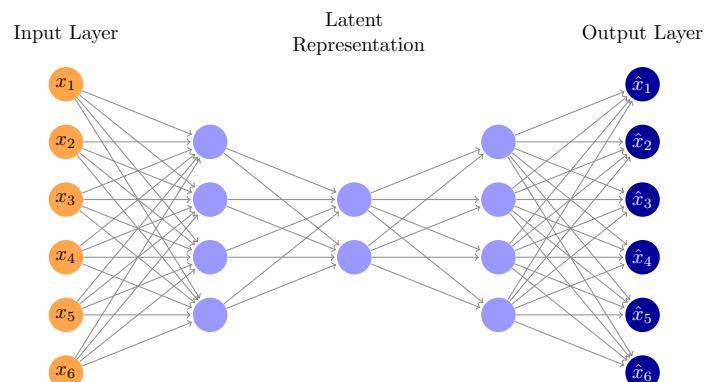
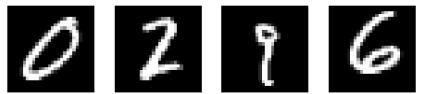


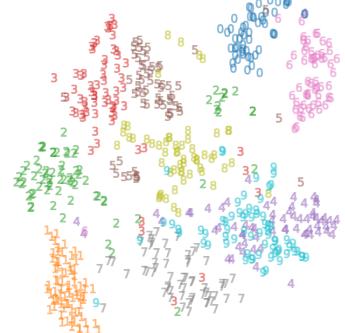
Figure 16.13: Illustration of an autoencoder with 3 hidden layers.



(a)



(b)



(c)

Figure 16.14: (a) Some MNIST digits. (b) Reconstruction of these images using a convolutional autoencoder. (c) t-SNE visualization of the 20-d embeddings. The colors correspond to class labels, which were not used during training. Generated by `ae_mnist_conv_jax.ipynb`.

For image data, we can make the encoder be a convolutional network, and the decoder be a transpose convolutional network. We can use this to compute low dimensional embeddings of image data. For example, suppose we fit such a model to some MNIST digits. We show the reconstruction abilities of such a model in Figure 16.14b. In Figure 16.14c, we show a 2d visualization of the 20-dimensional embedding space computed using t-SNE. The colors correspond to class labels, which were not used during training. We see fairly good separation, showing that images which are visually similar are placed close to each other in the embedding space, as desired. (See also Section 21.2.6, where we compare AEs with variational AEs.)

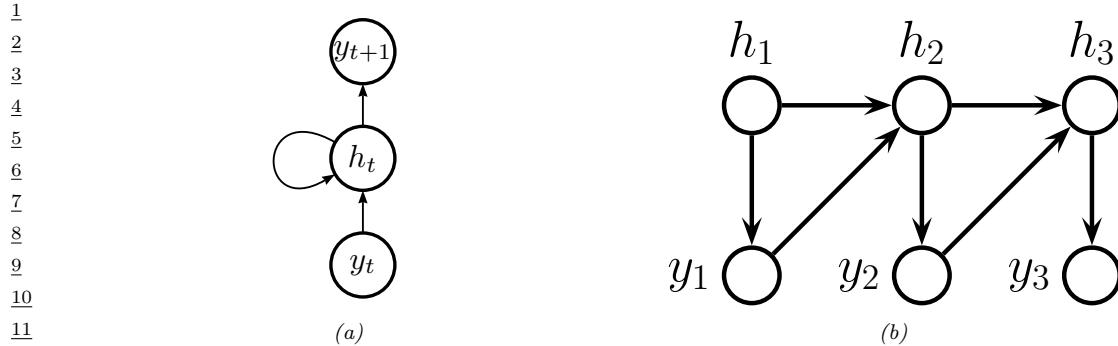


Figure 16.15: Illustration of a recurrent neural network (RNN). (a) With self-loop. (b) Unrolled in time.

16.3.4 Recurrent neural networks (RNN)

A **recurrent neural network (RNN)** is a network with a recurrent layer, as in Equation (16.15). This is illustrated in Figure 16.15. Formally this defines the following probability distribution over sequences:

$$p(\mathbf{y}_{1:T}) = \sum_{\mathbf{h}_{1:T}} p(\mathbf{y}_{1:T}, \mathbf{h}_{1:T}) = \sum_{\mathbf{h}_{1:T}} \mathbb{I}(\mathbf{h}_1 = \mathbf{h}_1^*) p(\mathbf{y}_1 | \mathbf{h}_1) \prod_{t=2}^T p(\mathbf{y}_t | \mathbf{h}_t) \mathbb{I}(\mathbf{h}_t = f(\mathbf{h}_{t-1}, \mathbf{y}_{t-1})) \quad (16.24)$$

where \mathbf{h}_t is the deterministic hidden state, computed from the last hidden state and last output using $f(\mathbf{h}_{t-1}, \mathbf{y}_{t-1})$. (At training time, \mathbf{y}_{t-1} is observed, but at prediction time, it is generated.)

In a vanilla RNN, the function f is a simple MLP. However, we can also use attention to selectively update parts of the state vector based on similarity between the input the previous state, as in the **GRU** (gated recurrent unit) model, and the **LSTM** (long short term memory) model. We can also make the model into a conditional sequence model, by feeding in extra inputs to the f function. See e.g., [Mur22, Ch. 15] for more details on RNNs.

16.3.5 Transformers

Consider the problem of classifying each word in a sentence, for example with its part of speech tag (noun, verb, etc). That is, we want to learn a mapping $f : \mathcal{X} \rightarrow \mathcal{Y}$, where $\mathcal{X} = \mathcal{V}^T$ is the set of input sequences defined over (word) vocabulary \mathcal{V} , T is the length of the sentence, and $\mathcal{Y} = \mathcal{T}^T$ is the set of output sequences, defined over (tag) vocabulary \mathcal{T} . To do well at this task, we need to learn a contextual embedding of each word. RNNs process one token at a time, so the embedding of the word at location t , \mathbf{z}_t , depends on the hidden state of the network, \mathbf{s}_t , which may be a lossy summary of all the previously seen words. We can create bidirectional RNNs so that future words can also affect the embedding of \mathbf{z}_t , but this dependence is still mediated via the hidden state. An alternative approach is to compute \mathbf{z}_t as a direct function of all the other words in the sentence, by using the attention operator discussed in Section 16.2.7 rather than using hidden state. This is called an (encoder-only) **transformer**, and is used by models such as BERT [Dev+19]. This idea is sketched in Figure 16.16.

47

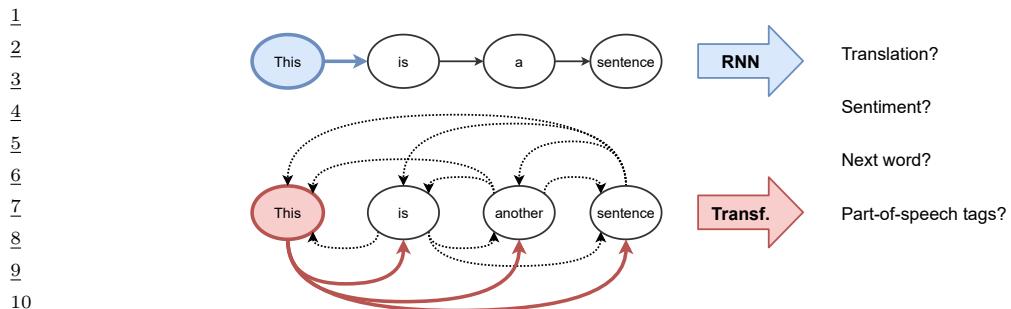


Figure 16.16: Visualizing the difference between an RNN and a transformer. From [Jos20]. Used with kind permission of Chaitanya Joshi.

It is also possible to create a decoder-only transformer, in which each output y_t only attends to all the previously generated outputs, $y_{1:t-1}$. This can be implemented using masked attention, and is useful for generative language models, such as GPT (see Section 22.4.1). We can combine the encoder and decoder to create a conditional sequence-to-sequence model, $p(y_{1:T_y} | x_{1:T_x})$, as proposed in the original transformer paper [Vas+17c]. See Supplementary Section 16.1.1 for the details.

It has been found that large transformers are very flexible sequence-to-sequence function approximators, if trained on enough data (see e.g., [Lin+21a] for a review in the context of NLP, and [Kha+21; Han+20; Zan21] for reviews in the context of computer vision). The reasons why they work so well are still not very clear. However, some initial insights can be found in e.g., [Rag+21; WGY21; Nel21; BP21]. See also Supplementary Section 16.1.2.5 where we discuss the connection with graph neural networks.

16.3.6 Graph neural networks (GNNs)

It is possible to define neural networks for working with graph-structured data. These are called **graph neural networks** or **GNNs**. See Supplementary Section 16.1.2 for details.

17 Bayesian neural networks

This chapter is coauthored with Andrew Wilson.

17.1 Introduction

Deep neural networks (DNNs) are usually trained using a (penalized) maximum likelihood objective to find a single setting of parameters. However, large flexible models like neural networks can represent many functions, corresponding to different parameter settings, which fit the training data well, yet generalize in different ways, a phenomenon known as **underspecification** (see e.g., [D'A+20], and Figure 17.10 for an illustration). Considering all of these different models together can lead to improved accuracy and uncertainty representation. This can be done by computing the posterior predictive distribution using Bayesian model averaging:

$$p(\mathbf{y}|\mathbf{x}, \mathcal{D}) = \int p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})p(\boldsymbol{\theta}|\mathcal{D})d\boldsymbol{\theta} \quad (17.1)$$

where $p(\boldsymbol{\theta}|\mathcal{D}) \propto p(\boldsymbol{\theta})p(\mathcal{D}|\boldsymbol{\theta})$.

The main challenges in applying Bayesian inference to DNNs are specifying suitable priors, and efficiently computing the posterior, which is challenging due to the large number of parameters and the large datasets. The application of Bayesian inference to DNNs is sometimes called **Bayesian deep learning** or **BDL**. By contrast, the term **deep Bayesian learning** or **DBL** refers to the use of deep models to help speedup Bayesian inference of “classical” models, usually by training amortized inference networks that can be used as part of a variational inference or importance sampling algorithm, as discussed in Section 10.3.6.) For more details on the topic of BDL, see e.g., [PS17; Wil20; WI20; Jos+22; Kha20].

17.2 Priors for BNNs

To perform Bayesian inference for the parameters of a DNN, we need to specify a prior $p(\boldsymbol{\theta})$. [Nal18; WI20; For21] discusses the issue of prior selection at length. Here we just give a brief summary of common approaches.

1 2 **17.2.1 Gaussian priors**

3 Consider an MLP with one hidden layer with activation function φ and a linear output:

4 5
$$f(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{W}_2 \varphi(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2 \quad (17.2)$$

6 (If the output is nonlinear, such as a softmax transform, we can fold it into the loss function during
7 training.) If we have two hidden layers this becomes

8 9
$$f(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{W}_3 (\varphi(\mathbf{W}_2 \varphi(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2)) + \mathbf{b}_3 \quad (17.3)$$

10 In general, with $L - 1$ hidden layers and a linear output, we have

11 12
$$f(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{W}_L (\cdots \varphi(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1)) + \mathbf{b}_L \quad (17.4)$$

13 We need to specify the priors for \mathbf{W}_l and \mathbf{b}_l for $l = 1 : L$. The most common choice is to use a
14 factored Gaussian prior:

15 16
$$\mathbf{W}_\ell \sim \mathcal{N}(\mathbf{0}, \alpha_\ell^2 \mathbf{I}), \mathbf{b}_\ell \sim \mathcal{N}(\mathbf{0}, \beta_\ell^2 \mathbf{I}) \quad (17.5)$$

17 The **Xavier initialization** or **Glorot initialization**, named after the first author of [GB10], is to
18 set

19 20 21
$$\alpha_\ell^2 = \frac{2}{n_{\text{in}} + n_{\text{out}}} \quad (17.6)$$

22 where n_{in} is the fan-in of a node in level ℓ (number of weights coming into a neuron), and n_{out} is
23 the fan-out (number of weights going out of a neuron). **LeCun initialization**, named after Yann
24 LeCun, corresponds to using

25 26 27
$$\alpha_\ell^2 = \frac{1}{n_{\text{in}}} \quad (17.7)$$

28 We can get a better understanding of these priors by considering the effect they have on the
29 corresponding distribution over functions that they define. To help understand this correspondence,
30 let us reparameterize the model as follows:

31 32
$$\mathbf{W}_\ell = \alpha_\ell \boldsymbol{\eta}_\ell, \boldsymbol{\eta}_\ell \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \mathbf{b}_\ell = \beta_\ell \boldsymbol{\epsilon}_\ell, \boldsymbol{\epsilon}_\ell \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (17.8)$$

33 Hence every setting of the prior hyperparameters specifies the following random function:

34 35
$$f(\mathbf{x}; \boldsymbol{\alpha}, \boldsymbol{\beta}) = \alpha_L \boldsymbol{\eta}_L (\cdots \varphi(\alpha_1 \boldsymbol{\eta}_1 \mathbf{x} + \beta_1 \boldsymbol{\epsilon}_1)) + \beta_L \boldsymbol{\epsilon}_L \quad (17.9)$$

36 To get a feeling for the effect of these hyperparameters, we can sample MLP parameters from this
37 prior and plot the resulting random functions. We use a sigmoid nonlinearity, so $\varphi(a) = \sigma(a)$. We
38 consider $L = 2$ layers, so \mathbf{W}_1 are the input-to-hidden weights, and \mathbf{W}_2 are the hidden-to-output
39 weights. We assume the input and output are scalars, so we are generating random nonlinear 1d
40 mappings $f : \mathbb{R} \rightarrow \mathbb{R}$.

41 Figure 17.1(a) shows some sampled functions where $\alpha_1 = 5$, $\beta_1 = 1$, $\alpha_2 = 1$, $\beta_2 = 1$. In
42 Figure 17.1(b) we increase α_1 ; this allows the first layer weights to get bigger, making the sigmoid-like
43 mapping wider.

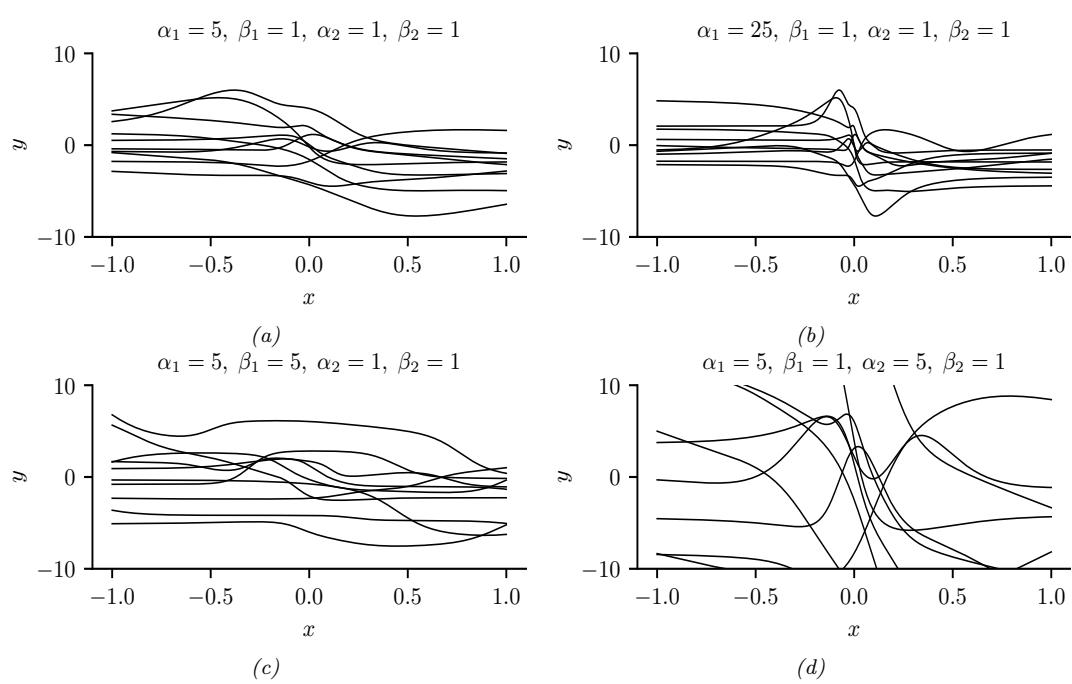


Figure 17.1: The effects of changing the hyperparameters on an MLP with one hidden layer. (a) Random functions sampled from a Gaussian prior with hyperparameters $\alpha_1 = 5, \beta_1 = 1, \alpha_2 = 1, \beta_2 = 1$. (b) Increasing α_1 by factor of 5. (c) Increasing β_1 by factor of 5. (d) Increasing α_2 by factor of 5. Generated by `mlp_priors_demo.ipynb`.

shape of the functions steeper. In Figure 17.1(c), we increase β_1 ; this allows the first layer biases to get bigger, which allows the center of the sigmoid to shift left and right more, away from the origin. In Figure 17.1(d), we increase α_2 ; this allows the second layer linear weights to get bigger, making the functions more “wiggly” (greater sensitivity to change in the input, and hence larger dynamic range).

The above results are specific to the case of sigmoidal activation functions. ReLU units can behave differently. For example, [WI20, App. E] show that for MLPs with ReLU units, if we set $\beta_\ell = 0$, so the bias terms are all zero, the effect of changing α_ℓ is just to rescale the output. To see this, note that Equation (17.9) simplifies to

$$f(\mathbf{x}; \boldsymbol{\alpha}, \boldsymbol{\beta} = \mathbf{0}) = \alpha_L \boldsymbol{\eta}_L(\cdots \varphi(\alpha_1 \boldsymbol{\eta}_1 \mathbf{x})) = \alpha_L \cdots \alpha_1 \boldsymbol{\eta}_L(\cdots \varphi(\boldsymbol{\eta}_1 \mathbf{x})) \quad (17.10)$$

$$= \alpha_L \cdots \alpha_1 f(\mathbf{x}; (\boldsymbol{\alpha} = \mathbf{1}, \boldsymbol{\beta} = \mathbf{0})) \quad (17.11)$$

where we used the fact that for ReLU, $\varphi(\alpha z) = \alpha \varphi(z)$ for any positive α , and $\varphi(\alpha z) = 0$ for any negative α (since the pre-activation $z \geq 0$). In general, it is the ratio of α and β that matters for determining what happens to input signals as they propagate forwards and backwards through a randomly initialized model; for details, see e.g., [Bah+20].

We see that initializing the model’s parameters at a particular random value is like sampling a

1 point from this prior over functions. In the limit of infinitely wide neural networks, we can derive
2 this prior distribution analytically: this is known as a **neural network Gaussian process**, and is
3 explained in Section 18.7.
4

56

7 17.2.2 Sparsity-promoting priors

8 Although Gaussian priors are simple and widely used, they are not the only option. For some
9 applications, it is useful to use **sparsity promoting priors**, such as the Laplace, which encourage
10 most of the weights (or channels in a CNN) to be zero (c.f., Section 15.2.5). For details, see [Hoe+21].
11

12

13 17.2.3 Learning the prior

14

15 We have seen how different priors for the parameters correspond to different priors over functions.
16 We could in principle set the hyperparameters (e.g., the α and β parameters of the Gaussian prior)
17 using grid search to optimize cross-validation loss. However, cross-validation can be slow, particularly
18 if we allow different priors for each layer of the network, as our grid search will grow exponentially
19 with the number of hyperparameters we wish to determine.

20 An alternative is to use gradient based methods to optimize the marginal likelihood
21

$$\log p(\mathcal{D}|\boldsymbol{\alpha}, \boldsymbol{\beta}) = \int \log p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta}|\boldsymbol{\alpha}, \boldsymbol{\beta})d\boldsymbol{\theta} \quad (17.12)$$

24

25 This approach is known as empirical Bayes (Section 3.8) or **evidence maximization**, since
26 $\log p(\mathcal{D}|\boldsymbol{\alpha}, \boldsymbol{\beta})$ is also called the evidence [Mac92a; WS93; Mac99b]. This can give rise to sparse models,
27 as we discussed in the context of automatic relevancy determination (Section 15.2.7). Unfortunately,
28 computing the marginal likelihood is computationally difficult for large neural networks.

29 Learning the prior is more meaningful if we can do it on a separate, but related dataset. In
30 [SZ+22] they propose to train a model on an initial, large dataset \mathcal{D}_1 (possibly unsupervised) to
31 get a point estimate, $\hat{\boldsymbol{\theta}}_1$, from which they can derive an approximate low-rank Gaussian posterior,
32 using the SWAG method (Section 17.3.8). They then use this informative prior when fine-tuning
33 the model on a downstream dataset \mathcal{D}_2 . The fine-tuning can either be a MAP estimate $\hat{\boldsymbol{\theta}}_2$ or some
34 approximate posterior, $p(\boldsymbol{\theta}_2|\mathcal{D}_2, \mathcal{D}_1)$, e.g., computed using MCMC (Section 17.3.7). They call this
35 technique “**Bayesian transfer learning**”. (See Section 19.5.1 for more details on transfer learning.)
36

37

38 17.2.4 Priors in function space

39

40 Typically, the relationship between the prior distribution over parameters and the functions preferred
41 by the prior is not transparent. In some cases, it can be possible to pick more informative priors
42 based on principles such as desired invariances that we want the function to satisfy (see e.g., [Nal18]).
43 [FBW21] introduces *residual pathway priors*, providing a mechanism for encoding high level concepts
44 into prior distributions, such as locality, independencies, and symmetries, without constraining model
45 flexibility. A different approach to encoding interpretable priors over functions leverages kernel
46 methods such as Gaussian processes (e.g., [Sun+19a]), as we discuss in Section 18.1.
47

17.2.5 Architectural priors

Beyond specifying the parametric prior, it is important to note that the architecture of the model can have an even larger effect on the induced distribution over functions, as argued in Wilson and Izmailov [WI20] and Izmailov et al. [Izm+21b]. For example, a CNN architecture encode prior knowledge about translation equivariance, due to its use of convolution, and hierarchical structure, due to its use of multiple layers. Other forms of inductive bias are induced by different architectures, such as RNNs. (Models such as transformers have weaker inductive bias, but consequently often need more data to perform well.) Thus we can think of the field of **neural architecture search** (reviewed in [EMH19]) as a form of structural prior learning.

In fact, with a suitable architecture, we can often get good results using random (untrained) models. For example, Ulyanov, Vedaldi, and Lempitsky [UVL18] showed that an untrained CNN with random parameters (sampled from a Gaussian) often works very well for low-level image processing tasks, such as image denoising, super-resolution and image inpainting. The resulting prior over functions has been called the **deep image prior**. Similarly, Pinto and Cox [PC12] showed that untrained CNNs with the right structure can do well at face recognition. Moreover, Zhang et al. [Zha+17] show that randomly initialized CNNs can process data to provide features that greatly improve the performance of other models, such as kernel methods.

17.3 Posteriors for BNNs

There are a large number of different approximate inference schemes that have been applied to Bayesian neural networks, with different strengths and limitations. In the sections below, we briefly describe some of these.

17.3.1 Monte Carlo dropout

Monte Carlo dropout (MCD) [GG16; KG17] is a very simple and widely used method for approximating the Bayesian predictive distribution. Usually stochastic dropout layers are added as a form of regularization, and are “turned off” at test time, as described in Section 16.2.6. However, the idea in MCD is to also perform random sampling at test time. More precisely, we drop out each hidden unit by sampling from a Bernoulli(p) distribution; we repeat this procedure S times, to create S distinct models. We then create an equally weighted average of the predictive distributions for each of these models:

$$p(y|\mathbf{x}, \mathcal{D}) \approx \frac{1}{S} \sum_{s=1}^S p(y|\mathbf{x}, \boldsymbol{\theta}^s) \quad (17.13)$$

where $\boldsymbol{\theta}^s$ is a version of the MAP parameter estimate where we randomly drop out some connections.

We give an example of this process in action in Figure 17.2. We see that it successfully captures uncertainty due to “out of distribution” inputs. (See Section 19.3.2 for more discussion of OOD detection.)

One drawback of MCD is that it is slow at test time. However this can be overcome by “distilling” the model’s predictions into a deterministic “student” network, as we discuss in sec:distillation.

A more fundamental problem is that MCD does not give proper uncertainty estimates, as argued in [Osb16; LF+21]. The problem is the following. Although MCD can be viewed as a form of variational

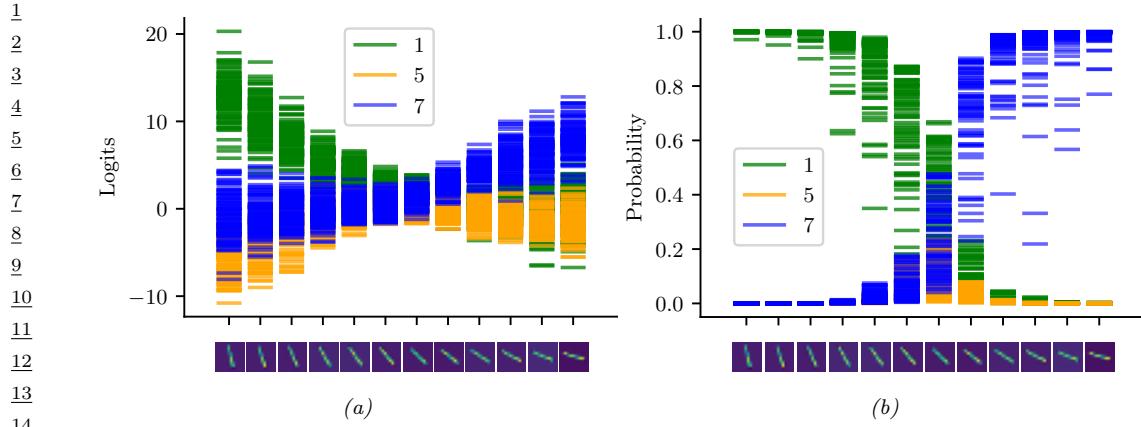


Figure 17.2: Illustration of MC dropout applied to the LeNet architecture. The inputs are some rotated images of the digit 1 from the MNIST dataset. (a) shows softmax inputs (logits) and the (b) shows the softmax outputs (probabilities). We see that the inputs are classified as digit 7 for the last three images (as shown by the probabilities), even though the model has high uncertainty (as shown by the logits). Adapted from Figure 4 of [GG16]. Generated by `mnist_classification_mc_dropout.ipynb`

inference [GG16], this is only true under a degenerate posterior approximation, corresponding to a mixture of two delta functions, one at 0 (for dropped out nodes) and one at the MLE. This posterior will not converge to the true posterior (which is a delta function at the MLE) even as the training set size goes to infinity, since we are always dropping out hidden nodes with a constant probability p [Osb16]. Fortunately this pathology can be fixed if the noise rate is optimized [GHK17]. For more details, see e.g., [HMG18; NHL19; LF+21].

17.3.2 Laplace approximation

In Section 7.4.3, we introduced the Laplace approximation, which computes a Gaussian approximation to the posterior, $p(\boldsymbol{\theta}|\mathcal{D})$, centered at the MAP estimate, $\boldsymbol{\theta}^*$. The posterior prediction matrix is equal to the Hessian of the negative log joint computed at the mode. The benefits of this approach are that it is simple, and it can be used to derive a Bayesian estimate from a pretrained model. A detailed explanation of the method in the context of DNNs can be found in [Dax+21]; here we just give a summary.

Let $\mathbf{f}(\mathbf{x}_n, \boldsymbol{\theta}) \in \mathbb{R}^C$ be the prediction function with C outputs, and $\boldsymbol{\theta} \in \mathbb{R}^P$ be the parameter vector. Let $\mathbf{r}(\mathbf{y}; \mathbf{f}) = \nabla_{\mathbf{f}} \log p(\mathbf{y}|\mathbf{f})$ be the residual¹, and $\Lambda(\mathbf{y}; \mathbf{f}) = -\nabla_{\mathbf{f}}^2 \log p(\mathbf{y}|\mathbf{f})$ be the per-input noise term. In addition, let $\mathbf{J} \in \mathbb{R}^{C \times P}$ be the Jacobian, $[\mathbf{J}_{\boldsymbol{\theta}}(\mathbf{x})]_{ci} = \frac{\partial f_c(\mathbf{x}, \boldsymbol{\theta})}{\partial \theta_i}$, and $\mathbf{H} \in \mathbb{R}^{C \times P \times P}$ be the Hessian, $[\mathbf{H}_{\boldsymbol{\theta}}(\mathbf{x})]_{cij} = \frac{\partial^2 f_c(\mathbf{x}, \boldsymbol{\theta})}{\partial \theta_i \partial \theta_j}$. Then the gradient and Hessian of the log likelihood are given by

¹ In the Gaussian case, this term becomes $\nabla_{\mathbf{f}} \|\mathbf{y} - \mathbf{f}\|^2 = 2\|\mathbf{y} - \mathbf{f}\|$, so can be interpreted as a residual error.

1 the following [IKB21]:
2

$$\nabla_{\boldsymbol{\theta}} \log p(\mathbf{y}|\mathbf{f}(\mathbf{x}, \boldsymbol{\theta})) = \mathbf{J}_{\boldsymbol{\theta}}(\mathbf{x})^T \mathbf{r}(\mathbf{y}; \mathbf{f}) \quad (17.14)$$

$$\nabla_{\boldsymbol{\theta}}^2 \log p(\mathbf{y}|\mathbf{f}(\mathbf{x}, \boldsymbol{\theta})) = \mathbf{H}_{\boldsymbol{\theta}}(\mathbf{x})^T \mathbf{r}(\mathbf{y}; \mathbf{f}) - \mathbf{J}_{\boldsymbol{\theta}}(\mathbf{x})^T \boldsymbol{\Lambda}(\mathbf{y}; \mathbf{f}) \mathbf{J}_{\boldsymbol{\theta}}(\boldsymbol{\theta}) \quad (17.15)$$

3
4 Since the network Hessian \mathbf{H} is usually intractable to compute, it is usually dropped, leaving only the
5 Jacobian term. This is called the **generalized Gauss-Newton** or **GGN** approximation [Sch02;
6 Mar20]. The GGN approximation is guaranteed to be positive definite. By contrast, this is not
7 true for the original Hessian in Equation (17.15), since the objective is not convex. Furthermore,
8 computing the Jacobian term is cheaper to compute than the Hessian.
9

10 Putting it all together, for a Gaussian prior, $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|\mathbf{m}_0, \mathbf{S}_0)$, the Laplace approximation
11 becomes $p(\boldsymbol{\theta}|\mathcal{D}) \approx (\mathcal{N}|\boldsymbol{\theta}^*, \boldsymbol{\Sigma}_{\text{GGN}})$, where
12

$$\boldsymbol{\Sigma}_{\text{GGN}}^{-1} = \sum_{n=1}^N \mathbf{J}_{\boldsymbol{\theta}^*}(\mathbf{x}_n)^T \boldsymbol{\Lambda}(\mathbf{y}_n; \mathbf{f}_n) \mathbf{J}_{\boldsymbol{\theta}^*}(\mathbf{x}_n) + \mathbf{S}_0^{-1} \quad (17.16)$$

13
14 Unfortunately inverting this matrix takes $O(P^3)$ time, so for models with many parameters, further
15 approximations are usually used. The simplest is to use a diagonal approximation, which takes $O(P)$
16 time and space. A more sophisticated approach is presented in [RBB18a], which leverages the **KFAC**
17 (Kronecker FActored Curvature) approximation of [MG15]. This approximates the covariance of each
18 layer using a kronecker product.
19

20 A limitation of the Laplace approximation is that the posterior covariance is derived from the
21 Hessian evaluated at the MAP parameters. This means Laplace forms a highly *local* approximation:
22 even if the non-Gaussian posterior could be well-described by a Gaussian distribution, the Gaussian
23 distribution *formed using Laplace* only captures the local characteristics of the posterior at the
24 MAP parameters — and may therefore suffer badly from local optima, providing overly compact
25 or diffuse representations. In addition, the curvature information is only used after the model has
26 been estimated, and not during the model optimization process. By contrast, variational inference
27 (Section 17.3.3) can provide more accurate approximations for comparable cost.
28

29 17.3.3 Variational inference 30

31
32 In fixed-form variational inference (Section 10.3), we choose a distribution for the posterior approxi-
33 mation $q_{\psi}(\boldsymbol{\theta})$ and minimize $D_{\text{KL}}(q \parallel p)$, with respect to ψ . We often choose a Gaussian approximate
34 posterior, $q_{\psi}(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$, which lets us use the reparameterization trick to create a low variance
35 estimator of the gradient of the ELBO (see Section 10.3.3). Despite the use of a Gaussian, the
36 parameters that minimize the KL objective are often different what we would find with the Laplace
37 approximation (Section 17.3.2).
38

39 Variational methods for neural networks date back to at least Hinton and Camp [HC93]. In deep
40 learning, [Gra11] revisited variational methods, using a Gaussian approximation with a diagonal
41 covariance matrix. This approximates the distribution of every parameter in the model by a univariate
42 Gaussian, where the mean is the point estimate, and the variance captures the uncertainty, as shown in
43 Figure 17.3. This approach was improved further in [Blu+15], who used the reparameterization trick
44 to compute lower variance estimates of the ELBO; they called their method **Bayes by backprop**
45 (**BBB**).
46

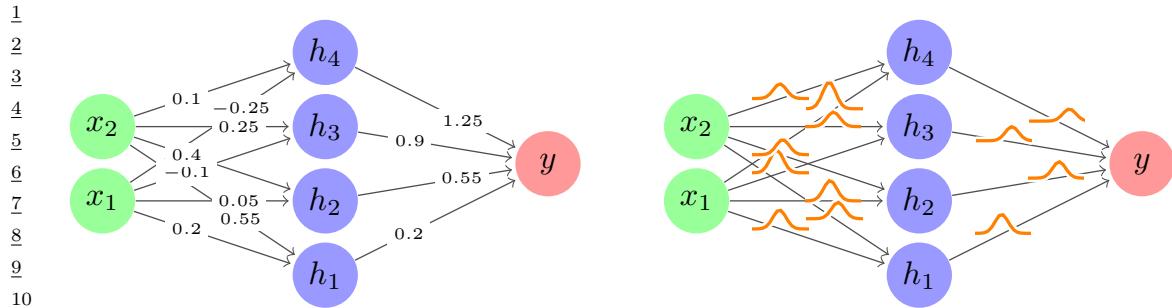


Figure 17.3: Illustration of an MLP with (Left) point estimate for each weight, (Right) a marginal distribution for each weight, corresponding to a fully factored posterior approximation.

In [Blu+15], they used a diagonal Gaussian posterior and vanilla SGD. In [Osa+19a], they used the **variational online Gauss-Newton (VOGN)** method of [Kha+18], for improved scalability. VOGN is a noisy version of natural gradient descent, where the extra noise emulates the effect of variational inference. In [Mis+18], they replaced the diagonal approximation with a low-rank plus diagonal approximation, and used VOGN for fitting. In [Tra+20b], they use a rank-one plus diagonal approximation known as **NAGVAC** (see Section 10.3.4.2). In this case, there are only 3 times as many parameters as when computing a point estimate (for the variational mean, variance, and rank-one vector), making the approach very scalable. In addition, in this case it is possible to analytically compute the natural gradient, which speeds up model fitting (see Section 6.4). Many other variational methods have also been proposed (see e.g., [LW16; Zha+18; Wu+19a; HHK19]). See also Section 17.5.4 for a discussion of online VI for DNNs.

17.3.4 Expectation propagation

Expectation propagation (EP) is similar to variational inference, except it locally optimizes $D_{\text{KL}}(p \parallel q)$ instead of $D_{\text{KL}}(q \parallel p)$, where p is the exact posterior and q is the approximate posterior. For details, see Section 10.7.

A special case of EP is the assumed density filtering (ADF) algorithm of Section 8.10, which is equivalent to the first pass of ADF. In Section 8.10.3 we show how to apply ADF to online logistic regression. In [HLA15a], they extend ADF to the case of BNNs; they called their method probabilistic backpropagation or **PBP**. They approximate every parameter in the model by a Gaussian factor, as in Figure 17.3. See Section 17.5.3 for the details.

17.3.5 Last layer methods and SNGP

A very simple approximation is to only “be Bayesian” about the weights in the final layer, and to use MAP estimates for all the other parameters. This is called the **neural-linear** approximation [RTS18]. In [KHH20] they show this can reduce overconfidence in predictions for inputs that are far from the training data. However, this approach ignores uncertainty introduced by the earlier feature extraction layers, where most of the parameters reside. We discuss a solution to this in Section 17.3.6.

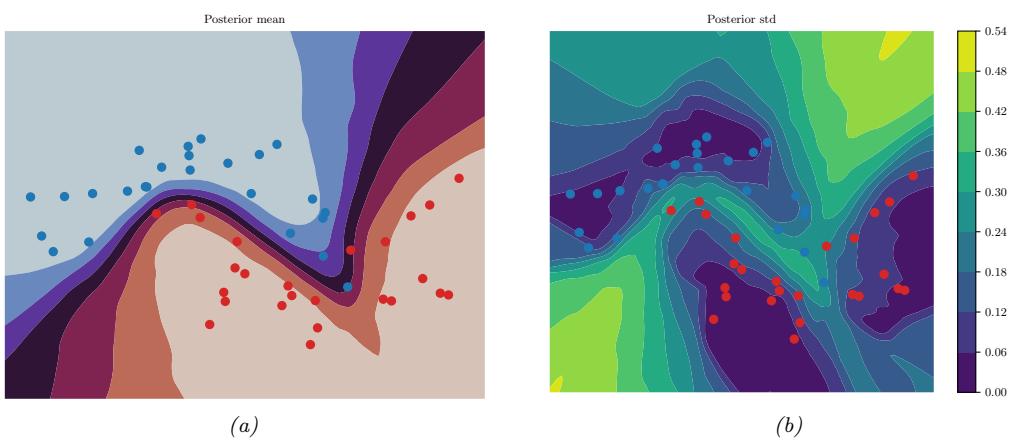


Figure 17.4: Illustration of an MLP fit to the two-moons dataset using HMC. (a) Posterior mean. (b) Posterior standard derivation. The uncertainty increases as we move away from the training data. Generated by [bnn_mlp_2d_hmc.ipynb](#).

17.3.6 SNGP

It is possible to combine DNNs with Gaussian Process (GP) models (Chapter 18), by using the DNN to act as a feature extractor, which is then fed into the kernel in the final layer. This is called “deep kernel learning” (see Section 18.6.6).

One problem with this is that the feature extractor may lose information which is not needed for classification accuracy, but which is needed for robust performance on out-of-distribution inputs (see Section 17.4.6.2). The basic problem is that, in a classification problem, there is no reduction in training accuracy (log likelihood) if points which are far away are projected close together, as long as they are on the correct side of the decision boundary. Thus the distances between two inputs can be erased by the feature extraction layers, so that OOD inputs appear to the final layer to be close to the training set.

One solution to this is to use the **SNGP** (spectrally normalized Gaussian process) method of [Liu+20d]. This constrains the feature extraction layers to be “distance preserving”, so that two inputs that are far apart in input space remain far apart after many layers of feature extraction, but using spectral normalization of the weights to bound the Lipschitz constant of the feature extractor. The overall approach ensures that information that is relevant for computing the confidence of a prediction, but which might be irrelevant to computing the label of a prediction, is not lost. This can help performance in tasks such as out-of-distribution detection (Section 17.4.6.2).

17.3.7 MCMC methods

Some of the earliest work on inference for BNNs was done by Radford Neal, who proposed to use Hamiltonian Monte Carlo (Section 12.5) to approximate the posterior [Nea96]. This is generally considered the gold standard method, since it does not make strong assumptions about the form of the posterior. For more recent work on scaling up HMC for BNNs, see e.g., [Izm+21b; CJ21].

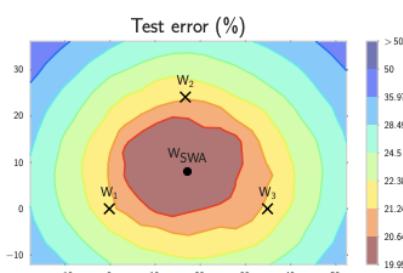


Figure 17.5: Illustration of stochastic weight averaging (SWA). The three crosses represent different SGD solutions. The star in the middle is the average of these parameter values. From Figure 1 of [Izm+18]. Used with kind permission of Andrew Wilson.

We give a simple example of vanilla HMC in Figure 17.4, where we fit a shallow MLP to a small 2d binary dataset. we plot the mean and standard deviation of the posterior predictive distribution, $p(y = 1|x; \mathcal{D})$. We see that the uncertainty is higher as we move away from the training data. (Compare to Bayesian logistic regression in 1d in Figure 15.8a.)

However, a significant limitation of standard MCMC procedures, including HMC, is that they require access to the full training set at each step. Stochastic gradient MCMC methods, such as SGLD, operate instead using mini-batches of data, offering a scalable alternative, as we discuss in Section 12.7.1. For an example of SGLD applied to an MLP, see Section 19.3.3.1.

17.3.8 Methods based on the SGD trajectory

In [MHB17; SL18; CS18], it was shown that, under some assumptions, the iterates produced by stochastic gradient descent (SGD), when run at a fixed learning rate, correspond to samples from a Gaussian approximation to the posterior centered at a local mode, $p(\theta|\mathcal{D}) \approx \mathcal{N}(\theta|\hat{\theta}, \Sigma)$. We can therefore use SGD to generate approximate posterior samples. This is similar to SG-MCMC methods, except we do not add explicit gradient noise, and the learning rate is held constant.

In [Izm+18], they noted that these SGD solutions (with fixed learning rate) surround the periphery of points of good generalization, as shown in Figure 17.5. This is in part because SGD does not converge to a local optimum unless the learning rate is annealed to 0. They therefore proposed to compute the average of several SGD samples, each one collected after a certain interval (e.g., one epoch of training), to get $\bar{\theta} = \frac{1}{S} \sum_{s=1}^S \theta_s$. They call this **stochastic weight averaging (SWA)**. They showed that the resulting point tends to correspond to a broader local minimum than the SGD solutions (c.f., Figure 17.9), resulting in better generalization performance.

The SWA approach is related to Polyak-Ruppert averaging, which is often used in convex optimization. The difference is that Polyak-Ruppert typically assumes the learning rate decays to zero, and uses an exponential moving average (EMA) of iterates, rather than an equal average; Polyak-Ruppert averaging is mainly used to reduce variance in the SGD estimate, rather than as a method to find points of better generalization.

The SWA approach is also related to **snapshot ensembles** [Hua+17a], and **fast geometric ensembles** [Gar+18c]; these methods save the parameters θ_s after increasing and decreasing the

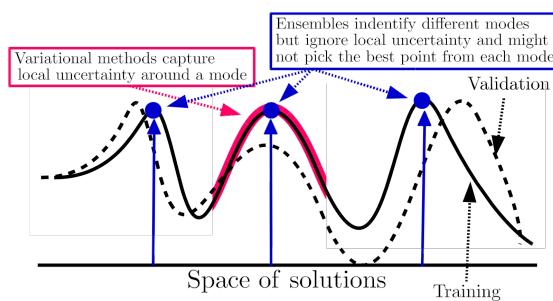


Figure 17.6: Cartoon illustration of the NLL as it varies across the parameter space. Subspace methods (red) model the local neighborhood around a local mode, whereas ensemble methods (blue) approximate the posterior using a set of distinct modes. From Figure 1 of [FHL19]. Used with kind permission of Balaji Lakshminarayanan.

learning rate multiple times in a cyclical fashion, and then computing the *average of the predictions* using $p(\mathbf{y}|\mathbf{x}, \mathcal{D}) \approx \frac{1}{S} \sum_{s=1}^S p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}_s)$, rather than computing the *average of the parameters* and predicting with a single model (which is faster). Moreover, by finding a flat region, representing a “center or mass” in the posterior, SWA can be seen as approximating the Bayesian model average in Equation 17.1 with a single model.

In [Mad+19], they proposed to fit a Gaussian distribution to the set of samples produced by SGD near a local mode. They use the SWA solution as the mean of the Gaussian. For the covariance matrix, they use a low-rank plus diagonal approximation of the form $p(\boldsymbol{\theta}|\mathcal{D}) = \mathcal{N}(\boldsymbol{\theta}|\bar{\boldsymbol{\theta}}, \boldsymbol{\Sigma})$, where $\boldsymbol{\Sigma} = (\boldsymbol{\Sigma}_{\text{diag}} + \boldsymbol{\Sigma}_{\text{lr}})/2$, $\boldsymbol{\Sigma}_{\text{diag}} = \text{diag}(\bar{\boldsymbol{\theta}}^2 - (\bar{\boldsymbol{\theta}})^2)$, $\bar{\boldsymbol{\theta}} = \frac{1}{S} \sum_{s=1}^S \boldsymbol{\theta}_s$, $\bar{\boldsymbol{\theta}}^2 = \frac{1}{S} \sum_{s=1}^S \boldsymbol{\theta}_s^2$, and $\boldsymbol{\Sigma}_{\text{lr}} = \frac{1}{S} \Delta \Delta^\top$ is the sample covariance matrix of the last K samples of $\Delta_i = (\boldsymbol{\theta}_i - \bar{\boldsymbol{\theta}})$, where $\bar{\boldsymbol{\theta}}_i$ is the running average of the parameters from the first i samples. They call this method **SWAG**, which stands for “stochastic weight averaging with Gaussian posterior”. This can be used to generate an arbitrary number of posterior samples at prediction time. They show that SWAG scales to large residual networks with millions of parameters, and large datasets such as ImageNet, with improved accuracy and calibration over conventional SGD training, and no additional training overhead.

17.3.9 Deep ensembles

Many conventional approximate inference methods focus on approximating the posterior $p(\boldsymbol{\theta}|\mathcal{D})$ in a local neighborhood around one of the posterior modes. While this is often not a major limitation in classical machine learning, modern deep neural networks have highly multi-modal posteriors, with parameters in different modes giving rise to very different functions. On the other hand, the functions in a neighborhood of a single mode may make fairly similar predictions. So using such a local approximation to compute the posterior predictive will underestimate uncertainty and generalize more poorly.

A simple alternative method is to train multiple models, and then to approximate the posterior

1 using an equally weighted mixture of delta functions,

2

$$\frac{4}{5} p(\boldsymbol{\theta}|\mathcal{D}) \approx \frac{1}{M} \sum_{m=1}^M \delta(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_m) \quad (17.17)$$

3

4 where M is the number of models, and $\hat{\boldsymbol{\theta}}_m$ is the MAP estimate for model m . See Figure 17.6 for a
5 sketch. This approach is called **deep ensembles** [LPB17; FHL19].

6 The models can differ in terms of their random seed used for initialization [LPB17], or hyper-
7 parameters [Wen+20c], or architecture [Zai+20], or all of the above. In addition, [DF21; TB22]
8 discusses how to add an explicit repulsive term to ensure functional diversity between the ensemble
9 members. This way, each member corresponds to a distinct prediction function. Combining these is
10 more effective than combining multiple samples from the same basin of attraction, especially in the
11 presence of dataset shift [Ova+19].

12

13 17.3.9.1 Multi-SWAG

14

15 We can further improve on this approach by fitting a Gaussian to each local mode using the SWAG
16 method from Section 17.3.8 to get a mixture of Gaussians approximation:

17

$$\frac{21}{22} p(\boldsymbol{\theta}|\mathcal{D}) \approx \frac{1}{M} \sum_{m=1}^M \mathcal{N}(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}}_m, \boldsymbol{\Sigma}_m) \quad (17.18)$$

23

24 This approach is known as **MultiSWAG** [WI20]. MultiSWAG performs a Bayesian model average
25 both across multiple basins of attraction, like deep ensembles, but also within each basin, and provides
26 an easy way to generate an arbitrary number of posterior samples, $S > M$, in an any-time fashion.

27

28 17.3.9.2 Deep ensembles with random priors

29

30 The standard way to fit each member of a deep ensemble is to initialize them each with a different
31 random set of parameters, but them to train them all on the same data. Unfortunately this can
32 result in the predictions from each ensemble member being rather similar, which reduces the benefit
33 of the approach. One way to increase diversity is to train each member on a different subset of the
34 data; this is called **bootstrap sampling**. Another approach is to define the i 'th ensemble member
35 $g_i(\mathbf{x})$ to be the addition of a trainable model $t_i(\mathbf{x})$ and a fixed, but random, **prior network**, $p_i(\mathbf{x})$,
36 to get

37

$$g_i(\mathbf{x}; \boldsymbol{\theta}_i) = t_i(\mathbf{x}; \boldsymbol{\theta}_i) + \beta p_i(\mathbf{x}) \quad (17.19)$$

38

39 where $\beta \geq 0$ controls the amount of data-independent variation between the members. The trainable
40 network learns to model the residual error between the true output and the value predicted by the
41 prior. This is called a **random prior deep ensemble** [OAC18]. See Figure 17.7 for an illustration.
42

43 17.3.9.3 Deep ensembles as approximate Bayesian inference

44

45 The posterior predictive distribution for a Bayesian neural network cannot be expressed in closed
46 form. Therefore all Bayesian inference approaches in deep learning are approximate. In this context,
47

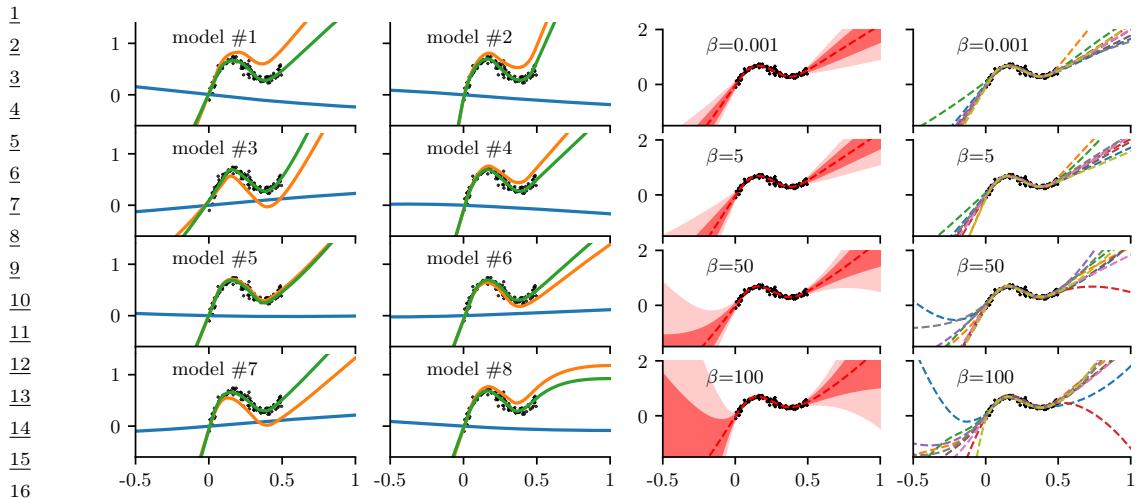


Figure 17.7: Deep ensemble with random priors. (a) Individual predictions from each member. Blue is the fixed random prior function, orange is the trainable function, green is the combination of the two. (b) Overall prediction from the ensemble, for increasingly large values of β . On the left we show (in red) the posterior mean and pointwise standard deviation, and on the right we show samples from the posterior. As β increases, we trust the random priors more, and pay less attention to the data, thus getting a more diffuse posterior.

Generated by [randomized_priors.ipynb](#).

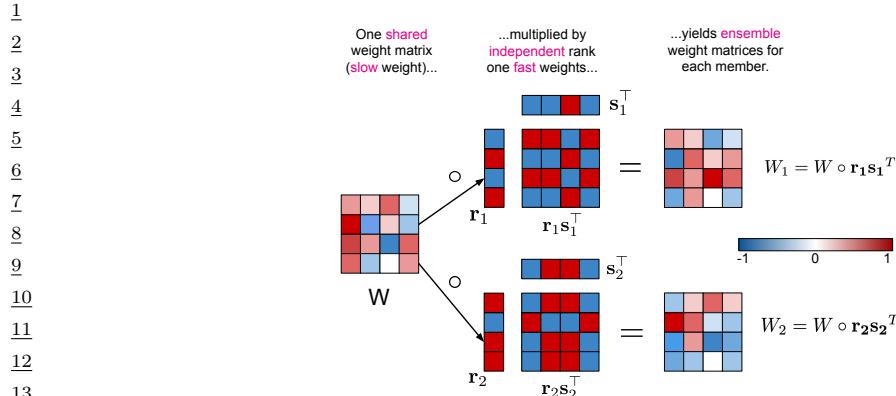
all approximate inference procedures fall onto a spectrum, representing how closely they approximate the true posterior predictive distribution. Deep ensembles can provide better approximations to a Bayesian model average than a single basin marginalization approach, because point masses from different basins of attraction represent greater functional diversity than standard Bayesian approaches which sample within a single basin.

17.3.9.4 Deep ensembles vs classical ensembles

Note that deep ensembles is slightly different to classical ensemble methods (see e.g., [Die00]), such as bagging and random forests, which obtain diversity of its predictors by training them on different subsets of the data (created using bootstrap resampling), or on different features. This data perturbation is necessary to get diversity when the base learner is a convex problem (such as a linear model, or shallow decision tree). In the deep ensemble approach, every model is trained on the same data, and the same input features. The diversity arises due to different starting parameters, different random seeds, and SGD noise, which induces different solutions due to the nonconvex loss. It is also possible to explicitly enforce diversity of the ensemble members, which can provably improve performance [TB22].

17.3.9.5 Deep ensembles vs mixtures of experts and stacking

If we use weighted combinations of the models, $p(\boldsymbol{\theta}|\mathcal{D}) = \sum_{m=1}^M p(m|\mathcal{D})p(\boldsymbol{\theta}|m, \mathcal{D})$, where $p(m|\mathcal{D})$ is the marginal likelihood of model m , then, in the large sample limit, this mixture will concentrate on



14 *Figure 17.8: Illustration of batch ensemble with 2 ensemble members. From Figure 2 of [WTB20]. Used with*
 15 *kind permission of Paul Vicol.*

16
17
18
19 the MAP model, so only one component will be selected. By contrast, in deep ensembles, we always
 20 use M equally weighted models. Thus we see that Bayes model averaging is not the same as model
 21 ensembling [Min00b]. Indeed, ensembling can enlarge the expressive power of the posterior predictive
 22 distribution compared to BMA [OCM21].

23 We can also make the mixing weights be conditional on the inputs:

24
25 $p(y|\mathbf{x}, \mathcal{D}) = \sum_m w_m(\mathbf{x}) p(y|\mathbf{x}, \theta_m)$ (17.20)
 26

27 If we constrain the weights to be non-zero and sum to one, this is called a **mixture of experts**.
 28 However, if we allow a general positive weighted combination, the approach is called **stacking** [Wol92;
 29 Bre96; Yao+18a; CAII20]. In stacking, the weights $w_m(\mathbf{x})$ are usually estimated on hold-out data,
 30 to make the method more robust to model misspecification.
 31

32 17.3.9.6 Batch ensemble 33

34 Deep ensembles require M times more memory and time than a single model. One way to reduce
 35 the memory cost is to share most of the parameters — which we call **slow weights**, \mathbf{W} — and then
 36 let each ensemble member m estimate its own local perturbation, which we will call **fast weights**,
 37 \mathbf{F}_m . We then define $\mathbf{W}_m = \mathbf{W} \odot \mathbf{F}_m$. For efficiency, we can define \mathbf{F}_m to be a rank-one matrix,
 38 $\mathbf{F}_m = \mathbf{s}_m \mathbf{r}_m^\top$, as illustrated in Figure 17.8. This is called **batch ensemble** [WTB20].

39 It is clear that the memory overhead is very small compared to naive ensembles, since we just need
 40 to store $2M$ vectors (\mathbf{s}_m^l and \mathbf{r}_m^l) for every layer, which is negligible compared to the quadratic cost
 41 of storing the shared weight matrix \mathbf{W}^l .

42 In addition to memory savings, batch ensemble can reduce the inference time by a constant factor
 43 by leveraging within-device parallelism. To see this, consider the output of one layer using ensemble
 44 m on example n :

45
46 $y_n^m = \varphi(\mathbf{W}_m^\top \mathbf{x}_n) = \varphi((\mathbf{W} \odot \mathbf{s}_m \mathbf{r}_m^\top)^\top \mathbf{x}_n) = \varphi((\mathbf{W}^\top (\mathbf{x}_n \odot \mathbf{s}_m) \odot \mathbf{r}_m))$ (17.21)
 47

We can vectorize this for a minibatch of inputs \mathbf{X} by replicating \mathbf{r}_m and \mathbf{s}_m along the B rows in the batch to form matrices, giving

$$\mathbf{Y}_m = \varphi(((\mathbf{X} \odot \mathbf{S}_m) \mathbf{W}) \odot \mathbf{R}_m) \quad (17.22)$$

This applies the same ensemble parameters m to every example in the minibatch of size B . To achieve diversity during training, we can divide the minibatch into M sub-batches, and use sub-batch m to train \mathbf{W}_m . (Note that this reduces the batch size for training each ensemble to B/M .) At test time, when we want to average over M models, we can replicate each input M times, leading to a batch size of BM .

In [WTB20], they show that this method outperforms MC dropout at negligible extra memory cost. However, the best combination was to combine batch ensemble with MC dropout; in some cases, this approached the performance of naive ensembles.

17.3.10 Approximating the posterior predictive distribution

Once we have approximated the parameter posterior, $q(\boldsymbol{\theta}) \approx p(\boldsymbol{\theta}|\mathcal{D})$, we can use it to approximate the posterior predictive distribution:

$$p(\mathbf{y}|\mathbf{x}, \mathcal{D}) = \int q(\boldsymbol{\theta}) p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) d\boldsymbol{\theta} \quad (17.23)$$

We usually approximate this integral using Monte Carlo:

$$p(\mathbf{y}|\mathbf{x}, \mathcal{D}) \approx \frac{1}{S} \sum_{s=1}^S p(\mathbf{y}|\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}^s)) \quad (17.24)$$

where $\boldsymbol{\theta}^s \sim q(\boldsymbol{\theta})$. We discuss some extensions of this approach below.

17.3.10.1 A linearized approximation

In [IKB21] they point out that samples from an approximate posterior, $q(\boldsymbol{\theta})$, can result in bad predictions when plugged into the model if the posterior puts probability density “in the wrong places”. This is because $\mathbf{f}(\mathbf{x}; \boldsymbol{\theta})$ is a highly nonlinear function of $\boldsymbol{\theta}$ that might behave quite differently when $\boldsymbol{\theta}$ is far from the MAP estimate on which $q(\boldsymbol{\theta})$ is centered. To avoid this problem, they propose to replace $\mathbf{f}(\mathbf{x}; \boldsymbol{\theta})$ with a linear approximation centered at the MAP estimate $\boldsymbol{\theta}^*$:

$$\mathbf{f}_{\text{lin}}^{\boldsymbol{\theta}^*}(\mathbf{x}, \boldsymbol{\theta}) = \mathbf{f}(\mathbf{x}, \boldsymbol{\theta}^*) + \mathbf{J}_{\boldsymbol{\theta}^*}(\boldsymbol{\theta} - \boldsymbol{\theta}^*) \quad (17.25)$$

Such a model is well behaved around $\boldsymbol{\theta}^*$, and so the approximation

$$p(\mathbf{y}|\mathbf{x}, \mathcal{D}) \approx \frac{1}{S} \sum_{s=1}^S p(\mathbf{y}|\mathbf{f}_{\text{lin}}^{\boldsymbol{\theta}^*}(\mathbf{x}, \boldsymbol{\theta}^s)) \quad (17.26)$$

often works better than Equation (17.24).

Note that $\mathbf{f}_{\text{lin}}^{\boldsymbol{\theta}^*}(\mathbf{x}, \boldsymbol{\theta})$ is a linear function of the parameters $\boldsymbol{\theta}$, but a nonlinear function of the inputs \mathbf{x} . Thus $p(\mathbf{y}|\mathbf{f}_{\text{lin}}^{\boldsymbol{\theta}^*}(\mathbf{x}, \boldsymbol{\theta}))$ is a generalized linear model (Section 15.1), so [IKB21] call this approximation the **GLM predictive distribution**.

1 **17.3.10.2 Distillation**

3 The MC approximation to the posterior predictive is S times slower than a standard, deterministic
4 plug-in approximation. One way to speed this up is to use **distillation** to approximate the
5 semi-parametric “teacher” model p_t from Equation (17.24) by a parametric “student” model p_s
6 by minimizing $\mathbb{E}[D_{\text{KL}}(p_t(\mathbf{y}|\mathbf{x}) \parallel p_s(\mathbf{y}|\mathbf{x}))]$ wrt p_s . This approach was first proposed in [HVD14],
7 who called the technique “**dark knowledge**”, because the teacher has “hidden” information in its
8 predictive probabilities (logits) than is not apparent in the raw one-hot labels.

9 In [Kor+15], this idea was used to distill the predictions from a teacher whose parameter posterior
10 was computed using HMC; this is called “**Bayesian dark knowledge**”. A similar idea was used in
11 [BPK16; GBP18], who distilled the predictive distribution derived from MC dropout (Section 17.3.1).
12 Since the parametric student is typically less flexible than the semi-parametric teacher, it may be
13 overconfident, and lack diversity in its predictions. To avoid this overconfidence, it is safer to make
14 the student be a mixture distribution c.f., [SG05]. See also [Tra+20a].

16 **17.3.11 Tempered and cold posteriors**

18 When working with BNNs for classification problems, the likelihood is usually taken to be

20 $p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) = \text{Cat}(y|\text{softmax}(f(\mathbf{x}; \boldsymbol{\theta})))$ (17.27)

21 where $f(\mathbf{x}; \boldsymbol{\theta}) \in \mathbb{R}^C$ returns the logits over the C class labels. This is the same as in multinomial
22 logistic regression (Section 15.3.2); the only difference is that f is a nonlinear function of $\boldsymbol{\theta}$.

24 However, in practice, it is often found (see e.g., [Zha+18; Wen+20b; LST21; Noc+21]) that BNNs
25 give better predictive accuracy if the likelihood function is scaled by some power α . That is, instead
26 of targeting the posterior $p(\boldsymbol{\theta}|\mathcal{D}) \propto p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})p(\boldsymbol{\theta})$, these methods target the **tempered posterior**,
27 $p_{\text{tempered}}(\boldsymbol{\theta}|\mathcal{D}) \propto p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})^\alpha p(\boldsymbol{\theta})$. In log space, we have

28 $\log p_{\text{tempered}}(\boldsymbol{\theta}|\mathcal{D}) = \alpha \log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) + \log p(\boldsymbol{\theta}) + \text{const}$ (17.28)

30 This is also called an **α -posterior** or **power posterior** [Med+21].

31 Another common method is to target the **cold posterior**, $p_{\text{cold}}(\boldsymbol{\theta}|\mathcal{D}) \propto p(\boldsymbol{\theta}|\mathbf{X}, \mathbf{y})^{1/T}$, or, in log
32 space,

34 $\log p_{\text{cold}}(\boldsymbol{\theta}|\mathcal{D}) = \frac{1}{T} \log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) + \frac{1}{T} \log p(\boldsymbol{\theta}) + \text{const}$ (17.29)

36 If $T < 1$, we say that the posterior is “cold”. Note that, in the case of a Gaussian prior, using the
37 cold prior is the same as using the tempered prior with a different hyperparameter, since

39 $\frac{1}{T} \log \mathcal{N}(\boldsymbol{\theta}|0, \sigma_{\text{cold}}^2 \mathbf{I}) = -\frac{1}{2T\sigma_{\text{cold}}^2} \sum_i \theta_i^2 + \text{const} = \mathcal{N}(\boldsymbol{\theta}|0, \sigma_{\text{tempered}}^2 \mathbf{I}) + \text{const}$ (17.30)

42 where $\sigma_{\text{tempered}}^2 = T\sigma_{\text{cold}}^2$. Thus both methods are effectively the same, and just reweight the
43 likelihood.

44 Cold posteriors in Bayesian neural network classifiers are a consequence of underrepresenting
45 aleatoric (label) uncertainty, as shown by [Kap+22]. On benchmarks such as CIFAR-100, we should
46 have essentially no uncertainty about the labels of the training images, yet Bayesian classifiers with

47

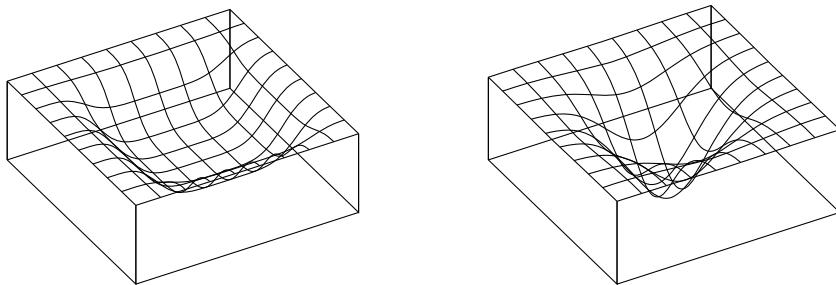


Figure 17.9: Flat vs sharp minima. From Figures 1 and 2 of [HS97]. Used with kind permission of Jürgen Schmidhuber.

softmax likelihoods have very high uncertainty for these points. Moreover, [Izm+21b] showed that the cold posterior effect in all the examples of [Wen+20b] when data augmentation is removed. [Kap+22] show that with the SGLD inference in [Wen+20b], data augmentation has the effect of raising the likelihood to a power $1/K$ for minibatches of size K . Cold posteriors exactly counteract this effect, more honestly representing our beliefs about aleatoric uncertainty, by sharpening the likelihood. However, tempering is not required, and [Kap+22] show that by using a Dirichlet observation model to explicitly represent (lack of) label noise, there is no cold posterior effect, even with data augmentation. The curation hypotheses of [Ait21] can be considered a special case of the above explanation, where curation has the effect of increasing our confidence about training labels.

In Section 14.1.3, we discuss generalized variational inference, which gives a general framework for understanding whether and how the likelihood or prior could benefit from tempering. Tempering is particularly useful if (as is usually the case) the model is misspecified [KJD21].

17.4 Generalization in Bayesian deep learning

In this section, we discuss why “being Bayesian” can improve predictive accuracy and generalization performance.

17.4.1 Sharp vs flat minima

Some optimization methods (in particular, second-order batch methods) are able to find “needles in haystacks”, corresponding to narrow but deep “holes” in the loss landscape, corresponding to parameter settings with very low loss. These are known as **sharp minima**, see Figure 17.9(right). From the point of view of minimizing the empirical loss, the optimizer has done a good job. However, such solutions generally correspond to a model that has overfit the data. It is better to find points that correspond to **flat minima**, as shown in Figure 17.9(left); such solutions are more robust and generalize better. To see why, note that flat minima correspond to regions in parameter space where there is a lot of posterior uncertainty, and hence samples from this region are less able to precisely memorize irrelevant details about the training set [AS17]. Put another way, the description length for sharp minima is large, meaning you need to use many bits of precision to specify the exact location in parameter space to avoid incurring large loss, whereas the description length for flat minima is

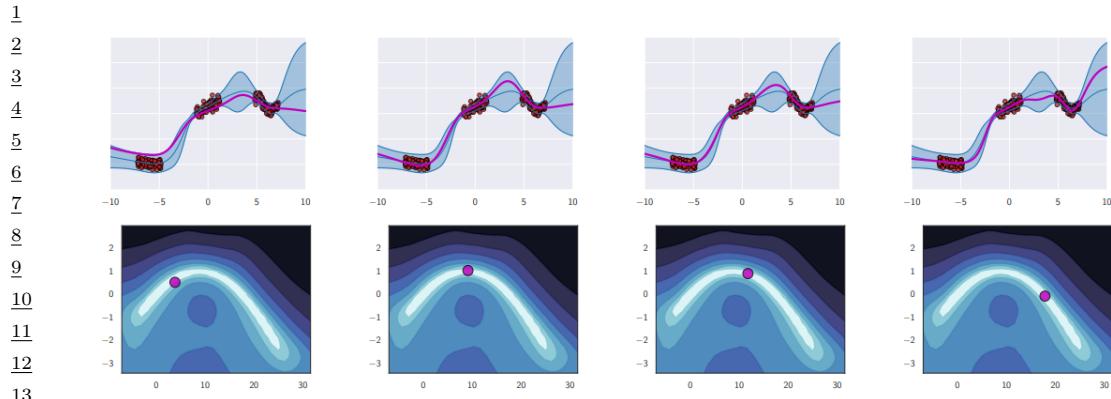


Figure 17.10: Diversity of high performing functions sampled from the posterior. Top row: we show predictions on the 1d input domain for 4 different functions. We see that they extrapolate in different ways outside of the support of the data. Bottom row: we show a 2d subspace spanning two distinct modes (MAP estimates), and connected by a low-loss curved path computed as in [Gar+18c]. From Figure 8 of [WI20]. Used with kind permission of Andrew Wilson.

19

20

21

22 less, resulting in better generalization [Mac03].

23 SGD often finds such flat minima by virtue of the addition of noise, which prevents it from
 24 “entering” narrow regions of the loss landscape (see Section 12.5.7). In addition, in higher dimensional
 25 spaces, flat regions occupy a much greater volume, and are thus much more easily discoverable by
 26 optimization procedures. More precisely, the analysis in [SL18] shows that the probability of entering
 27 any given basin of attraction \mathcal{A} around a minimum is given by $p_{SGD}(\boldsymbol{\theta} \in \mathcal{A}) \propto \int_{\mathcal{A}} e^{-\mathcal{L}(\boldsymbol{\theta})} d\boldsymbol{\theta}$. Note
 28 that this is integrating over the volume of space corresponding to \mathcal{A} , and hence is proportional to
 29 the model evidence (marginal likelihood) for that region, as explained in Section 3.9.1. Since the
 30 evidence is parameterization invariant (since we marginalize out the parameters), this means that
 31 SGD will avoid regions that have low evidence (corresponding to sharp minima) regardless of how we
 32 parameterize the model (contrary to the claims in [Din+17]).

33 In fact, several papers have shown that we can view SGD as approximately sampling from the
 34 Bayesian posterior (see Section 17.3.8). The SWA method (Section 17.3.8) can be seen as finding a
 35 center of mass in the posterior based on these SGD samples, finding solutions that generalize better
 36 than picking a single SGD point.

37 If we must use a single solution, a flat one will help us better approximate the Bayesian model
 38 average in the integral of Equation (17.1). However, by attempting to perform a more complete
 39 Bayesian model average, we will select for flatness without having to deal with the messiness of
 40 having to worry about flatness definitions, or the effects of reparametrization, or unknown implicit
 41 regularization, as the model average will automatically weight regions with the greatest volume.

42

43 17.4.2 Mode connectivity and the loss landscape

44

45 In DNNs there are often many low-loss solutions, which provide complementary explanations of
 46 the data. Moreover, in [Gar+18c] they showed that two independently trained SGD solutions can
 47

be connected by a curve in a subspace, along which the training loss remains near-zero, known as **mode connectivity**. Despite having the same training loss, these different parameter settings give rise to very different functions, as illustrated in Figure 17.10, where we show predictions on a 1d regression problem coming from different points in parameter space obtained by interpolating along a mode connecting curve between two distinct MAP estimates. Using a Bayesian model average, we can combine these functions together to provide much better performance over a single flat solution [Izm+19].

Recently, it has been discovered [Ben+21b] that there are in fact large multidimensional simplexes of low loss solutions, which can be combined together for significantly improved performance. These results further motivate the Bayesian approach (Equation (17.1)), where we perform a posterior weighted model average.

17.4.3 Effective dimensionality of a model

Modern DNNs have millions of parameters, but these parameters are often not well-determined by the data, i.e., there can be a lot of posterior uncertainty. By averaging over the posterior, we reduce the chance of overfitting, because we do not use “degrees of freedom” that are not needed or warranted.

To quantify the number of degrees of freedom, or **effective dimensionality** [Mac92b], we follow [MBW20] and define

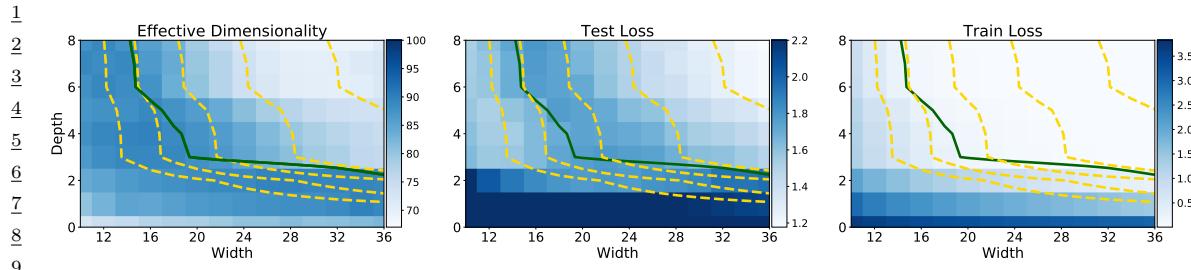
$$N_{\text{eff}}(\mathbf{H}, c) = \sum_{i=1}^k \frac{\lambda_i}{\lambda_i + c}, \quad (17.31)$$

where λ_i are the eigenvalues of the Hessian matrix \mathbf{H} computed at a local mode, and $c > 0$ is a regularization parameter. Intuitively, the effective dimension counts the number of well-determined parameters. A “flat minimum” will have many directions in parameter space that are not well-determined, and hence will have low effective dimensionality. This means that we can perform Bayesian inference in a low dimensional subspace [Izm+19]: Since there is functional homogeneity in all directions but those defining the effective dimension, neural networks can be significantly compressed.

This compression perspective can also be used to understand why the effective dimension can be a good proxy for generalization. If two models have similar training loss, but one has lower effective dimension, then it is providing a better compression for the data at the same fidelity. In Figure 17.11 we show that for CNNs with low training loss (above the green partition), the effective dimensionality closely tracks generalization performance. We also see that the number of parameters alone is not a strong determinant of generalization. Indeed, models with more parameters can have a lower number of effective parameters. We also see that wide but shallow models overfit, while depth helps provide lower effective dimensionality, leading to a better compression of the data. It is depth that makes modern neural networks distinctive, providing hierarchical inductive biases making it possible to discover more regularity in the data.

17.4.4 The hypothesis space of DNNs

Zhang et al. [Zha+17] showed that CNNs can fit CIFAR-10 images with random labels with zero training error, but can still generalize well on the noise-free test set. It has been claimed that this



¹⁰ Figure 17.11: Left: Effective dimensionality as a function of model width and depth for a CNN on CIFAR-100.
¹¹ Center: Test loss as a function of model width and depth. Right: Train loss as a function of model width and
¹² depth. Yellow level curves represent equal parameter counts ($1e5$, $2e5$, $4e5$, $1.6e6$). The green curve separates
¹³ models with near-zero training loss. Effective dimensionality serves as a good proxy for generalization for
¹⁴ models with low train loss. We see wide but shallow models overfit, providing low train loss, but high test
¹⁵ loss and high effective dimensionality. For models with the same train loss, lower effective dimensionality
¹⁶ can be viewed as a better compression of the data at the same fidelity. Thus depth provides a mechanism for
¹⁷ compression, which leads to better generalization. From Figure 2 of [MBW20]. Used with kind permission of
Andrew Wilson.

¹⁸

¹⁹

²⁰

²¹

result contradicts a classical understanding of generalization, because it shows that neural networks are capable of significantly overfitting the data, but can still generalize well on structured inputs.

We can resolve this paradox by taking a Bayesian perspective. In particular, we know that modern CNNs are very flexible, so they can fit almost pattern (since they are in fact universal approximators). However, their architecture encodes a prior over what kinds of patterns they expect to see in the data (see Section 17.2.5). Image datasets with random labels *can* be represented by this function class, but such solutions receive very low marginal likelihood, since they strongly violate the prior assumptions [WI20]. By contrast, image datasets where the output labels are consistent with patterns in the input get much higher marginal likelihood.

This phenomenon is not unique to DNNs. For example, it also occurs with Gaussian processes (Chapter 18). Such models are also universal approximators, but they allocate most of their probability mass to a small range of solutions (depending on the chosen kernel). They can also fit image datasets with random labels, but such data receives a low marginal likelihood [WI20].

In general, we can distinguish the support of a model, i.e., the set of functions it can represent, from the distribution over that support, i.e., the inductive bias which leads it to prefer some functions over others. We would like to use models where the support is large, so we can capture the complexity of real-world data, but also where the inductive bias places probability mass on the kinds of functions we expect to see. If we succeed at this, the posterior will quickly converge on the true function after seeing a small amount of data. This idea is sketched in Figure 17.12.

⁴¹

⁴² 17.4.5 PAC-Bayes

⁴³

⁴⁴ **PAC-Bayes** [McA99; LC02; Gue19; Alq21; GSZ21] provides a promising mechanism to derive non-vacuous generalization bounds for large *stochastic networks* [Ney+17; NBS18; DR17], with parameters sampled from a probability distribution. In particular, the difference between the train

⁴⁷

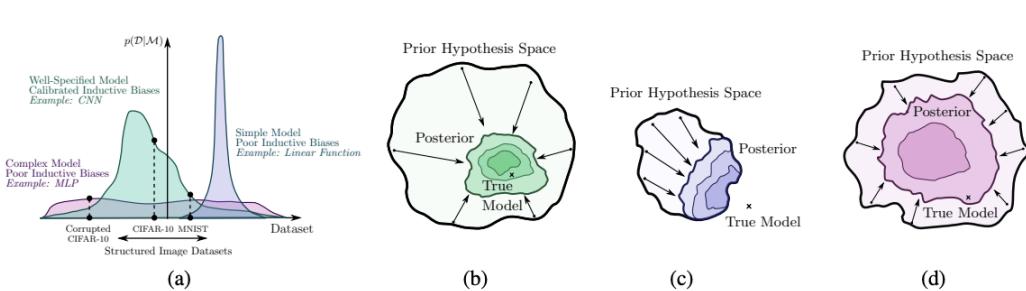


Figure 17.12: Illustration of the behavior of different kinds of model families and the prior distribution they induce over datasets. (a) The purple model is a simple linear model that has small support, and can only represent a few kinds of datasets. The pink model is an unstructured MLP: this has support over a large range of datasets with a fairly uninformative (broad) prior. Finally the green model is a CNN; this has support over a large range of datasets but the prior is more concentrated on certain kinds of datasets that have compositional structure. (b) The posterior for the green model (CNN) rapidly collapses to the true model, since it is consistent with the data. (c) The posterior for the purple model (linear) also rapidly collapses, but to a solution which cannot represent the true model. (d) The posterior for the pink model (MLP) collapses very slowly (as a function of dataset size). From Figure 2 of [WI20]. Used with kind permission of Andrew Wilson.

error and the generalization error can be expressed as

$$\sqrt{\frac{D_{\text{KL}}(Q \parallel P) + c}{2(N - 1)}}, \quad (17.32)$$

where c is a constant and N is the number of training points. P is the prior distribution over the parameters and Q is an arbitrary distribution, which can be chosen to optimize the bound.

The perspective in this chapter is largely complementary, and in some ways orthogonal, to the PAC-Bayes literature. Our focus has been on Bayesian marginalization, particularly multi-modal marginalization, and a prescriptive approach to model construction. In contrast, PAC-Bayes bounds are about bounding the empirical risk of a single sample, rather than marginalization, and are not currently prescriptive: what we would do to improve the bounds, such as reducing the number of model parameters, or using highly compact priors, does not typically improve generalization. Moreover, while we have seen Bayesian model averaging over multimodal posteriors has a significant effect on generalization, it has a minimal logarithmic effect on PAC-Bayes bounds. In general, because the bounds are loose, albeit non-vacuous in some cases, there is often room to make modeling choices that improve PAC-Bayes bounds without improving generalization, making it hard to derive a prescription for model construction from the bounds.

17.4.6 Out-of-Distribution generalization for BNNs

Bayesian methods are often assumed to be more robust in the context of distribution shift (discussed in Chapter 19), because they capture more uncertainty than methods based on point estimation. However, there are some subtleties, some of which we discuss below.

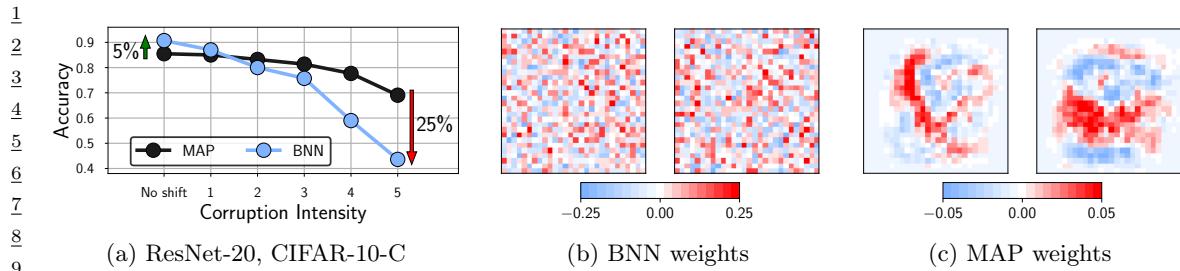


Figure 17.13: **Bayesian neural networks under covariate shift.** (a): Performance of a ResNet-20 on the pixelate corruption in CIFAR-10-C. For the highest degree of corruption, a Bayesian model average underperforms a MAP solution by 25% (44% against 69%) accuracy. See Izmailov et al. [Izm+21b] for details. (b): Visualization of the weights in the first layer of a Bayesian fully-connected network on MNIST sampled via HMC. (c): The corresponding MAP weights. We visualize the weights connecting the input pixels to a neuron in the hidden layer as a 28×28 image, where each weight is shown in the location of the input pixel it interacts with. This is Figure 1 of Izmailov et al. [Izm+21a].

17
18

17.4.6.1 BMA can give poor results with default priors

Many approximate inference methods, especially deep ensembles, are significantly less overconfident (more well calibrated) in the presence of some kinds of covariate shifts [Ova+19]. However, in [Izm+21b], it was noted that HMC, which arguably offers the most accurate approximation to the posterior, often works poorly under distribution shift.

Rather than an idiosyncracy of HMC, Izmailov et al. [Izm+21a] show this lack of robustness is a foundational issue of Bayesian model averaging under covariate shift, caused by degeneracies in the training data, and a poor choice of prior. As an illustrative special case, MNIST digits all have black corner pixels. Weights in the first layer of a neural network connected to these pixels are multiplied by zero, and thus can take any value without affecting the outputs of the network. Classical MAP training or deep ensembles of MAP solutions with a Gaussian prior will therefore drive these parameters to zero, since they don't help with the data fit, and the resulting network will be robust to corruptions on these pixels. On the other hand, the posterior for these parameters will be the same as the prior, and so a Bayesian model average will multiply corruptions by random numbers sampled from the prior, leading to degraded predictive performance.

Figure 17.13(b, c) visualizes this example, showing the first-layer weights of a fully-connected network for the MAP solution and a BNN posterior sample, on MNIST. The MAP weights corresponding to zero intensity pixels near the boundary are near zero, while the BNN weights look noisy, sampled from a Gaussian prior.

Izmailov et al. [Izm+21a] prove that this issue is a special case of a much more general problem, whenever there are linear dependencies in the input features of the training data, both for fully-connected and convolutional networks. In this case, the data live on a hyperplane. If a covariate or domain shift, moves orthogonal to this hyperplane, the posterior will be the same as the prior in the direction of the shift. The posterior model average will thus be highly vulnerable to shifts that do not particularly affect the underlying semantic structure of the problem (such as corruptions), whereas the MAP solution will be entirely robust to such shifts.

By introducing a prior over parameters which is aligned with the principal components of the

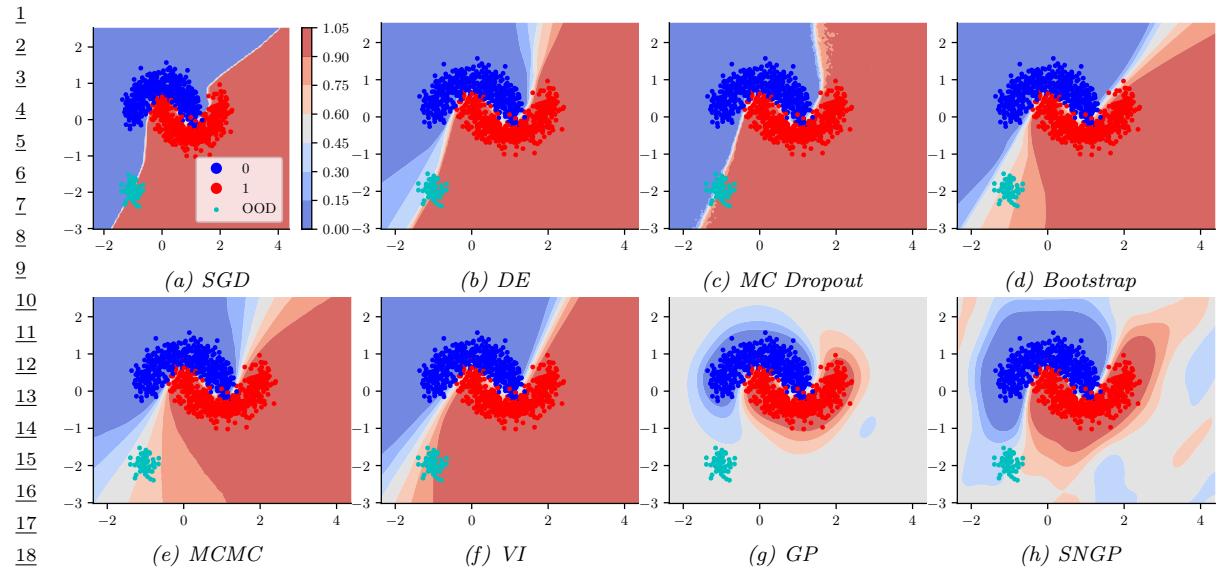


Figure 17.14: Predictions made by various (B)NNs when presented with the training data shown in blue and red. The green blob is an example of some OOD inputs. Methods are: (a) Standard SGD; (b) Deep Ensemble of 10 models with different random initializations; (c) MC Dropout with 50 samples; (d) Bootstrap training, where each of the 10 models is initialized identically but given different versions of the data, obtained by resampling with replacement; (e) MCMC using NUTS algorithm with 3000 warmup steps and 3000 samples; (f) Variational inference; (g) Gaussian process classifier using RBF kernel; (h) SNGP. The model is an MLP with 8,16,16,8 units in the hidden layers and ReLU activation. The output layer has 1 neuron with sigmoid activation. Generated by [makemoons_comparison.ipynb](#)

training inputs, we can substantially improve the generalization accuracy of Bayesian neural networks in out-of-distribution settings. Izmailov et al. [Izm+21a] propose the following *EmpCov* prior: $p(w^1) = \mathcal{N}(0, \alpha \Sigma + \epsilon I)$, where w^1 are the first layer weights, $\Sigma = \frac{1}{n-1} \sum_{i=1}^n x_i x_i^T$ is the empirical covariance of the training input features x_i , $\alpha > 0$ determines the scale of the prior, and ϵ is a small positive constant to ensure the covariance matrix is positive definite. With this improved prior they are able to obtain a method that is much more robust to distribution shift.

17.4.6.2 BNNs can be overconfident on OOD inputs

An important problem in practice is how a predictive model will behave when it is given an input that is “out of distribution” or OOD. Ideally we would like the model to express that it is not confident in its prediction, so that the system can abstain from predicting (see Section 19.3.3). Using “exact” inference methods, such as MCMC, for BNNs can give this behavior in some cases. For example, in Section 19.3.3.1 we showed that an MLP which was fit to MNIST using SGLD would be less overconfident than a point estimate (computed using SGD) when presented with inputs from fashion MNIST. However, this behavior does not always occur reliably.

To illustrate the problem, consider the 2d nonlinear binary classification dataset shown in Figure 17.14. In addition to the two training classes, we have highlighted (in green) a set of OOD inputs

1 that are far from the support of the training set. Intuitively we would expect the model to predict
 2 a probability of 0.5 (corresponding to “don’t know”) for such inputs that are far from the training
 3 set. However we see that the only methods that do so are the Gaussian Process (GP) classifier (see
 4 Section 18.4) and the SNGP model (Section 17.3.6), which contains a GP layer on top of the feature
 5 extractor.
 6

7 The lesson we learn from this simple example is that “being Bayesian” only helps if we are using a
 8 good hypothesis class. If we only consider a single MLP classifier, with standard Gaussian priors on
 9 the weights, it is extremely unlikely that we will learn the kind of compact decision boundary shown
 10 in Figure 17.14g, because that function has negligible support under our prior (c.f. Section 17.4.4).
 11 Instead we should embrace the power of Bayes to avoid overfitting and use as complex a model class
 12 as we can afford.

13

14 17.4.7 Model Selection for BNNs

15 Historically, the marginal likelihood (aka Bayesian evidence) has been used for selecting between
 16 trained neural architectures and learning hyperparameters [Mac92a]. It has recently seen a resurgence
 17 for this purpose, with scalable Laplace approximations [Imm+21; Dax+21]. Moreover, the variational
 18 lower bound to the marginal likelihood, the ELBO, is often used in VAEs to train millions of
 19 parameters in the decoder network [KW14].

20 However, [Lot+22] show many ways in which the question the marginal likelihood answers, “what
 21 is the likelihood of generating the training data from my prior?” different from the question we
 22 care about when comparing the generalization of trained models, “what is the probability that the
 23 posterior could generate withheld points from the data distribution?”. For example, eventually the
 24 posterior predictive will become insensitive to increases in prior variance, because the prior becomes
 25 weak, but the marginal likelihood will become increasingly negative — leaving room for the marginal
 26 likelihood to prefer models that have significantly worse generalization. They show how the marginal
 27 likelihood can be trivially modified to be more aligned with generalization for model selection, and
 28 hyperparameter learning. These questions apply to model selection in general, but particularly neural
 29 networks, which are often underspecified by the data, and require approximate inference.

30

31

32 17.5 Online inference

33

34 In Section 17.3, we have focused on batch or offline inference. However, an important application
 35 of Bayesian inference is in sequential settings, where the data arrives in a continuous stream, and
 36 the model has to “keep up”. This is called **sequential Bayesian inference**, and is one approach to
 37 **online learning** (see Section 19.7.5). In this section, we discuss some algorithmic approaches to
 38 this problem in the context of DNNs. These methods are widely used for continual learning, which
 39 we discuss Section 19.7.

40

41

42 17.5.1 Sequential Laplace for DNNs

43

44 In [RBB18b], they extended the Laplace method of Section 17.3.2 to the sequential setting. Specifically,
 45 let $p(\theta|\mathcal{D}_{1:t-1}) \approx \mathcal{N}(\theta|\mu_{t-1}, \Lambda_{t-1}^{-1})$ be the approximate posterior from the previous step; we assume
 46 the precision matrix is Kronecker factored. We now compute the new mean by solving the MAP

47

1 problem

2

$$\boldsymbol{\mu}_t = \operatorname{argmax} \log p(\mathcal{D}_t | \boldsymbol{\theta}) + \log p(\boldsymbol{\theta} | \mathcal{D}_{t-1}) \quad (17.33)$$

3

$$= \operatorname{argmax} \log p(\mathcal{D}_t | \boldsymbol{\theta}) - \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\mu}_{t-1}) \boldsymbol{\Lambda}_{t-1}^{-1} (\boldsymbol{\theta} - \boldsymbol{\mu}_{t-1}) \quad (17.34)$$

4

Once we have computed $\boldsymbol{\mu}_t$, we compute the approximate Hessian at this point, and get the new
5 posterior precision

6

$$\boldsymbol{\Lambda}_t = \lambda \mathbf{H}(\boldsymbol{\mu}_t) + \boldsymbol{\Lambda}_{t-1} \quad (17.35)$$

7

where $\lambda \geq 0$ is a weighting factor that trades off how much the model pays attention to the new data
8 vs old data.

9

Now suppose we use a diagonal approximation to the posterior prediction matrix. From Equation
10 (17.34), we see that this amounts to adding a quadratic penalty to each new MAP estimate, to
11 encourage it to remain close to the parameters from previous tasks. This approach is called **elastic**
12 **weight consolidation (EWC)** [Kir+17].

13 17.5.2 Extended Kalman Filtering for DNNs

14

In Section 29.7.2, we showed how Kalman filtering can be used to incrementally compute the
15 exact posterior for the weights of a linear regression model with known variance, i.e., we compute
16 $p(\boldsymbol{\theta} | \mathcal{D}_{1:t}, \sigma^2)$, where $\mathcal{D}_{1:t} = \{(\mathbf{u}_i, y_i) : i = 1 : t\}$ is the data seen so far, and

17

$$p(y_t | \mathbf{u}_t, \boldsymbol{\theta}, \sigma^2) = \mathcal{N}(y_t | \boldsymbol{\theta}^\top \mathbf{u}_t, \sigma^2) \quad (17.36)$$

18

is the linear regression likelihood. The application of KF to this model is known as recursive least
19 squares.

20 Now consider the case of nonlinear regression:

21

$$p(y_t | \mathbf{u}_t, \boldsymbol{\theta}, \sigma^2) = \mathcal{N}(y_t | f(\boldsymbol{\theta}, \mathbf{u}_t), \sigma^2) \quad (17.37)$$

22

where $f(\boldsymbol{\theta}, \mathbf{u}_t)$ is some nonlinear function, such as an MLP. We can use the extended Kalman filter
23 (Section 8.5.2) to approximately compute $p(\boldsymbol{\theta}_t | \mathcal{D}_{1:t}, \sigma^2)$, where $\boldsymbol{\theta}_t$ is the hidden state (see e.g., [SW89;
24 PF03]). To see this, note that we can set the dynamics model to the identity function, $f(\boldsymbol{\theta}_t) = \boldsymbol{\theta}_t$, so
25 the parameters are propagated through unchanged, and the observation model to the input-dependent
26 function $f(\boldsymbol{\theta}_t) = f(\boldsymbol{\theta}_t, \mathbf{u}_t)$. We set the observation noise to $\mathbf{R}_t = \sigma^2$, and the dynamics noise to
27 $\mathbf{Q}_t = q\mathbf{I}$, where q is a small constant, to allow the parameters to slowly drift according to artificial
28 **process noise**. (In practice it can be useful to anneal q from a large initial value to something near
29 0.)

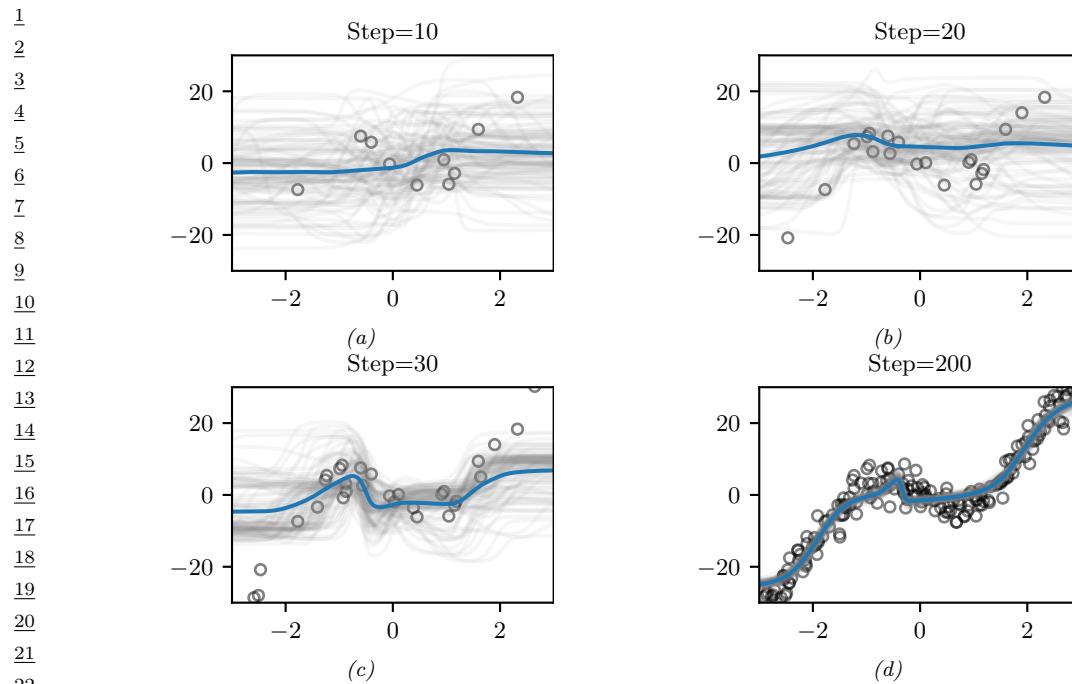
30 17.5.2.1 Example

31

We now give an example of this process in action. We sample a synthetic dataset from the true
32 function

33

$$h^*(u) = x - 10 \cos(u) \sin(u) + u^3 \quad (17.38)$$



23 *Figure 17.15: Sequential Bayesian inference for the parameters of an MLP using the extended Kalman
24 filter. We show results after seeing the first 10, 20, 30 and 200 observations. (For a video of this, see
25 <https://bit.ly/3wXnWaM>.) Generated by `ekf_mlp.ipynb`.*

28 and add Gaussian noise with $\sigma = 3$. We then fit this with an MLP with one hidden layer with H
29 hidden units, so the model has the form
30

$$31 \quad f(\boldsymbol{\theta}, \mathbf{u}) = \mathbf{W}_2 \tanh(\mathbf{W}_1 \mathbf{u} + \mathbf{b}_1) + \mathbf{b}_2 \quad (17.39)$$

33 where $\mathbf{W}_1 \in \mathbb{R}^{H \times 1}$, $\mathbf{b}_1 \in \mathbb{R}^H$, $\mathbf{W}_2 \in \mathbb{R}^{1 \times H}$, $\mathbf{b}_2 \in \mathbb{R}^1$. We set $H = 6$, so there are $D = 19$ parameters
34 in total.

35 Given the data, we sequentially compute the posterior, starting from a vague Gaussian prior,
36 $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta} | \mathbf{0}, \boldsymbol{\Sigma}_0)$, where $\boldsymbol{\Sigma}_0 = 100\mathbf{I}$. (In practice we cannot start from the prior mean, which is
37 $\boldsymbol{\theta}_0 = \mathbf{0}$, since linearizing the model around this point results in a zero gradient, so we use an initial
38 random sample for $\boldsymbol{\theta}_0$.) The results are shown in Figure 17.15. We can see that the model adapts
39 to the data, without having to specify any learning rate. In addition, we see that the predictions
40 become gradually more confident, as the posterior concentrates on the MLE.
41

42 17.5.2.2 Setting the variance terms 43

44 In the above example, we set the variance terms by hand. In general we need to estimate the noise
45 variance σ , which determines \mathbf{R}_t and hence the learning rate, as well as the strength of the prior $\boldsymbol{\Sigma}_0$,
46 which controls the amount of regularization. Some methods for doing this are discussed in [FNG00].
47

17.5.2.3 Reducing the computational complexity

The naive EKF method described above takes $O(N_z^3)$ time, which is prohibitive for large neural networks. A simple approximation, known as the **decoupled EKF**, was proposed in [PF91; SPD92] (see [PF03] for a review). This partitions the weights into G groups or blocks, and estimates the relevant matrices for each group g independently. If $G = 1$, this reduces the standard global EKF. If we put each weight into its own group, we get a fully diagonal approximation. In practice this does not work any better than SGD, since it ignores correlations between the parameters. A useful compromise is to put all the weights corresponding to each neuron into its own group; this is called “node decoupled EKF”.

Another approach to increasing computational efficiency is to leverage the fact that the effective dimensionality of a DNN is often quite low (see Section 17.4.3). Indeed we can approximate the model parameters by using a low dimensional vector of coefficients that specify the point in a linear manifold corresponding to weight space; the basis set defining this linear manifold can either be chosen randomly [Li+18b; Lar+21], or can be estimated using PCA applied to the SGD iterates [Izm+19]. We can exploit this observation to perform EKF in this low-dimensional subspace, which significantly speeds up inference, as discussed in [DMKM22].

17.5.3 Assumed Density Filtering for DNNs

In Section 8.10.3, we discussed how to use assumed density filtering (ADF) to perform online (binary) logistic regression. In this section, we generalize this to nonlinear predictive models, such as DNNs. The key is to perform Gaussian moment matching of the hidden activations at each layer of the model. This provides an alternative to the EKF approach in Section 17.5.2, which is based on linearization of the network.

We will assume the following likelihood:

$$p(\mathbf{y}_t | \mathbf{u}_t, \mathbf{w}_t) = \text{Expfam}(\mathbf{y}_t | \ell^{-1}(f(\mathbf{u}_t; \mathbf{w}_t))) \quad (17.40)$$

where $f(\mathbf{x}; \mathbf{w})$ is the DNN, ℓ^{-1} is the inverse link function, and $\text{Expfam}()$ is some exponential family distribution. For example, if f is linear and we are solving a binary classification problem, we can write

$$p(y_t | \mathbf{u}_t, \mathbf{w}_t) = \text{Ber}(y_t | \sigma(\mathbf{u}_t^\top \mathbf{w}_t)) \quad (17.41)$$

We discussed using ADF to fit this model in Section 8.10.3.

In [HLA15b], they propose **Probabilistic backpropagation (PBP)**, which is an instance of ADF applied to MLPs. The basic idea is to approximate the posterior over the weights in each layer using a fully factorized distribution

$$p(\mathbf{w}_t | \mathcal{D}_{1:t}) \approx p_t(\mathbf{w}_t) = \prod_{l=1}^L \prod_{i=1}^{D_l} \prod_{j=1}^{D_{l-1}+1} \mathcal{N}(w_{ijl} | \mu_{ijl}^t, \tau_{ijl}^t) \quad (17.42)$$

where L is the number of layers, and D_l is the number of neurons in layer l . (The **expectation backpropagation** algorithm of [SHM14] is a special case of this, where the variances are fixed to $\tau = 1$.)

1 Suppose the parameters are static, so $\mathbf{w}_t = \mathbf{w}_{t-1}$. Then the new posterior, after conditioning on
2 the t 'th observation, is given by
3

4

$$\hat{p}_t(\mathbf{w}) = \frac{1}{Z_t} p(\mathbf{y}_t | \mathbf{u}_t, \mathbf{w}) \mathcal{N}(\mathbf{w} | \boldsymbol{\mu}^{t-1}, \boldsymbol{\Sigma}^{t-1}) \quad (17.43)$$

5

6 where $\boldsymbol{\Sigma}^{t-1} = \text{diag}(\boldsymbol{\tau}^{t-1})$. We then project $\hat{p}_t(\mathbf{w})$ instead the space of factored Gaussians to compute
7 the new (approximate) posterior, $p_t(\mathbf{w})$. This can be done by computing the following means and
8 variances [Min01a]:
9

10

$$\mu_{ijl}^t = \mu_{ijl}^{t-1} + \tau_{ijl}^{t-1} \frac{\partial \ln Z_t}{\partial \mu_{ijl}^{t-1}} \quad (17.44)$$

11

12

$$\tau_{ijl}^t = \tau_{ijl}^{t-1} - (\tau_{ijl}^{t-1})^2 \left[\left(\frac{\partial \ln Z_t}{\partial \mu_{ijl}^{t-1}} \right)^2 - 2 \frac{\partial \ln Z_t}{\partial \tau_{ijl}^{t-1}} \right] \quad (17.45)$$

13

14 In the forwards pass, we compute Z_t by propagating the input \mathbf{u}_t through the model. Since we have
15 a Gaussian distribution over the weights, instead of a point estimate, this induces an (approximately)
16 Gaussian distribution over the values of the hidden units. For certain kinds of activation functions
17 (such as ReLU), the relevant integrals (to compute the means and variances) can be solved analytically,
18 as in GP-neural networks (Section 18.7). The result is that we get a Gaussian distribution over the
19 final layer of the form $\mathcal{N}(\boldsymbol{\eta}_t | \boldsymbol{\mu}, \boldsymbol{\Sigma})$, where $\boldsymbol{\eta}_t = f(\mathbf{u}_t; \mathbf{w}_t)$ is the output of the neural network before
20 the GLM link function induced by $p_t(\mathbf{w}_t)$. Hence we can approximate the partition function using
21

22

$$Z_t \approx \int p(\mathbf{y}_t | \boldsymbol{\eta}_t) \mathcal{N}(\boldsymbol{\eta}_t | \boldsymbol{\mu}, \boldsymbol{\Sigma}) d\boldsymbol{\eta}_t \quad (17.46)$$

23

24 We now discuss how to compute this integral. In the case of probit classification, with $y \in \{-1, +1\}$,
25 we have $p(y|\mathbf{x}, \mathbf{w}) = \Phi(y\eta)$, where Φ is the cdf of the standard normal. We can then use the following
26 analytical result
27

28

$$\int \Phi(y\eta) \mathcal{N}(h|\mu, \sigma) d\eta = \Phi \left(\frac{y\mu}{\sqrt{1+\sigma}} \right) \quad (17.47)$$

29

30 In the case of logistic classification, with $y \in \{0, 1\}$, we have $p(y|\mathbf{x}, \mathbf{w}) = \text{Ber}(y|\sigma(\eta))$; in this case,
31 we can use the probit approximation from Section 15.3.5. For the multiclass case, where $\mathbf{y} \in \{0, 1\}^C$
32 (one-hot encoding), we have $p(\mathbf{y}|\mathbf{x}, \mathbf{w}) = \text{Cat}(\mathbf{y}|\text{softmax}(\boldsymbol{\eta}))$. A variational lower bound to $\log Z_t$ for
33 this case is given in [GDFY16].
34

35 Once we have computed Z_t , we can take gradients and update the Gaussian posterior moments,
36 before moving to the next step.
37

38

39 17.5.4 Online variational inference for DNNs

40

41 A natural approach to online learning is to use variational inference, where the prior is the posterior
42 from the previous step. This is known as **streaming variational Bayes** [Bro+13]. In more detail,
43

1 at step t , we compute
 2

$$\underline{\psi}_t = \operatorname{argmin}_{\psi} \underbrace{\mathbb{E}_{q(\theta|\psi)} [\ell_t(\theta)] + D_{\text{KL}}(q(\theta|\psi) \| q(\theta|\psi_{t-1}))}_{-\mathcal{L}_t(\psi)} \quad (17.48)$$

$$= \operatorname{argmin}_{\psi} \mathbb{E}_{q(\theta|\psi)} [\ell_t(\theta) + \log q(\theta|\psi) - \log q(\theta|\psi_{t-1})] \quad (17.49)$$

8 where $\ell_t(\theta) = -\log p(\mathcal{D}_t|\theta)$ is the negative log likelihood (or, more generally, some loss function) of
 9 the data batch at step t .
 10

11 When applied to DNNs, this approach is called **variational continual learning** or **VCL** [Ngu+18].
 12 (We discuss continual learning in Section 19.7.) An efficient implementation of this, known as **FOO-**
 13 **VB** (“Fixed-point Operator for Online Variational Bayes”) is given in [Zen+21].

14 One problem with the VCL objective in Equation (17.48) is that the KL term can cause the
 15 model to become too sparse, which can prevent the model from adapting or learning new tasks.
 16 This problem is called **variational overpruning** [TT17]. More precisely, the reason this happens
 17 as is as follows: some weights might not be needed to fit a given dataset, so their posterior will be
 18 equal to the prior; but sampling from these high-variance weights will add noise to the likelihood; to
 19 reduce this, the optimization method will prefer to set the bias term to a large negative value, so
 20 the corresponding unit is “turned off”, and thus has no effect on the likelihood. Unfortunately, these
 21 “dead units” become stuck, so there is not enough network capacity to learn the next task.
 22

23 In [LST21], they propose a solution to this, known as **generalized variational continual**
 24 **learning** or GVCL. The first step is to downweight the KL term by a factor $\beta < 1$ to get
 25

$$\mathcal{L}_t = \mathbb{E}_{q(\theta|\psi)} [\ell_t(\theta)] + \beta D_{\text{KL}}(q(\theta|\psi) \| q(\theta|\psi_{t-1})) \quad (17.50)$$

26 Interestingly, one can show that in the limit of $\beta \rightarrow 0$, this recovers several standard methods that
 27 use a Laplace approximation based on the Hessian. In particular if we use a diagonal variational
 28 posterior, this reduces to online EWC method of [Sch+18]; if we use a block-diagonal and Kronecker
 29 factored posterior, this reduces to the online structured Laplace method of [RBB18b]; and if we use
 30 a low-rank posterior precision matrix, this reduces to the SOLA method of [Yin+20].
 31

32 The second step is to replace the prior and posterior by using tempering, which is useful when
 33 the model is misspecified, as discussed in Section 17.3.11. In the case of Gaussians, raising the
 34 distribution to the power λ is equivalent to tempering with a temperature of $\tau = 1/\lambda$, which is the
 35 same as scaling the covariance by λ^{-1} . Thus the GVCL objective becomes
 36

$$\mathcal{L}_t = \mathbb{E}_{q(\theta|\psi)} [\ell_t(\theta)] + \beta D_{\text{KL}}(q(\theta|\psi)^\lambda \| q(\theta|\psi_{t-1})^\lambda) \quad (17.51)$$

37 This can be optimized using SGD, assuming the posterior is reparameterizable (see Section 10.3.3).
 38

40 17.6 Hierarchical Bayesian neural networks

41 In some problems, we have multiple related datasets, such as a set of medical images from different
 42 hospitals. Some aspects of the data (e.g., the shape of healthy vs diseased cells) is generally the same
 43 across datasets, but other aspects may be unique or idiosyncratic (e.g., each hospital may use a
 44 different colored die for staining). To model this, we can use a hierarchical Bayesian model, in which
 45 we allow the parameters for each dataset to be different (to capture random effects), while coming from
 46 47

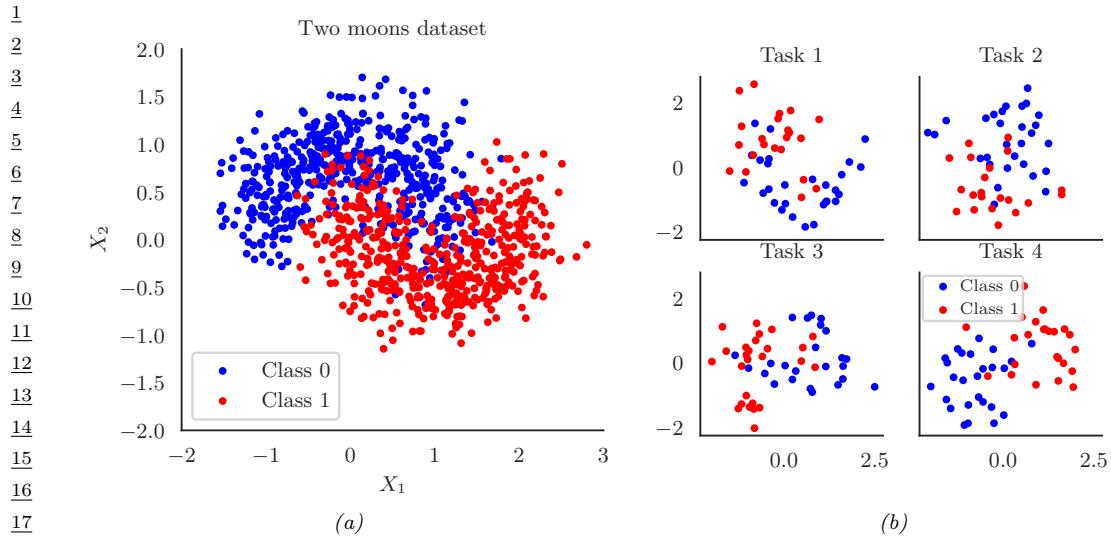


Figure 17.16: (a) Two moons synthetic dataset. (b) Multi-task version, where we rotate the data to create 18 related tasks (groups). Each dataset has 50 training and 50 test points. Here we show the first 4 tasks. Generated by [bnn_hierarchical.ipynb](#).

a common prior (to capture shared effects). This is the setup we considered in Section 15.5, where we discuss hierarchical Bayesian GLMs. In this section, we extend this to nonlinear predictors based on neural networks. (The setup is very similar to domain generalization, discussed in Section 19.6.2, except here we care about performance on all the domains, not just a held-out target domain.)

17.6.1 Example: multi-moons classification

In this section, we consider an example² where we want to solve multiple related nonlinear binary classification problems coming from J different environments or distributions. We assume that each environment has its own unique decision boundary $p(y|\mathbf{x}, \mathbf{w}^j)$, so this is a form of concept shift (see Section 19.2.3). However we assume the overall shape of each boundary is similar to a common shared boundary, denote $p(y|\mathbf{x}, \mathbf{w}^0)$. We only have a small number N_j of examples from each environment, $\mathcal{D}^j = \{(\mathbf{x}_n^j, y_n^j) : n = 1 : N_j\}$, but we can utilise their common structure to do better than fitting J separate models.

To illustrate this, we create some synthetic 2d data for the $J = 18$ tasks. We start with the two-moons dataset, illustrated in Figure 17.16a. Each task is obtained by rotating the 2d inputs by a different amount, to create 18 related classification problems (see Figure 17.16b). See Figure 17.16b for the training data for 4 tasks.

To handle the nonlinear decision boundary, we use a multilayer perceptron. Since the dataset is low-dimensional (2d input), we use a shallow model with just 2 hidden layers, each with 5 neurons. We could fit a separate MLP to each task, but since we have limited data per task ($N_j = 50$ examples

² 2. This example is from https://tweicki.io/blog/2018/08/13/hierarchical_bayesian_neural_network/. For a real-world example of a similar approach applied to a gesture recognition task, see [Jos+17].

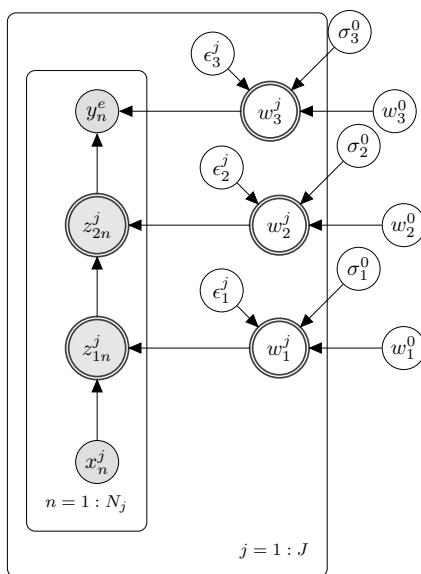


Figure 17.17: Illustration of a hierarchical Bayesian MLP with 2 hidden layers. There are J different models, each with N_j observed samples, and a common set of global shared parent parameters denoted with the 0 superscript. Nodes which are shaded are observed. Nodes with double ringed circles are deterministic functions of their parents.

for training), this works poorly, as we show below. We could also pool all the data and fit a single model, but this does even worse, since the datasets come from different underlying distributions, so mixing the data together from different “concepts” confuses the model. Instead we adopt a hierarchical Bayesian approach.

Our modeling assumptions are shown in Figure 17.17. In particular, we assume the weight from unit i to unit k in layer l for environment j , denoted $w_{i,k,l}^j$, comes from a common prior value $w_{i,k,l}^0$, with a random offset. We use the non-centered parameterization from Section 12.6.5 to write

$$w_{i,k,l}^j = w_{i,k,l}^0 + \epsilon_{i,k,l}^j \times \sigma_l^0 \quad (17.52)$$

where $\epsilon_{i,k,l}^j \sim \mathcal{N}(0, 1)$. By allowing a different σ_l^0 per layer l , we let the model control the degree of shrinkage to the prior for each layer separately. (We could also make the σ_l^j parameters be environment specific, which would allow for different amounts of distribution shift from the common parent.) For the hyper-parameters, we put $\mathcal{N}(0, 1)$ priors on $w_{i,k,l}^0$, and $\mathcal{N}_+(1)$ priors on σ_l^0 .

We compute the posterior $p(\epsilon_{1:L}^{1:J}, \mathbf{w}_{1:L}^0, \sigma_{1:L}^0 | \mathcal{D})$ using HMC (Section 12.5). We then evaluate this model using a fresh set of labeled samples from each environment. The average classification accuracy on the train and test sets for the non-hierarchical model (one MLP per environment, fit separately) is 86% and 83%. For the hierarchical model, this improves to 91% and 89% respectively.

To see why the hierarchical model works better, we will plot the posterior predictive distribution in 2d. Figure 17.18 shows the results for the non-hierarchical models; we see that the method fails to

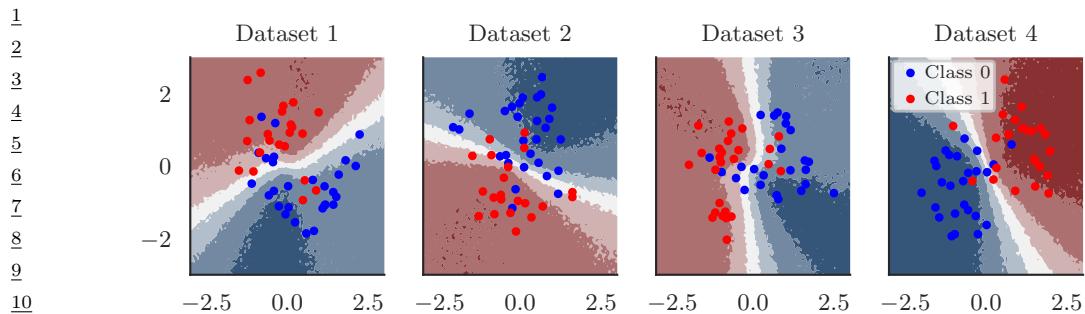


Figure 17.18: Results of fitting separate MLPs on each dataset. Generated by [bnn_hierarchical.ipynb](#).

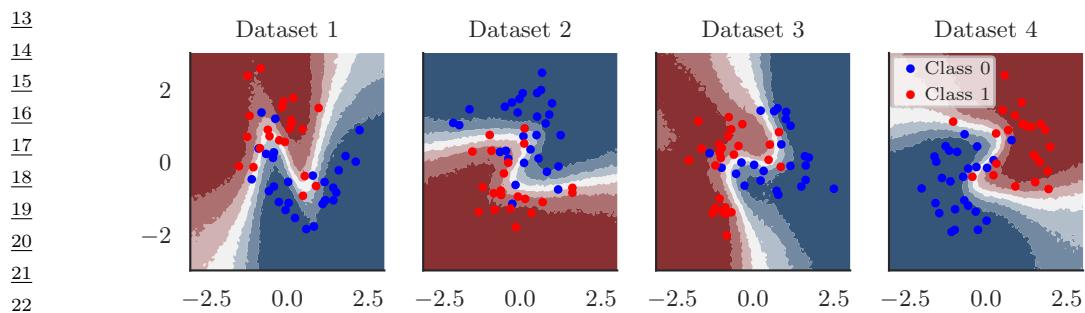


Figure 17.19: Results of fitting hierarchical MLP on all datasets jointly. Generated by [bnn_hierarchical.ipynb](#).

learn the common underlying Z-shaped decision boundary. By contrast, Figure 17.19 shows that the hierarchical method has correctly recovered the common pattern, while still allowing group variation.

29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

18 Gaussian processes

This chapter is co-authored with Andrew Wilson.

18.1 Introduction

Deep neural networks are a family of flexible function approximators of the form $f(\mathbf{x}; \boldsymbol{\theta})$, where the dimensionality of $\boldsymbol{\theta}$ (i.e., the number of parameters) is fixed, and independent of the size N of the training set. However, such parametric models can overfit when N is small, and can underfit when N is large, due to their fixed capacity. In order to create models whose capacity automatically adapts to the amount of data, we turn to **nonparametric models**.

There are many approaches to building nonparametric models for classification and regression (see e.g., [Was06]). In this chapter, we consider a Bayesian approach in which we represent uncertainty about the input-output mapping f by defining a prior distribution over functions, and then updating it given data. In particular, we will use a **Gaussian process** to represent the prior $p(f)$; we then use Bayes rule to derive the posterior $p(f|\mathcal{D})$, which is another GP, as we explain below. More details on GPs can be found the excellent book [RW06], as well as the interative tutorial at <https://distill.pub/2019/visual-exploration-gaussian-processes>. See also Chapter 31 for other examples of Bayesian nonparametric models.

18.1.1 GPs: What and why?

To explain GPs in more detail, recall that a Gaussian random vector of length N , $\mathbf{f} = [f_1, \dots, f_N]$, is defined by its mean $\boldsymbol{\mu} = \mathbb{E}[\mathbf{f}]$ and its covariance $\boldsymbol{\Sigma} = \text{Cov}[\mathbf{f}]$. Now consider a function $f : \mathcal{X} \rightarrow \mathbb{R}$ evaluated at a set of inputs, $\mathbf{X} = \{\mathbf{x}_n \in \mathcal{X}\}_{n=1}^N$. Let $\mathbf{f}_X = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)]$ be the set of unknown function values at these points. If \mathbf{f}_X is jointly Gaussian for any set of $N \geq 1$ points, then we say that $f : \mathcal{X} \rightarrow \mathbb{R}$ is a **Gaussian process**. Such a process is defined by its **mean function** $m(\mathbf{x}) \in \mathbb{R}$ and a **covariance function**, $\mathcal{K}(\mathbf{x}, \mathbf{x}') \geq 0$, which is any positive definite **Mercer kernel** (see Section 18.2). For example, we might use an RBF kernel of the form $\mathcal{K}(\mathbf{x}, \mathbf{x}') \propto \exp(-\|\mathbf{x} - \mathbf{x}'\|^2)$ (see Section 18.2.1.1 for details).

We denote the corresponding GP by

$$f(\mathbf{x}) \sim GP(m(\mathbf{x}), \mathcal{K}(\mathbf{x}, \mathbf{x}')) \tag{18.1}$$

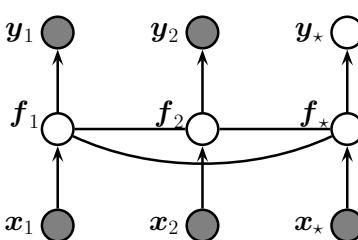


Figure 18.1: A Gaussian process for 2 training points, \mathbf{x}_1 and \mathbf{x}_2 , and 1 testing point, \mathbf{x}_ , represented as a graphical model representing $p(\mathbf{y}, \mathbf{f}_X | \mathbf{X}) = \mathcal{N}(\mathbf{f}_X | m(\mathbf{X}), \mathcal{K}(\mathbf{X})) \prod_i p(y_i | f_i)$. The hidden nodes $f_i = f(\mathbf{x}_i)$ represent the value of the function at each of the data points. These hidden nodes are fully interconnected by undirected edges, forming a Gaussian graphical model; the edge strengths represent the covariance terms $\Sigma_{ij} = \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)$. If the test point \mathbf{x}_* is similar to the training points \mathbf{x}_1 and \mathbf{x}_2 , then the value of the hidden function f_* will be similar to f_1 and f_2 , and hence the predicted output y_* will be similar to the training values y_1 and y_2 .*

18

19

20 where

$$21 \quad m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})] \tag{18.2}$$

$$23 \quad \mathcal{K}(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))^\top] \tag{18.3}$$

24 This means that, for any finite set of points $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, we have

$$26 \quad p(\mathbf{f}_X | \mathbf{X}) = \mathcal{N}(\mathbf{f}_X | \boldsymbol{\mu}_X, \mathbf{K}_{X,X}) \tag{18.4}$$

28 where $\boldsymbol{\mu}_X = (m(\mathbf{x}_1), \dots, m(\mathbf{x}_N))$ and $\mathbf{K}_{X,X}(i, j) \triangleq \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)$.

29 A GP can be used to define a prior over functions. We can evaluate this prior at any set of points
30 we choose. However, to learn about the function from data, we have to update this prior with a
31 likelihood function. We typically assume we have a set of N iid observations $\mathcal{D} = \{(\mathbf{x}_i, y_i) : i = 1 : N\}$,
32 where $y_i \sim p(y|f(\mathbf{x}_i))$, as shown in Figure 18.1. If we use a Gaussian likelihood, we can compute the
33 posterior $p(f|\mathcal{D})$ in closed form, as we discuss in Section 18.3. For other kinds of likelihoods, we will
34 need to use approximate inference, as we discuss in Section 18.4. In many cases f is not directly
35 observed, and instead forms part of a latent variable model, both in supervised and unsupervised
36 settings such as in Section 28.3.7.

37 The generalization properties of a Gaussian process are controlled by its covariance function
38 (kernel), which we describe in Section 18.2. These kernels live in a reproducing kernel Hilbert space
39 (RKHS), described in Section 18.3.7.1.

40 GPs were originally designed for spatial data analysis, where the input is 2d. This special case
41 is called **kriging**. However, they can be applied to higher dimensional inputs. In addition, while
42 they have been traditionally limited to small datasets, it is now possible to apply GPs to problems
43 with millions of points, with essentially exact inference. We discuss these scalability advances in
44 Section 18.5.

45 Moreover, while Gaussian processes have historically been considered smoothing interpolators, GPs
46 now routinely perform representation learning, through covariance function learning, and multilayer
47

models. These advances have clearly illustrated that GPs and neural networks are not competing, but complementary, and can be combined for better performance than would be achieved by deep learning alone. We describe GPs for representation learning in Section 18.6.

The connections between Gaussian processes and neural networks can also be further understood by considering infinite limits of neural networks that converge to Gaussian processes with particular covariance functions, which we describe in Section 18.7.

So Gaussian processes are non-parametric models which can scale and do representation learning. But why, in the age of deep learning, should we want to use a Gaussian process? There are several compelling reasons to prefer a GP, including:

- Gaussian processes typically provide well-calibrated predictive distributions, with a good characterization of epistemic (model) uncertainty — uncertainty arising from not knowing which of many solutions is correct. For example, as we move away from the data, there are a greater variety of consistent solutions, and so we expect greater uncertainty.
- Gaussian processes are often state-of-the-art for continuous regression problems, especially spatiotemporal problems, such as weather interpolation and forecasting. In regression, Gaussian process inference can also typically be performed in closed form.
- The marginal likelihood of a Gaussian process provides a powerful mechanism for flexible kernel learning. Kernel learning enables us to provide long-range extrapolations, but also tells us interpretable properties of the data that we didn't know before, towards scientific discovery.
- Gaussian processes are often used as a probabilistic surrogate for objectives in optimization, in a procedure known as **Bayesian optimization** (Section 6.8). To maximize an objective, we wish to move where there is a high expected value, but also to explore where we have large uncertainty. The ability for a Gaussian process to provide closed form inference in regression, in conjunction with high quality uncertainty representations, make them particularly impactful in this setting. Bayesian optimization has a wide range of applications, including A/B testing, experimental design, protein engineering, hyperparameter tuning, and AutoML. See Section 6.8 for details.

18.2 Mercer kernels

The generalization properties of Gaussian processes boil down to how we encode prior knowledge about the similarity of two input vectors. If we know that \mathbf{x}_i is similar to \mathbf{x}_j , then we can encourage the model to make the predicted output at both locations (i.e., $f(\mathbf{x}_i)$ and $f(\mathbf{x}_j)$) to be similar.

To define similarity, we introduce the notion of a **kernel function**. The word “kernel” has many different meanings in mathematics; here we consider a **Mercer kernel**, also called a **positive definite kernel**. This is any symmetric function $\mathcal{K} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}^+$ such that

$$\sum_{i=1}^N \sum_{j=1}^N \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) c_i c_j \geq 0 \quad (18.5)$$

for any set of N (unique) points $\mathbf{x}_i \in \mathcal{X}$, and any choice of numbers $c_i \in \mathbb{R}$. We assume $\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) > 0$, so that we can only achieve equality in the above equation if $c_i = 0$ for all i .

1 Another way to understand this condition is the following. Given a set of N datapoints, let us
2 define the **Gram matrix** as the following $N \times N$ similarity matrix:
3

$$\underline{4} \quad \mathbf{K} = \begin{pmatrix} \mathcal{K}(\mathbf{x}_1, \mathbf{x}_1) & \cdots & \mathcal{K}(\mathbf{x}_1, \mathbf{x}_N) \\ \underline{5} & \vdots & \\ \underline{6} \mathcal{K}(\mathbf{x}_{N_D}, \mathbf{x}_1) & \cdots & \mathcal{K}(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix} \quad (18.6)$$

7 We say that \mathcal{K} is a Mercer kernel iff the Gram matrix is positive definite for any set of (distinct)
8 inputs $\{\mathbf{x}_i\}_{i=1}^N$.
9

10 The most widely used kernel for real-valued inputs is the **radial basis function** or **RBF** kernel,
11 defined by
12

$$\underline{13} \quad \mathcal{K}(\mathbf{x}, \mathbf{x}'; \ell) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\ell^2}\right) \quad (18.7)$$

14 Here ℓ corresponds to the length-scale of the kernel, i.e., the distance over which we expect differences
15 to matter. This is also sometimes known as the **bandwidth** parameter. The RBF kernel measures
16 similarity between two vectors in \mathbb{R}^D using (scaled) Euclidean distance. In Section 18.2.1, we will
17 discuss several other kinds of kernel.
18

19 18.2.1 Some popular Mercer kernels

20 In the sections below, we describe some popular Mercer kernels. More details can be found at [Wil14]
21 and <https://www.cs.toronto.edu/~duvenaud/cookbook/>. See also Section 18.6 where we discuss
22 how to learn kernels from data.
23

24 18.2.1.1 Stationary kernels for real-valued vectors

25 For real-valued inputs, $\mathcal{X} = \mathbb{R}^D$, it is common to use **stationary kernels**, which are functions of
26 the form $\mathcal{K}(\mathbf{x}, \mathbf{x}') = \mathcal{K}(\mathbf{r})$, where $\mathbf{r} = \mathbf{x} - \mathbf{x}'$; thus the output only depends on the relative difference
27 between the inputs. Furthermore, in many cases, all that matters is the magnitude of the difference:
28

$$\underline{29} \quad r = \|\mathbf{r}\|_2 = \|\mathbf{x} - \mathbf{x}'\| \quad (18.8)$$

30 We give some examples below. (See also Figure 18.3 and Figure 18.4 for some visualizations of these
31 kernels.)
32

33 Squared exponential (RBF) kernel

34 The **squared exponential** (SE) kernel, also sometimes called the **exponentiated quadratic** kernel,
35 is defined as
36

$$\underline{37} \quad \mathcal{K}(r; \ell) = \exp\left(-\frac{r^2}{2\ell^2}\right) \quad (18.9)$$

38 Here ℓ corresponds to the length-scale of the kernel, i.e., the distance over which we expect differences
39 to matter.
40

41

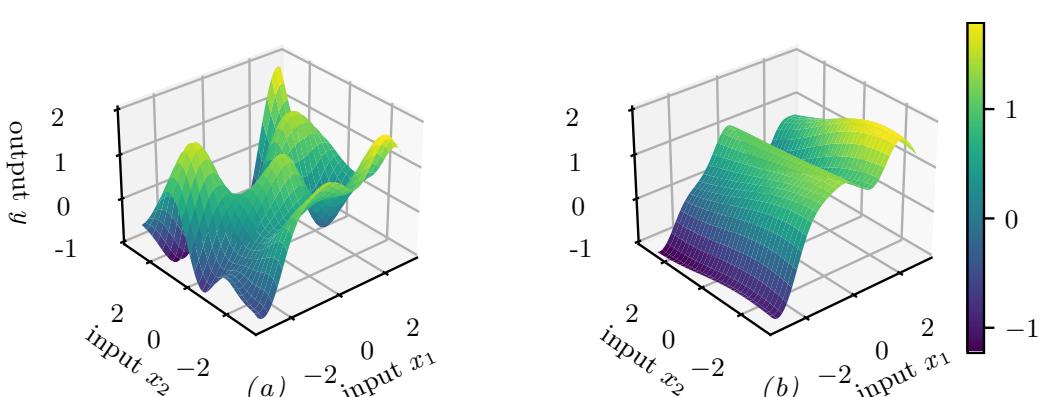


Figure 18.2: Function samples from a GP with an ARD kernel. (a) $\ell_1 = \ell_2 = 1$. Both dimensions contribute to the response. (b) $\ell_1 = 1, \ell_2 = 5$. The second dimension is essentially ignored. Adapted from Figure 5.1 of [RW06]. Generated by `gpr_demo_ard.ipynb`.

From Equation (18.8) we can rewrite this kernel as

$$\mathcal{K}(\mathbf{x}, \mathbf{x}'; \ell) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\ell^2}\right) \quad (18.10)$$

This is the RBF kernel we encountered earlier. It is also sometimes called the **Gaussian kernel**.

See Figure 18.3(f) and Figure 18.4(f) for a visualization in 1D.

ARD kernel

We can generalize the RBF kernel by replacing Euclidean distance with Mahalanobis distance, as follows:

$$\mathcal{K}(\mathbf{r}; \Sigma, \sigma^2) = \sigma^2 \exp\left(-\frac{1}{2} \mathbf{r}^\top \Sigma^{-1} \mathbf{r}\right) \quad (18.11)$$

where $\mathbf{r} = \mathbf{x} - \mathbf{x}'$. If Σ is diagonal, this can be written as

$$\mathcal{K}(\mathbf{r}; \ell, \sigma^2) = \sigma^2 \exp\left(-\frac{1}{2} \sum_{d=1}^D \frac{1}{\ell_d^2} r_d^2\right) = \prod_{d=1}^D \mathcal{K}(r_d; \ell_d, \sigma^{2/d}) \quad (18.12)$$

where

$$\mathcal{K}(r; \ell, \tau^2) = \tau^2 \exp\left(-\frac{1}{2} \frac{1}{\ell^2} r^2\right) \quad (18.13)$$

We can interpret σ^2 as the overall variance, and ℓ_d as defining the **characteristic length scale** of dimension d . If d is an irrelevant input dimension, we can set $\ell_d = \infty$, so the corresponding dimension will be ignored. This is known as **Automatic Relevance Determination** or **ARD** (Section 15.2.7). Hence the corresponding kernel is called the **ARD kernel**. See Figure 18.2 for an illustration of some 2d functions sampled from a GP using this prior.

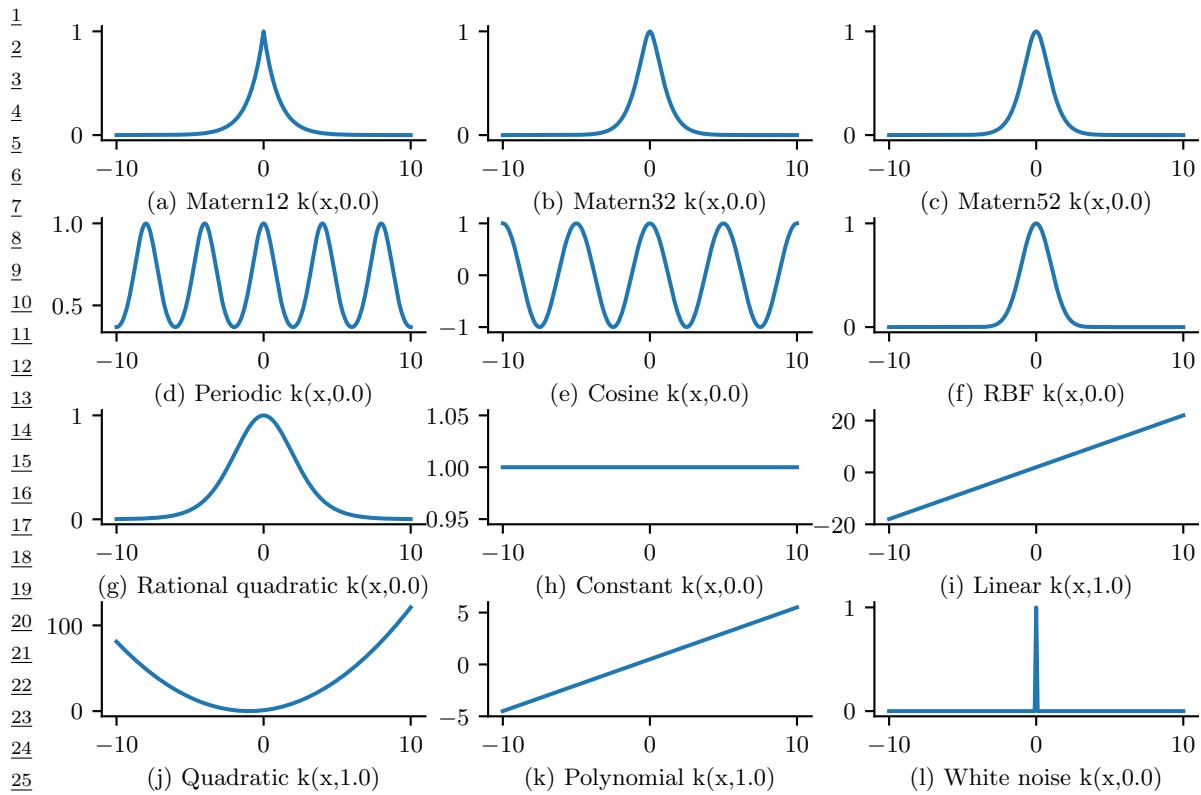


Figure 18.3: GP kernels evaluated at $k(x, 0)$ as a function of x . Generated by [gpKernelPlot.ipynb](#).

Matern kernels

The SE kernel gives rise to functions that are infinitely differentiable, and therefore are very smooth.

For many applications, it is better to use the **Matern kernel**, which gives rise to “rougher” functions, which can better model local “wiggles” without having to make the overall length scale very small.

The Matern kernel has the following form:

$$\mathcal{K}(r; \nu, \ell) = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu}r}{\ell} \right)^\nu K_\nu \left(\frac{\sqrt{2\nu}r}{\ell} \right) \quad (18.14)$$

where K_ν is a modified Bessel function and ℓ is the length scale. Functions sampled from this GP are k -times differentiable iff $\nu > k$. As $\nu \rightarrow \infty$, this approaches the SE kernel.

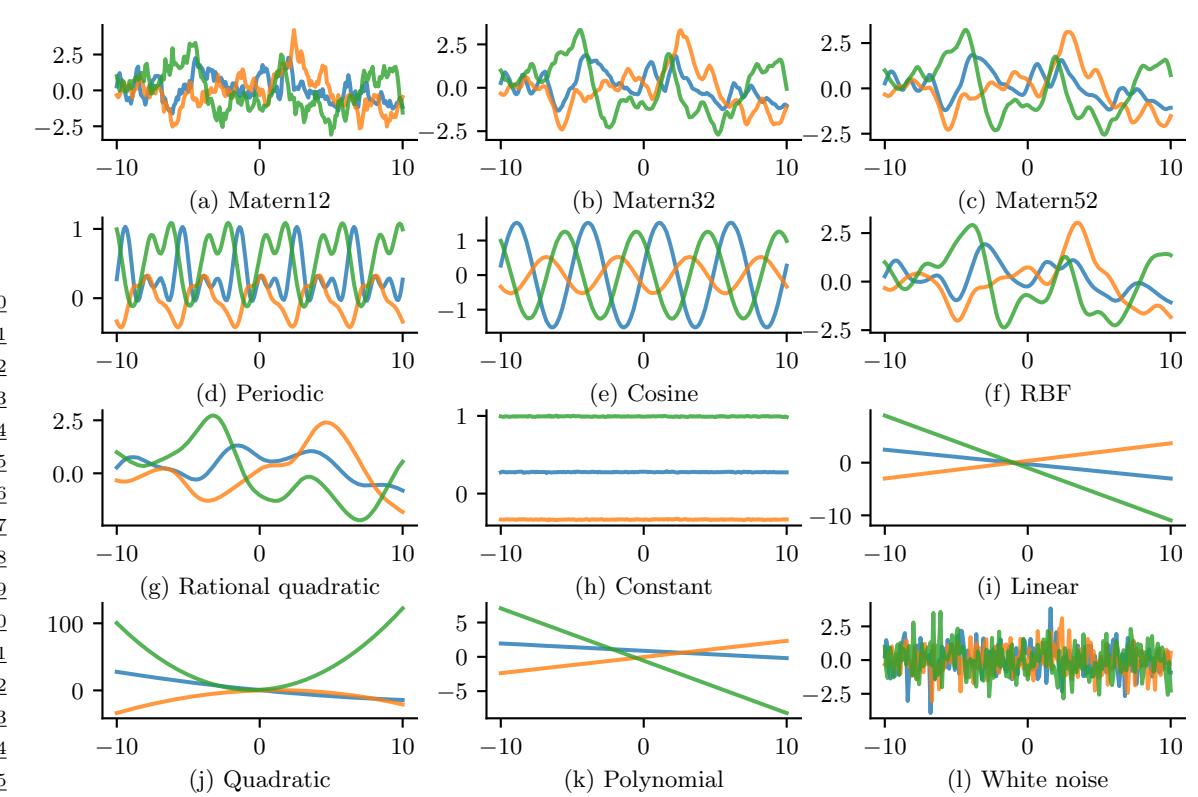


Figure 18.4: GP samples drawn using different kernels. Generated by [gpKernelPlot.ipynb](#).

For values $\nu \in \{\frac{1}{2}, \frac{3}{2}, \frac{5}{2}\}$, the function simplifies as follows:

$$\mathcal{K}(r; \frac{1}{2}, \ell) = \exp\left(-\frac{r}{\ell}\right) \quad (18.15)$$

$$\mathcal{K}(r; \frac{3}{2}, \ell) = \left(1 + \frac{\sqrt{3}r}{\ell}\right) \exp\left(-\frac{\sqrt{3}r}{\ell}\right) \quad (18.16)$$

$$\mathcal{K}(r; \frac{5}{2}, \ell) = \left(1 + \frac{\sqrt{5}r}{\ell} + \frac{5r^2}{3\ell^2}\right) \exp\left(-\frac{\sqrt{5}r}{\ell}\right) \quad (18.17)$$

See Figure 18.3(a-c) and Figure 18.4(a-c) for a visualization.

The value $\nu = \frac{1}{2}$ corresponds to the **Ornstein-Uhlenbeck process**, which describes the velocity of a particle undergoing Brownian motion. The corresponding function is continuous but not differentiable, and hence is very “jagged”.

1
2 **Periodic kernels**

3 One way to create a periodic 1d random function is to map x to the 2d space $\mathbf{u}(x) = (\cos(x), \sin(x))$,
4 and then use an SE kernel in \mathbf{u} -space:

5

$$\mathcal{K}(x, x') = \exp\left(-\frac{2 \sin^2((x - x')/2)}{\ell^2}\right) \quad (18.18)$$

6 which follows since $(\cos(x) - \cos(x'))^2 + (\sin(x) - \sin(x'))^2 = 4 \sin^2((x - x')/2)$. We can generalize this
7 by specifying the period p to get the **periodic kernel**, also called the **exp-sine-squared kernel**:

8

$$\mathcal{K}_{\text{per}}(r; \ell, p) = \exp\left(-\frac{2}{\ell^2} \sin^2(\pi \frac{r}{p})\right) \quad (18.19)$$

9 where p is the period and ℓ is the length scale. See Figure 18.3(d-e) and Figure 18.4(d-e) for a
10 visualization.

11 A related kernel is the **cosine kernel**:

12

$$\mathcal{K}(r; p) = \cos\left(2\pi \frac{r}{p}\right) \quad (18.20)$$

13
14 **Rational quadratic kernel**

15 We define the **rational quadratic** kernel to be

16

$$\mathcal{K}_{RQ}(r; \ell, \alpha) = \left(1 + \frac{r^2}{2\alpha\ell^2}\right)^{-\alpha} \quad (18.21)$$

17 We recognize this is proportional to a Student T density. Hence it can be interpreted as a scale
18 mixture of SE kernels of different characteristic lengths. In particular, let $\tau = 1/\ell^2$, and assume
19 $\tau \sim \text{Ga}(\alpha, \ell^2)$. Then one can show that

20

$$\mathcal{K}_{RQ}(r) = \int p(\tau|\alpha, \ell^2) \mathcal{K}_{SE}(r|\tau) d\tau \quad (18.22)$$

21 As $\alpha \rightarrow \infty$, this reduces to a SE kernel.

22 See Figure 18.3(g) and Figure 18.4(g) for a visualization.

23
24 **18.2.1.2 Making new kernels from old**

25 Given two valid kernels $\mathcal{K}_1(\mathbf{x}, \mathbf{x}')$ and $\mathcal{K}_2(\mathbf{x}, \mathbf{x}')$, we can create a new kernel using any of the following
26 methods:

27

$$\mathcal{K}(\mathbf{x}, \mathbf{x}') = c\mathcal{K}_1(\mathbf{x}, \mathbf{x}'), \text{ for any constant } c > 0 \quad (18.23)$$

28

$$\mathcal{K}(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})\mathcal{K}_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}'), \text{ for any function } f \quad (18.24)$$

29

$$\mathcal{K}(\mathbf{x}, \mathbf{x}') = q(\mathcal{K}_1(\mathbf{x}, \mathbf{x}')) \text{ for any function polynomial } q \text{ with nonneg. coef.} \quad (18.25)$$

30

$$\mathcal{K}(\mathbf{x}, \mathbf{x}') = \exp(\mathcal{K}_1(\mathbf{x}, \mathbf{x}')) \quad (18.26)$$

31

$$\mathcal{K}(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{A} \mathbf{x}', \text{ for any psd matrix } \mathbf{A} \quad (18.27)$$

For example, suppose we start with the linear kernel $\mathcal{K}(\mathbf{x}, \mathbf{x}') = \mathbf{x}\mathbf{x}'$. We know this is a valid Mercer kernel, since the corresponding Gram matrix is just the (scaled) covariance matrix of the data. From the above rules, we can see that the **polynomial kernel**

$$\mathcal{K}(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}')^M \quad (18.28)$$

is a valid Mercer kernel. This contains all monomials of order M . For example, if $M = 2$ and the inputs are 2d, we have

$$(\mathbf{x}^\top \mathbf{x}')^2 = (x_1 x'_1 + x_2 x'_2)^2 = (x_1 x'_1)^2 + (x_2 x'_2)^2 + (x_1 x'_1)(x_2 x'_2) \quad (18.29)$$

We can generalize this to contain all terms up to degree M by using the kernel $\mathcal{K}(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}' + c)^M$. For example, if $M = 2$ and the inputs are 2d, we have

$$\begin{aligned} (\mathbf{x}^\top \mathbf{x}' + 1)^2 &= (x_1 x'_1)^2 + (x_1 x'_1)(x_2 x'_2) + (x_1 x'_1) \\ &\quad + (x_2 x'_2)(x_1 x'_1) + (x_2 x'_2)^2 + (x_2 x'_2) \\ &\quad + (x_1 x'_1) + (x_2 x'_2) + 1 \end{aligned} \quad (18.30)$$

We can also use the above rules to establish that the Gaussian kernel is a valid kernel. To see this, note that

$$\|\mathbf{x} - \mathbf{x}'\|^2 = \mathbf{x}^\top \mathbf{x} + (\mathbf{x}')^\top \mathbf{x}' - 2\mathbf{x}^\top \mathbf{x}' \quad (18.31)$$

and hence

$$\mathcal{K}(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2 / 2\sigma^2) = \exp(-\mathbf{x}^\top \mathbf{x} / 2\sigma^2) \exp(\mathbf{x}^\top \mathbf{x}' / \sigma^2) \exp(-(\mathbf{x}')^\top \mathbf{x}' / 2\sigma^2) \quad (18.32)$$

is a valid kernel.

18.2.1.3 Combining kernels by addition and multiplication

We can also combine kernels using addition or multiplication:

$$\mathcal{K}(\mathbf{x}, \mathbf{x}') = \mathcal{K}_1(\mathbf{x}, \mathbf{x}') + \mathcal{K}_2(\mathbf{x}, \mathbf{x}') \quad (18.33)$$

$$\mathcal{K}(\mathbf{x}, \mathbf{x}') = \mathcal{K}_1(\mathbf{x}, \mathbf{x}') \times \mathcal{K}_2(\mathbf{x}, \mathbf{x}') \quad (18.34)$$

Multiplying two positive-definite kernels together always results in another positive definite kernel. This is a way to get a conjunction of the individual properties of each kernel, as illustrated in Figure 18.5.

In addition, adding two positive-definite kernels together always results in another positive definite kernel. This is a way to get a disjunction of the individual properties of each kernel, as illustrated in Figure 18.6.

For an example of combining kernels to forecast some timeseries data, see Section 18.8.1.

18.2.1.4 Kernels for structured inputs

Kernels are particularly useful when the inputs are structured objects, such as strings and graphs, since it is often hard to “featurize” variable-sized inputs. For example, we can define a **string kernel** which compares strings in terms of the number of n-grams they have in common [Lod+02; BC17].

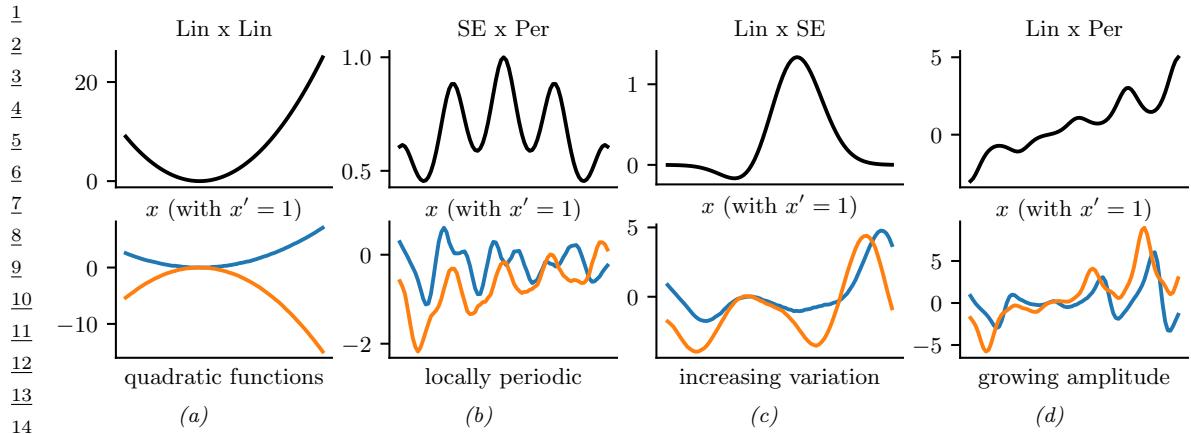


Figure 18.5: Examples of 1d structures obtained by multiplying elementary kernels. Top row shows $\mathcal{K}(x, x' = 1)$. Bottom row shows some functions sampled from $GP(f|0, \mathcal{K})$. Adapted from Figure 2.2 of [Duv14]. Generated by [combining_kernels_by_multiplication.ipynb](#).

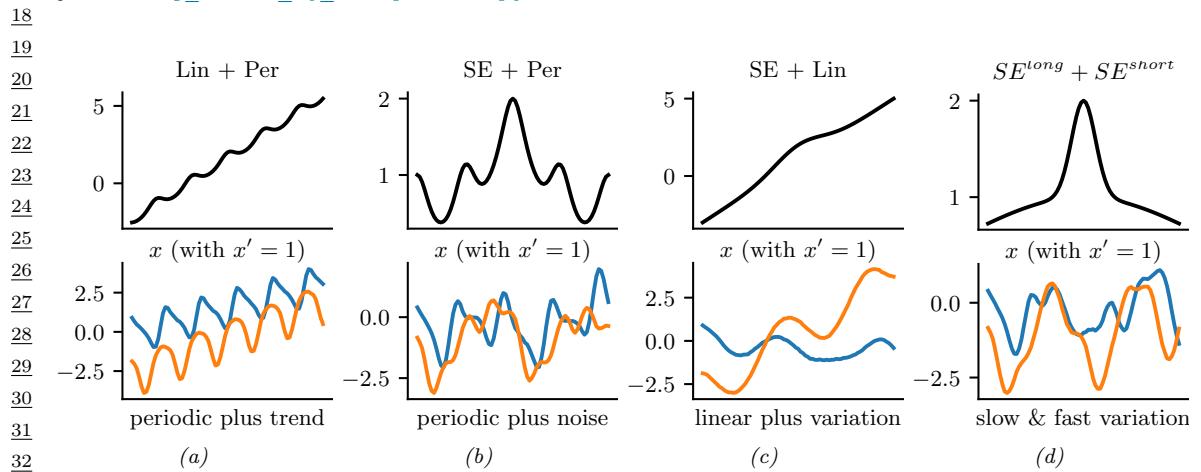


Figure 18.6: Examples of 1d structures obtained by summing elementary kernels. Top row shows $\mathcal{K}(x, x' = 1)$. Bottom row shows some functions sampled from $GP(f|0, \mathcal{K})$. Adapted from Figure 2.2 of [Duv14]. Generated by [combining_kernels_by_summation.ipynb](#).

We can also define kernels on graphs [KJM19]. For example, the **random walk kernel** conceptually performs random walks on two graphs simultaneously, and then counts the number of paths that were produced by both walks. This can be computed efficiently as discussed in [Vis+10]. For more details on graph kernels, see [KJM19].

For a review of kernels on structured objects, see e.g., [Gär03].

1

18.2.2 Mercer's theorem

2

3 Recall that any positive definite matrix \mathbf{K} can be represented using an eigendecomposition of the
4 form $\mathbf{K} = \mathbf{U}^\top \mathbf{\Lambda} \mathbf{U}$, where $\mathbf{\Lambda}$ is a diagonal matrix of eigenvalues $\lambda_i > 0$, and \mathbf{U} is a matrix containing
5 the eigenvectors. Now consider element (i, j) of \mathbf{K} :

6

$$k_{ij} = (\mathbf{\Lambda}^{\frac{1}{2}} \mathbf{U}_{:,i})^\top (\mathbf{\Lambda}^{\frac{1}{2}} \mathbf{U}_{:,j}) \quad (18.35)$$

7 where $\mathbf{U}_{:,i}$ is the i 'th column of \mathbf{U} . If we define $\phi(\mathbf{x}_i) = \mathbf{U}_{:,i}$, then we can write

8

$$k_{ij} = \sum_{m=1}^M \lambda_m \phi_m(\mathbf{x}_i) \phi_m(\mathbf{x}_j) \quad (18.36)$$

9

10 where M is the rank of the kernel matrix. Thus we see that the entries in the kernel matrix can be
11 computed by performing an inner product of some feature vectors that are implicitly defined by the
12 eigenvectors of the kernel matrix.

13 This idea can be generalized to apply to kernel functions, not just kernel matrices, as we now show.
14 First, we define an **eigenfunction** $\phi()$ of a kernel \mathcal{K} with eigenvalue λ wrt measure μ as a function
15 that satisfies

$$\int \mathcal{K}(\mathbf{x}, \mathbf{x}') \phi(\mathbf{x}) d\mu(\mathbf{x}) = \lambda \phi(\mathbf{x}') \quad (18.37)$$

16 We usually sort the eigenfunctions in order of decreasing eigenvalue, $\lambda_1 \geq \lambda_2 \geq \dots$. The eigenfunctions
17 are orthogonal wrt μ :

18

$$\int \phi_i(\mathbf{x}) \phi_j(\mathbf{x}) d\mu(\mathbf{x}) = \delta_{ij} \quad (18.38)$$

19

20 where δ_{ij} is the Kronecker delta. With this definition in hand, we can state **Mercer's theorem**.
21 Informally, it says that any positive definite kernel function can be represented as the following
22 infinite sum:

23

$$\mathcal{K}(\mathbf{x}, \mathbf{x}') = \sum_{m=1}^{\infty} \lambda_m \phi_m(\mathbf{x}) \phi_m(\mathbf{x}') \quad (18.39)$$

24

25 where ϕ_m are eigenfunctions of the kernel, and λ_m are the corresponding eigenvalues. This is the
26 functional analog of Equation (18.36).

27 A **degenerate kernel** has only a finite number of non-zero eigenvalues. In this case, we can
28 rewrite the kernel function as an inner product between two finite-length vectors. For example,
29 consider the **quadratic kernel** $\mathcal{K}(\mathbf{x}, \mathbf{x}') = \langle \mathbf{x}, \mathbf{x}' \rangle^2$. In 2d, we have

30

$$\mathcal{K}(\mathbf{x}, \mathbf{x}') = (x_1 x'_1 + x_2 x'_2)^2 = x_1^2(x'_1)^2 + 2x_1 x_2 x'_1 x'_2 + x_2^2(x'_2)^2 \quad (18.40)$$

31

32 If we define $\phi(x_1, x_2) = [x_1^2, \sqrt{2}x_1 x_2, x_2^2] \in \mathbb{R}^3$, then we can write this as $\mathcal{K}(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^\top \phi(\mathbf{x}')$.
33 Thus we see that this kernel is degenerate.

34

35 Now consider the RBF kernel. In this case, the corresponding feature representation is infinite
36 dimensional (see Section 18.2.3.1 for details). However, by working with kernel functions, we can
37 avoid having to deal with infinite dimensional vectors.

38 From the above, we see that we can replace inner product operations in an explicit (possibly infinite
39 dimensional) feature space with a call to a kernel function, i.e., we replace $\phi(\mathbf{x})^\top \phi(\mathbf{x})$ with $\mathcal{K}(\mathbf{x}, \mathbf{x}')$.
40 This is called the **kernel trick**.

41

1 **18.2.3 Kernels from Spectral Densities**

3 Consider the case of a **shift-invariant kernel**, also known as a *stationary kernel*, which satisfies
4 $\mathcal{K}(\mathbf{x}, \mathbf{x}') = \mathcal{K}(\boldsymbol{\delta})$, where $\boldsymbol{\delta} = \mathbf{x} - \mathbf{x}'$, for $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$. Let us further assume that $\mathcal{K}(\boldsymbol{\delta})$ is positive definite.
5 In this case, **Bochner's theorem** tells us that we can represent $\mathcal{K}(\boldsymbol{\delta})$ by its Fourier transform:
6

7
$$\mathcal{K}(\boldsymbol{\delta}) = \int_{\mathbb{R}^d} p(\boldsymbol{\omega}) e^{j\boldsymbol{\omega}^\top \boldsymbol{\delta}} d\boldsymbol{\omega} \quad (18.41)$$

8

10 where $j = \sqrt{-1}$, $e^{j\theta} = \cos(\theta) + j \sin(\theta)$, and $p(\boldsymbol{\omega})$, the **spectral density**, is a distribution over
11 frequencies.

12 We can easily derive and gain intuitions into several kernels from spectral densities. If we take the
13 Fourier transform of an RBF kernel we find the spectral density $p(\boldsymbol{\omega}) = \sqrt{2\pi\ell^2} \exp(-2\pi^2\boldsymbol{\omega}^2\ell^2)$. Thus
14 the spectral density is also Gaussian, but with a bandwidth *inversely* proportional to the length-scale
15 hyperparameter ℓ . That is, as ℓ becomes large, the spectral density collapses onto a point mass.
16 This result is intuitive: as we increase the length-scale our model treats points as correlated over
17 large distances, and becomes very smooth and slowly varying, and thus low-frequency. In general,
18 since the Gaussian distribution has relatively light tails, we can see that RBF kernels won't generally
19 support high frequency solutions.

20 We can instead use a Student- t spectral density, which has heavy tails that will provide greater
21 support for higher frequencies. Taking the inverse Fourier transform of this spectral density, we
22 recover the Matern kernel, with degrees of freedom ν corresponding to the degrees of freedom in
23 the spectral density. Indeed, the smaller we make ν the less smooth and higher frequency are the
24 associated fits to data using a Matern kernel.

25 We can also derive **spectral mixture kernels** by modelling the spectral density as a scale-location
26 mixture of Gaussians and taking the inverse Fourier transform [WA13]. Since scale-location mixtures
27 of Gaussians are dense in the set of distributions, and can therefore approximate any spectral density,
28 this kernel can approximate any stationary kernel to arbitrary precision. The spectral mixture kernel
29 thus forms a powerful approach to kernel learning, which we discuss further in Section 18.6.5.
30

31 **18.2.3.1 Random feature kernels**

33 Although the power of kernels resides in the ability to avoid working with featurized representations
34 of the inputs, such kernelized methods can take $O(N^3)$ time, in order to invert the Gram matrix \mathbf{K} ,
35 This can make it difficult to use such methods on large scale data. Fortunately, we can approximate
36 the feature map for many kernels using a randomly chosen finite set of M basis functions, thus
37 reducing the cost to $O(NM + M^3)$.

38 We will show how to do this for shift-invariant kernels by returning to Bochner's theorem in
39 Eq. (18.41):
40

41
$$\mathcal{K}(\boldsymbol{\delta}) = \int_{\mathbb{R}^d} p(\boldsymbol{\omega}) e^{j\boldsymbol{\omega}^\top \boldsymbol{\delta}} d\boldsymbol{\omega} \quad (18.42)$$

42

44 where $j = \sqrt{-1}$, $e^{j\theta} = \cos(\theta) + j \sin(\theta)$, and the spectral density $p(\boldsymbol{\omega})$ is a distribution over frequencies.

45 In the case of a Gaussian RBF kernel, we have seen that the spectral density is a Gaussian
46 distribution. Hence we can easily compute a Monte Carlo approximation to this integral by sampling
47

random Gaussian vectors. This yields the following approximation: $\mathcal{K}(\mathbf{x}, \mathbf{x}') \approx \phi(\mathbf{x})^\top \phi(\mathbf{x}')$, where the (real-valued) feature vector is given by

$$\phi(\mathbf{x}) = \sqrt{\frac{1}{D}} [\sin(\mathbf{z}_1^\top \mathbf{x}), \dots, \sin(\mathbf{z}_D^\top \mathbf{x}), \cos(\mathbf{z}_1^\top \mathbf{x}), \dots, \cos(\mathbf{z}_D^\top \mathbf{x})] \quad (18.43)$$

$$= \sqrt{\frac{1}{D}} [\sin(\mathbf{Z}^\top \mathbf{x}), \cos(\mathbf{Z}^\top \mathbf{x})] \quad (18.44)$$

Here $\mathbf{Z} = (1/\sigma)\mathbf{G}$, and $\mathbf{G} \in \mathbb{R}^{d \times D}$ is a random Gaussian matrix, where the entries are sampled iid from $\mathcal{N}(0, 1)$. The representation in Equation (18.44) are called **random Fourier features (RFF)** [RR08] or “weighted sums of random kitchen sinks” [RR09]. (One can obtain an even better approximation by ensuring that the rows of \mathbf{Z} are random but orthogonal; this is called **orthogonal random features** [Yu+16].)

One can create similar random feature representations for other kinds of kernels. We can then use such features for supervised learning by defining $f(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{W}\varphi(\mathbf{Z}\mathbf{x}) + \mathbf{b}$, where \mathbf{Z} is a random Gaussian matrix, and the form of φ depends on the chosen kernel. This is equivalent to a one layer MLP with random input-to-hidden weights; since we only optimize the hidden-to-output weights $\boldsymbol{\theta} = (\mathbf{W}, \mathbf{b})$, the model is equivalent to a linear model with fixed random features. If we use enough random features, we can approximate the performance of a kernelized prediction model, but the computational cost is now $O(N)$ rather than $O(N^2)$.

18.3 GPs with Gaussian likelihoods

In this section, we discuss GPs for regression, using a Gaussian likelihood. In this case, all the computations can be performed in closed form, using standard linear algebra methods. We extend this framework to non-Gaussian likelihoods later in the chapter.

18.3.1 Predictions using noise-free observations

Suppose we observe a training set $\mathcal{D} = \{(\mathbf{x}_n, y_n) : n = 1 : N\}$, where $y_n = f(\mathbf{x}_n)$ is the noise-free observation of the function evaluated at \mathbf{x}_n . If we ask the GP to predict $f(\mathbf{x})$ for a value of \mathbf{x} that it has already seen, we want the GP to return the answer $f(\mathbf{x})$ with no uncertainty. In other words, it should act as an **interpolator** of the training data. Here we assume the observed function values are noiseless. We will consider the case of noisy observations shortly.

Now we consider the case of predicting the outputs for new inputs that may not be in \mathcal{D} . Specifically, given a test set \mathbf{X}_* of size $N_* \times D$, we want to predict the function outputs $\mathbf{f}_* = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_{N_*})]$. By definition of the GP, the joint distribution $p(\mathbf{f}_X, \mathbf{f}_* | \mathbf{X}, \mathbf{X}_*)$ has the following form

$$\begin{pmatrix} \mathbf{f}_X \\ \mathbf{f}_* \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \boldsymbol{\mu}_X \\ \boldsymbol{\mu}_* \end{pmatrix}, \begin{pmatrix} \mathbf{K}_{X,X} & \mathbf{K}_{X,*} \\ \mathbf{K}_{X,*}^\top & \mathbf{K}_{*,*} \end{pmatrix} \right) \quad (18.45)$$

where $\boldsymbol{\mu}_X = (m(\mathbf{x}_1), \dots, m(\mathbf{x}_{N_D}))$, $\boldsymbol{\mu}_* = (m(\mathbf{x}_1^*), \dots, m(\mathbf{x}_{N_*}^*))$, $\mathbf{K}_{X,X} = \mathcal{K}(\mathbf{X}, \mathbf{X})$ is $N_D \times N_D$, $\mathbf{K}_{X,*} = \mathcal{K}(\mathbf{X}, \mathbf{X}_*)$ is $N_D \times N_*$, and $\mathbf{K}_{*,*} = \mathcal{K}(\mathbf{X}_*, \mathbf{X}_*)$ is $N_* \times N_*$. See Figure 18.7 for a static illustration, and <http://www.infinitecuriosity.org/vizgp/> for an interactive visualization.

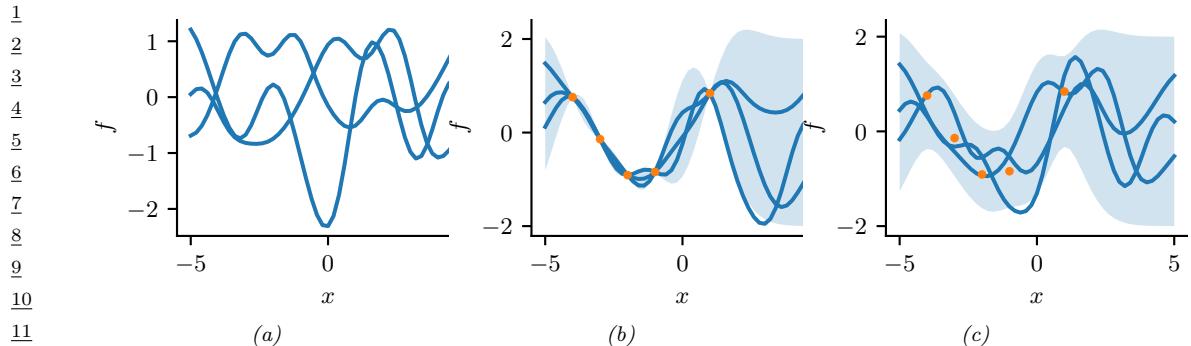


Figure 18.7: Left: some functions sampled from a GP prior with RBF kernel. Middle: some samples from a GP posterior, after conditioning on 5 noise-free observations. Right: some samples from a GP posterior, after conditioning on 5 noisy observations. The shaded area represents $\mathbb{E}[f(\mathbf{x})] \pm 2\sqrt{\text{Var}[f(\mathbf{x})]}$. Adapted from Figure 2.2 of [RW06]. Generated by [gpr_demo_noise_free.ipynb](#).

By the standard rules for conditioning Gaussians (Section 2.2.5.4), the posterior has the following form

$$p(\mathbf{f}_* | \mathbf{X}_*, \mathcal{D}) = \mathcal{N}(\mathbf{f}_* | \boldsymbol{\mu}_{*|X}, \boldsymbol{\Sigma}_{*|X}) \quad (18.46)$$

$$\boldsymbol{\mu}_{*|X} = \boldsymbol{\mu}_* + \mathbf{K}_{X,*}^\top \mathbf{K}_{X,X}^{-1} (\mathbf{f}_X - \boldsymbol{\mu}_X) \quad (18.47)$$

$$\boldsymbol{\Sigma}_{*|X} = \mathbf{K}_{*,*} - \mathbf{K}_{X,*}^\top \mathbf{K}_{X,X}^{-1} \mathbf{K}_{X,*} \quad (18.48)$$

This process is illustrated in Figure 18.7. On the left we show some samples from the prior, $p(f)$, where we use an RBF kernel (Section 18.2.1.1) and a zero mean function. On the right, we show samples from the posterior, $p(f|\mathcal{D})$. We see that the model perfectly interpolates the training data, and that the predictive uncertainty increases as we move further away from the observed data.

Note that the cost of the above method for sampling N_* points is $O(N_*^3)$. This can be reduced to $O(N_*)$ time using the methods in [Ple+18; Wil+20a].

18.3.2 Predictions using noisy observations

In Section 18.3.1, we showed how to do GP regression when the training data was noiseless. Now let us consider the case where what we observe is a noisy version of the underlying function, $y_n = f(\mathbf{x}_n) + \epsilon_n$, where $\epsilon_n \sim \mathcal{N}(0, \sigma_y^2)$. In this case, the model is not required to interpolate the data, but it must come “close” to the observed data. The covariance of the observed noisy responses is

$$\text{Cov}[y_i, y_j] = \text{Cov}[f_i, f_j] + \text{Cov}[\epsilon_i, \epsilon_j] = \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) + \sigma_y^2 \delta_{ij} \quad (18.49)$$

where $\delta_{ij} = \mathbb{I}(i = j)$. In other words

$$\text{Cov}[\mathbf{y} | \mathbf{X}] = \mathbf{K}_{X,X} + \sigma_y^2 \mathbf{I}_N \quad (18.50)$$

The joint density of the observed data and the latent, noise-free function on the test points is given

1
2 by

3
4
$$\begin{pmatrix} \mathbf{y} \\ \mathbf{f}_* \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \boldsymbol{\mu}_X \\ \boldsymbol{\mu}_* \end{pmatrix}, \begin{pmatrix} \mathbf{K}_{X,X} + \sigma_y^2 \mathbf{I} & \mathbf{K}_{X,*} \\ \mathbf{K}_{X,*}^\top & \mathbf{K}_{*,*} \end{pmatrix} \right) \quad (18.51)$$

5

6 Hence the posterior predictive density at a set of test points \mathbf{X}_* is

7
8
$$p(\mathbf{f}_* | \mathcal{D}, \mathbf{X}_*) = \mathcal{N}(\mathbf{f}_* | \boldsymbol{\mu}_{*|X}, \boldsymbol{\Sigma}_{*|X}) \quad (18.52)$$

9
10
$$\boldsymbol{\mu}_{*|X} = \boldsymbol{\mu}_* + \mathbf{K}_{X,*}^\top (\mathbf{K}_{X,X} + \sigma_y^2 \mathbf{I})^{-1} (\mathbf{y} - \boldsymbol{\mu}_X) \quad (18.53)$$

11
12
$$\boldsymbol{\Sigma}_{*|X} = \mathbf{K}_{*,*} - \mathbf{K}_{X,*}^\top (\mathbf{K}_{X,X} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{K}_{X,*} \quad (18.54)$$

13 In the case of a single test input, this simplifies as follows

14
15
$$p(f_* | \mathcal{D}, \mathbf{x}_*) = \mathcal{N}(f_* | m_* + \mathbf{k}_*^\top (\mathbf{K}_{X,X} + \sigma_y^2 \mathbf{I})^{-1} (\mathbf{y} - \boldsymbol{\mu}_X), k_{**} - \mathbf{k}_*^\top (\mathbf{K}_{X,X} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{k}_*) \quad (18.55)$$

16 where $\mathbf{k}_* = [\mathcal{K}(\mathbf{x}_*, \mathbf{x}_1), \dots, \mathcal{K}(\mathbf{x}_*, \mathbf{x}_N)]$ and $k_{**} = \mathcal{K}(\mathbf{x}_*, \mathbf{x}_*)$. If the mean function is zero, we can
17 write the posterior mean as follows:

18
19
$$\boldsymbol{\mu}_{*|X} = \mathbf{k}_*^\top \underbrace{\mathbf{K}_\sigma^{-1} \mathbf{y}}_{\boldsymbol{\alpha}} = \sum_{n=1}^N \mathcal{K}(\mathbf{x}_*, \mathbf{x}_n) \alpha_n \quad (18.56)$$

20

21 where

22
23
$$\mathbf{K}_\sigma = \mathbf{K}_{X,X} + \sigma_y^2 \mathbf{I} \quad (18.57)$$

24
25
$$\boldsymbol{\alpha} = \mathbf{K}_\sigma^{-1} \mathbf{y} \quad (18.58)$$

26 Fitting this model amounts to computing $\boldsymbol{\alpha}$ in Equation (18.58). This is usually done by computing
27 the Cholesky decomposition of \mathbf{K}_σ , as described in Section 18.3.6. Once we have computed $\boldsymbol{\alpha}$, we
28 can compute predictions for each test point in $O(N)$ time for the mean, and $O(N^2)$ time for the
29 variance.

31 18.3.3 Weight space vs function space

32 In this section, we show how Bayesian linear regression is a special case of a GP.

33 Consider the linear regression model $y = f(\mathbf{x}) + \epsilon$, where $f(\mathbf{x}) = \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x})$ and $\epsilon \sim \mathcal{N}(0, \sigma_y^2)$. If we
34 use a Gaussian prior $p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \boldsymbol{\Sigma}_w)$, then the posterior is as follows (see Section 15.2.1 for the
35 derivation):

36
37
$$p(\mathbf{w} | \mathcal{D}) = \mathcal{N}(\mathbf{w} | \frac{1}{\sigma_y^2} \mathbf{A}^{-1} \boldsymbol{\Phi}^\top \mathbf{y}, \mathbf{A}^{-1}) \quad (18.59)$$

38

39 where $\boldsymbol{\Phi}$ is the $N \times D$ design matrix, and

40
41
$$\mathbf{A} = \sigma_y^{-2} \boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \boldsymbol{\Sigma}_w^{-1} \quad (18.60)$$

42 The posterior predictive distribution for $f_* = f(\mathbf{x}_*)$ is therefore

43
44
$$p(f_* | \mathcal{D}, \mathbf{x}_*) = \mathcal{N}(f_* | \frac{1}{\sigma_y^2} \boldsymbol{\phi}_*^\top \mathbf{A}^{-1} \boldsymbol{\Phi}^\top \mathbf{y}, \boldsymbol{\phi}_*^\top \mathbf{A}^{-1} \boldsymbol{\phi}_*) \quad (18.61)$$

45

where $\phi_* = \phi(\mathbf{x}_*)$. This views the problem of inference and prediction in **weight space**.

We now show that this is equivalent to the predictions made by a GP using a kernel of the form $\mathcal{K}(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^\top \Sigma_w \phi(\mathbf{x}')$. To see this, let $\mathbf{K} = \Phi \Sigma_w \Phi^\top$, $\mathbf{k}_* = \Phi \Sigma_w \phi_*$, and $k_{**} = \phi_*^\top \Sigma_w \phi_*$. Using this notation, and the matrix inversion lemma, we can rewrite Equation (18.61) as follows

$$p(f_* | \mathcal{D}, \mathbf{x}_*) = \mathcal{N}(f_* | \boldsymbol{\mu}_{*|X}, \boldsymbol{\Sigma}_{*|X}) \quad (18.62)$$

$$\boldsymbol{\mu}_{*|X} = \phi_*^\top \Sigma_w \Phi^\top (\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y} = \mathbf{k}_*^\top (\mathbf{K}_{X,X} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y} \quad (18.63)$$

$$\boldsymbol{\Sigma}_{*|X} = \phi_*^\top \Sigma_w \phi_* - \phi_*^\top \Sigma_w \Phi^\top (\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \Phi \Sigma_w \phi_* = k_{**} - \mathbf{k}_*^\top (\mathbf{K}_{X,X} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{k}_* \quad (18.64)$$

which matches the results in Equation (18.55), assuming $m(\mathbf{x}) = 0$. A non-zero mean can be captured by adding a constant feature with value 1 to $\phi(\mathbf{x})$.

Thus we can derive a GP from Bayesian linear regression. Note, however, that linear regression assumes $\phi(\mathbf{x})$ is a finite length vector, whereas a GP allows us to work directly in terms of kernels, which may correspond to infinite length feature vectors (see Section 18.2.2). That is, a GP works in **function space**.

18.3.4 Semi-parametric GPs

So far, we have mostly assumed the mean of the GP is 0, and have relied on its interpolation abilities to model the mean function. Sometimes it is useful to fit a global linear model for the mean, and use the GP to model the residual errors, as follows:

$$g(\mathbf{x}) = f(\mathbf{x}) + \boldsymbol{\beta}^\top \phi(\mathbf{x}) \quad (18.65)$$

where $f(\mathbf{x}) \sim \text{GP}(0, \mathcal{K}(\mathbf{x}, \mathbf{x}'))$, and $\phi()$ are some fixed basis functions. This combines a parametric and a non-parametric model, and is known as a **semi-parametric model**.

If we assume $\boldsymbol{\beta} \sim \mathcal{N}(\mathbf{b}, \mathbf{B})$, we can integrate these parameters out to get a new GP [O'H78]:

$$g(\mathbf{x}) \sim \text{GP}(\phi(\mathbf{x})^\top \mathbf{b}, \mathcal{K}(\mathbf{x}, \mathbf{x}') + \phi(\mathbf{x})^\top \mathbf{B} \phi(\mathbf{x}')) \quad (18.66)$$

Let $\mathbf{H}_X = \phi(\mathbf{X})^\top$ be the $D \times N$ matrix of training examples, and $\mathbf{H}_* = \phi(\mathbf{X}_*)^\top$ be the $D \times N_*$ matrix of test examples. The corresponding predictive distribution for test inputs \mathbf{X}_* has the following form [RW06, p28]:

$$\mathbb{E}[g(\mathbf{X}_*) | \mathcal{D}] = \mathbf{H}_*^\top \bar{\boldsymbol{\beta}} + \mathbf{K}_{X,*}^\top \mathbf{K}_\sigma^{-1} (\mathbf{y} - \mathbf{H}_X^\top \bar{\boldsymbol{\beta}}) = \mathbb{E}[f(\mathbf{X}_*) | \mathcal{D}] + \mathbf{R}^\top \bar{\boldsymbol{\beta}} \quad (18.67)$$

$$\text{Cov}[g(\mathbf{X}_*) | \mathcal{D}] = \text{Cov}[f(\mathbf{X}_*) | \mathcal{D}] + \mathbf{R}^\top (\mathbf{B}^{-1} + \mathbf{H}_X \mathbf{K}_\sigma^{-1} \mathbf{H}_X^\top)^{-1} \mathbf{R} \quad (18.68)$$

$$\bar{\boldsymbol{\beta}} = (\mathbf{B}^{-1} + \mathbf{H}_X \mathbf{K}_\sigma^{-1} \mathbf{H}_X^\top)^{-1} (\mathbf{H}_X \mathbf{K}_\sigma^{-1} \mathbf{y} + \mathbf{B}^{-1} \mathbf{b}) \quad (18.69)$$

$$\mathbf{R} = \mathbf{H}_* - \mathbf{H}_X \mathbf{K}_\sigma^{-1} \mathbf{K}_{X,*} \quad (18.70)$$

These results can be interpreted as follows: the mean is the usual mean from the GP, plus a global offset from the linear model, using $\bar{\boldsymbol{\beta}}$; and the covariance is the usual covariance from the GP, plus an additional positive term due to the uncertainty in $\boldsymbol{\beta}$.

In the limit of an uninformative prior for the regression parameters, as $\mathbf{B} \rightarrow \infty \mathbf{I}$, this simplifies to

$$\mathbb{E}[g(\mathbf{X}_*) | \mathcal{D}] = \mathbb{E}[f(\mathbf{X}_*) | \mathcal{D}] + \mathbf{R}^\top (\mathbf{H}_X \mathbf{K}_\sigma^{-1} \mathbf{H}_X^\top)^{-1} \mathbf{H}_X \mathbf{K}_\sigma^{-1} \mathbf{y} \quad (18.71)$$

$$\text{Cov}[g(\mathbf{X}_*) | \mathcal{D}] = \text{Cov}[f(\mathbf{X}_*) | \mathcal{D}] + \mathbf{R}^\top (\mathbf{H}_X \mathbf{K}_\sigma^{-1} \mathbf{H}_X^\top)^{-1} \mathbf{R} \quad (18.72)$$

18.3.5 Marginal likelihood

Most kernels have some free parameters. For example, the RBF-ARD kernel (Section 18.2.1.1) has the form

$$\mathcal{K}(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2}\sum_{d=1}^D \frac{1}{\ell_d^2}(x_d - x'_d)^2\right) = \prod_{d=1}^D \mathcal{K}_{\ell_d}(x_d, x'_d) \quad (18.73)$$

where each ℓ_d is a length scale for feature dimension d . Let these (and the observation noise variance σ_y^2 , if present) be denoted by $\boldsymbol{\theta}$. We can compute the likelihood of these parameters as follows:

$$p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = p(\mathcal{D}|\boldsymbol{\theta}) = \int p(\mathbf{y}|\mathbf{f}_X, \boldsymbol{\theta})p(\mathbf{f}_X|\mathbf{X}, \boldsymbol{\theta})d\mathbf{f}_X \quad (18.74)$$

Since we are integrating out the function f , we often call $\boldsymbol{\theta}$ hyperparameters, and the quantity $p(\mathcal{D}|\boldsymbol{\theta})$ the marginal likelihood.

Since f is a GP, we can compute the above integral using the marginal likelihood for the corresponding Gaussian. This gives

$$\log p(\mathcal{D}|\boldsymbol{\theta}) = -\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu}_X)^\top \mathbf{K}_\sigma^{-1}(\mathbf{y} - \boldsymbol{\mu}_X) - \frac{1}{2} \log |\mathbf{K}_\sigma| - \frac{N}{2} \log(2\pi) \quad (18.75)$$

The first term term is the square of the Mahalanobis distance between the observations and the predicted values: better fits will have smaller distance. The second term is the log determinant of the covariance matrix, which measures model complexity: smoother functions will have smaller determinants, so $-\log |\mathbf{K}_\sigma|$ will be larger (less negative) for simpler functions. The marginal likelihood measures the tradeoff between fit and complexity.

In Section 18.6.1, we discuss how to learn the kernel parameters from data by maximizing the marginal likelihood wrt $\boldsymbol{\theta}$.

18.3.6 Computational and numerical issues

In this section, we discuss computational and numerical issues which arise when implementing the above equations. For notational simplicity, we assume the prior mean is zero, $m(\mathbf{x}) = 0$.

The posterior predictive mean is given by $\mu_* = \mathbf{k}_*^\top \mathbf{K}_\sigma^{-1} \mathbf{y}$. For reasons of numerical stability, it is unwise to directly invert \mathbf{K}_σ . A more robust alternative is to compute a Cholesky decomposition, $\mathbf{K}_\sigma = \mathbf{L}\mathbf{L}^\top$, which takes $O(N^3)$ time. Given this, we can compute

$$\mu_* = \mathbf{k}_*^\top \mathbf{K}_\sigma^{-1} \mathbf{y} = \mathbf{k}_*^\top \mathbf{L}^{-\top} (\mathbf{L}^{-1} \mathbf{y}) = \mathbf{k}_*^\top \boldsymbol{\alpha} \quad (18.76)$$

Here $\boldsymbol{\alpha} = \mathbf{L}^\top \backslash (\mathbf{L} \backslash \mathbf{y})$, where we have used the backslash operator to represent backsubstitution.

We can compute the variance in $O(N^2)$ time for each test case using

$$\sigma_*^2 = k_{**} - \mathbf{k}_*^\top \mathbf{L}^{-T} \mathbf{L}^{-1} \mathbf{k}_* = k_{**} - \mathbf{v}^\top \mathbf{v} \quad (18.77)$$

where $\mathbf{v} = \mathbf{L} \backslash \mathbf{k}_*$.

1 Finally, the log marginal likelihood (needed for kernel learning, Section 18.6) can be computed
2 using
3

4

$$\log p(\mathbf{y}|\mathbf{X}) = -\frac{1}{2}\mathbf{y}^\top \boldsymbol{\alpha} - \sum_{n=1}^N \log L_{nn} - \frac{N}{2} \log(2\pi) \quad (18.78)$$

5

6 We see that overall cost is dominated by $O(N^3)$. We discuss faster, but approximate, methods in
7 Section 18.5.
8

11 18.3.7 Kernel ridge regression

12

13 The term **ridge regression** refers to linear regression with an ℓ_2 penalty on the regression weights:
14

15

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{n=1}^N (y_n - f(\mathbf{x}_n; \mathbf{w}))^2 + \lambda \|\mathbf{w}\|_2^2 \quad (18.79)$$

16

17 where $f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^\top \mathbf{x}$. The solution for this is
18

19

$$\mathbf{w}^* = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y} = \left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top + \lambda \mathbf{I} \right)^{-1} \left(\sum_{n=1}^N \mathbf{x}_n y_n \right) \quad (18.80)$$

20

21 In this section, we consider a function space version of this:
22

23

$$f^* = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \sum_{n=1}^N (y_n - f(\mathbf{x}_n))^2 + \lambda \|f\|^2 \quad (18.81)$$

24

25 For this to make sense, we have to define the function space \mathcal{F} and the norm $\|f\|$. If we use a
26 function space derived from a positive definite kernel function \mathcal{K} , the resulting method is called
27 **kernel ridge regression** (KRR). We will see that the resulting estimate $f^*(\mathbf{x}_*)$ is equivalent to
28 the posterior mean of a GP. We give the details below.

34 18.3.7.1 Reproducing kernel Hilbert spaces

35

36 In this section, we briefly introduce the relevant mathematical “machinery” needed to explain KRR.
37

38 Let $\mathcal{F} = \{f : \mathcal{X} \rightarrow \mathbb{R}\}$ be a space of real-valued functions. Elements of this space (i.e., functions)
39 can be added and scalar multiplied as if they were vectors. That is, if $f \in \mathcal{F}$ and $g \in \mathcal{F}$, then
40 $\alpha f + \beta g \in \mathcal{F}$ for $\alpha, \beta \in \mathbb{R}$. We can also define an **inner product** for \mathcal{F} , which is a mapping $\langle f, g \rangle \in \mathbb{R}$
41 which satisfies the following:

42

$$\langle \alpha f_1 + \beta f_2, g \rangle = \alpha \langle f_1, g \rangle + \beta \langle f_2, g \rangle \quad (18.82)$$

43

$$\langle f, g \rangle = \langle g, f \rangle \quad (18.83)$$

$$\langle f, f \rangle \geq 0 \quad (18.84)$$

$$\langle f, f \rangle = 0 \text{ iff } f(x) = 0 \text{ for all } x \in \mathcal{X} \quad (18.85)$$

44

1 We define the norm of a function using
2

$$\underline{3} \quad \|f\| \triangleq \sqrt{\langle f, f \rangle} \quad (18.86)$$

5 A function space \mathcal{H} with an inner product operator is called a **Hilbert space**. (We also require
6 that the function space be complete, which means that every Cauchy sequence of functions $f_i \in \mathcal{H}$
7 has a limit that is also in \mathcal{H} .)

8 The most common Hilbert space is the space known as L^2 . To define this, we need to specify a
9 **measure** μ on the input space \mathcal{X} ; this is a function that assigns any (suitable) subset A of \mathcal{X} to a
10 positive number, such as its volume. This can be defined in terms of the density function $w : \mathcal{X} \rightarrow \mathbb{R}$,
11 as follows:

$$\underline{12} \quad \mu(A) = \int_A w(x) dx \quad (18.87)$$

15 Thus we have $\mu(dx) = w(x)dx$. We can now define $L^2(\mathcal{X}, \mu)$ to be the space of functions $f : \mathcal{X} \rightarrow \mathbb{R}$
16 that satisfy
17

$$\underline{18} \quad \int_{\mathcal{X}} f(x)^2 w(x) dx < \infty \quad (18.88)$$

20 This is known as the set of **square-integrable functions**. This space has an inner product defined
21 by
22

$$\underline{23} \quad \langle f, g \rangle = \int_{\mathcal{X}} f(x)g(x)w(x) dx \quad (18.89)$$

26 We define a **Reproducing Kernel Hilbert Space** or **RKHS** as follows. Let \mathcal{H} be a Hilbert
27 space of functions $f : \mathcal{X} \rightarrow \mathbb{R}$. We say that \mathcal{H} is an RKHS endowed with inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ if there
28 exists a (symmetric) **kernel function** $\mathcal{K} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ with the following properties:

- 29 • For every $\mathbf{x} \in \mathcal{X}$, $\mathcal{K}(\mathbf{x}, \cdot) \in \mathcal{H}$.
- 30 • \mathcal{K} satisfies the **reproducing property**:

$$\underline{32} \quad \langle f(\cdot), \mathcal{K}(\cdot, \mathbf{x}') \rangle = f(\mathbf{x}') \quad (18.90)$$

33 The reason for the term “reproducing property” is as follows. Let $f(\cdot) = \mathcal{K}(\mathbf{x}, \cdot)$. Then we have that
34

$$\underline{35} \quad \langle \mathcal{K}(\mathbf{x}, \cdot), \mathcal{K}(\cdot, \mathbf{x}') \rangle = \mathcal{K}(\mathbf{x}, \mathbf{x}') \quad (18.91)$$

37 18.3.7.2 Complexity of a function in an RKHS

39 The main utility of RKHS from the point of view of machine learning is that it allows us to define a
40 notion of a function’s “smoothness” or “complexity” in terms of its norm, as we now discuss.

41 Suppose we have a positive definite kernel function \mathcal{K} . From Mercer’s theorem we have $\mathcal{K}(\mathbf{x}, \mathbf{x}') =$
42 $\sum_{i=1}^{\infty} \lambda_i \phi_i(\mathbf{x}) \phi_i(\mathbf{x}')$. Now consider a Hilbert space \mathcal{H} defined by functions of the form $f(\mathbf{x}) =$
43 $\sum_{i=1}^{\infty} f_i \phi_i(\mathbf{x})$, with $\sum_{i=1}^{\infty} f_i^2 / \lambda_i < \infty$. The inner product of two functions in this space is

$$\underline{45} \quad \langle f, g \rangle_{\mathcal{H}} = \sum_{i=1}^{\infty} \frac{f_i g_i}{\lambda_i} \quad (18.92)$$

1 Hence the (squared) norm is given by
2

3
4 $\|f\|_{\mathcal{H}}^2 = \langle f, f \rangle_{\mathcal{H}} = \sum_{i=1}^{\infty} \frac{f_i^2}{\lambda_i}$ (18.93)
5
6

7 This is analogous to the quadratic form $\mathbf{f}^T \mathbf{K}^{-1} \mathbf{f}$ which occurs in some GP objectives (see Equa-
8 tion (18.102)). Thus the smoothness of the function is controlled by the properties of the corresponding
9 kernel.

10

11 18.3.7.3 Representer theorem

12 In this section, we consider the problem of (regularized) empirical risk minimization in function space.
13 In particular, consider the following problem:

14
15 $f^* = \underset{f \in \mathcal{H}_{\mathcal{K}}}{\operatorname{argmin}} \sum_{n=1}^N \ell(y_n, f(\mathbf{x}_n)) + \frac{\lambda}{2} \|f\|_{\mathcal{H}}^2$ (18.94)
16
17

18 where $\mathcal{H}_{\mathcal{K}}$ is an RKHS with kernel \mathcal{K} and $\ell(y, \hat{y}) \in \mathbb{R}$ is a loss function. Then one can show [KW70;
19 SHS01] the following result:

20
21
22 $f^*(x) = \sum_{n=1}^N \alpha_n \mathcal{K}(x, x_n)$ (18.95)
23
24

25 where $\alpha_n \in \mathbb{R}$ are some coefficients that depend on the training data. This is called the **representer**
26 **theorem**.

27 Now consider the special case where the loss function is squared loss, and $\lambda = \sigma_y^2$. We want to
28 minimize

29
30 $\mathcal{L}(f) = \frac{1}{2\sigma_y^2} \sum_{n=1}^N (y_n - f(\mathbf{x}_n))^2 + \frac{1}{2} \|f\|_{\mathcal{H}}^2$ (18.96)
31
32

33 Substituting in Equation (18.95), and using the fact that $\langle \mathcal{K}(\cdot, \mathbf{x}_i), \mathcal{K}(\cdot, \mathbf{x}_j) \rangle = \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)$, we obtain
34

35
36 $\mathcal{L}(f) = \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\alpha} + \frac{1}{2\sigma_y^2} \|\mathbf{y} - \mathbf{K} \boldsymbol{\alpha}\|^2$ (18.97)
37

38 $= \frac{1}{2} \boldsymbol{\alpha}^T (\mathbf{K} + \frac{1}{\sigma_y^2} \mathbf{K}^2) \boldsymbol{\alpha} - \frac{1}{\sigma_y^2} \mathbf{y}^T \mathbf{K} \boldsymbol{\alpha} + \frac{1}{2\sigma_y^2} \mathbf{y}^T \mathbf{y}$ (18.98)
39

40 Minimizing this wrt $\boldsymbol{\alpha}$ gives $\hat{\boldsymbol{\alpha}} = (\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y}$, which is the same as Equation (18.58). Furthermore,
41 the prediction for a test point is

42
43 $\hat{f}(\mathbf{x}_*) = \mathbf{k}_*^T \boldsymbol{\alpha} = \mathbf{k}_*^T (\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y}$ (18.99)
44

45 This is known as **kernel ridge regression** [Vov13]. We see that the result matches the posterior
46 predictive mean of a GP in Equation (18.56).

47

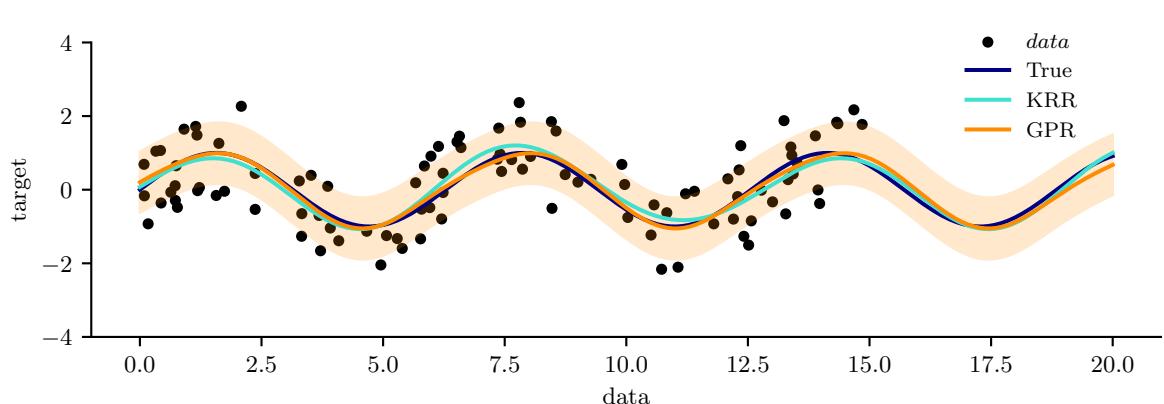


Figure 18.8: Kernel ridge regression (KRR) compared to Gaussian process regression (GPR) using the same kernel. Generated by [krr_vs_gpr.ipynb](#).

18.3.7.4 Example of KRR vs GPR

In this section, we compare KRR with GP regression on a simple 1d problem. Since the underlying function is believed to be periodic, we use the periodic kernel from Equation (18.19). To capture the fact that the observations are noisy, we add to this a **white noise kernel**

$$\mathcal{K}(\mathbf{x}, \mathbf{x}') = \sigma_y^2 \delta(\mathbf{x} - \mathbf{x}') \quad (18.100)$$

as in Equation (18.49). Thus there are 3 GP hyper-parameters: the kernel length scale ℓ , the kernel periodicity p , and the noise level σ_y^2 . We can optimize these by maximizing the marginal likelihood using gradient descent (see Section 18.6.1). For KRR, we also have 3 hyper-parameters (ℓ , p and $\lambda = \sigma_y^2$); we optimize these using grid search combined with cross validation (which in general is slower than gradient based optimization). The resulting model fits are shown in Figure 18.8, and are very similar, as is to be expected.

18.4 GPs with non-Gaussian likelihoods

So far, we have focused on GPs for regression using Gaussian likelihoods. In this case, the posterior is also a GP, and all computation can be performed analytically. However, if the likelihood is non-Gaussian, we can no longer compute the posterior exactly. We can create variety of different “classical” models by changing the form of the likelihood, as we show in Table 18.1. In the sections below, we briefly discuss some approximate inference methods.

18.4.1 Binary classification

In this section, we consider binary classification using GPs. If we use the sigmoid link function, we have $p(y_n = 1|\mathbf{x}_n) = \sigma(y_n f(\mathbf{x}_n))$. If we assume $y_n \in \{-1, +1\}$, then we have $p(y_n|\mathbf{x}_n) = \sigma(y_n f_n)$, since $\sigma(-z) = 1 - \sigma(z)$. If we use the probit link, we have $p(y_n = 1|\mathbf{x}_n) = \Phi(y_n f(\mathbf{x}_n))$, where $\Phi(z)$ is the cdf of the standard normal. More generally, let $p(y_n|\mathbf{x}_n) = \text{Ber}(y_n|\varphi(f_n))$. The overall log joint

Model	Likelihood	Section
Regression	$\mathcal{N}(f_i, \sigma_y^2)$	Section 18.3.2
Robust regression	$\mathcal{T}_\nu(f_i, \sigma_y^2)$	Section 18.4.4
Binary classification	$\text{Ber}(\sigma(f_i))$	Section 18.4.1
Multiclass classification	$\text{Cat}(\text{softmax}(f_i))$	Section 18.4.2
Poisson regression	$\text{Poi}(\exp(f_i))$	Section 18.4.3

Table 18.1: Summary of GP models with a variety of likelihoods.

$\log p(y_i f_i)$	$\frac{\partial}{\partial f_i} \log p(y_i f_i)$	$\frac{\partial^2}{\partial f_i^2} \log p(y_i f_i)$
$\log \sigma(y_i f_i)$	$t_i - \pi_i$	$-\pi_i(1 - \pi_i)$
$\log \Phi(y_i f_i)$	$\frac{y_i \phi(f_i)}{\Phi(y_i f_i)}$	$-\frac{\phi_i^2}{\Phi(y_i f_i)^2} - \frac{y_i f_i \phi(f_i)}{\Phi(y_i f_i)}$

Table 18.2: Likelihood, gradient and Hessian for binary logistic/ probit GP regression. We assume $y_i \in \{-1, +1\}$ and define $t_i = (y_i + 1)/2 \in \{0, 1\}$ and $\pi_i = \sigma(f_i)$ for logistic regression, and $\pi_i = \Phi(f_i)$ for probit regression. Also, ϕ and Φ are the pdf and cdf of $\mathcal{N}(0, 1)$. From [RW06, p43].

has the form

$$\mathcal{L}(\mathbf{f}_X) = \log p(\mathbf{y}|\mathbf{f}_X) + \log p(\mathbf{f}_X|\mathbf{X}) \quad (18.101)$$

$$= \log p(\mathbf{y}|\mathbf{f}_X) - \frac{1}{2} \mathbf{f}_X^\top \mathbf{K}_{X,X}^{-1} \mathbf{f}_X - \frac{1}{2} \log |\mathbf{K}_{X,X}| - \frac{N}{2} \log 2\pi \quad (18.102)$$

The simplest approach to approximate inference is to use a Laplace approximation (Section 7.4.3). The gradient and Hessian of the log joint are given by

$$\nabla \mathcal{L} = \nabla \log p(\mathbf{y}|\mathbf{f}_X) - \mathbf{K}_{X,X}^{-1} \mathbf{f}_X \quad (18.103)$$

$$\nabla^2 \mathcal{L} = \nabla^2 \log p(\mathbf{y}|\mathbf{f}_X) - \mathbf{K}_{X,X}^{-1} = -\mathbf{\Lambda} - \mathbf{K}_{X,X}^{-1} \quad (18.104)$$

where $\mathbf{\Lambda} \triangleq -\nabla^2 \log p(\mathbf{y}|\mathbf{f}_X)$ is a diagonal matrix, since the likelihood factorizes across examples. Expressions for the gradient and Hessian of the log likelihood for the logit and probit case are shown in Table 18.2. At convergence, the Laplace approximation of the posterior takes the following form:

$$p(\mathbf{f}_X|\mathcal{D}) \approx q(\mathbf{f}_X) = \mathcal{N}(\hat{\mathbf{f}}, (\mathbf{K}_{X,X}^{-1} + \mathbf{\Lambda})^{-1}) \quad (18.105)$$

where $\hat{\mathbf{f}}$ is the MAP estimate. See [RW06, Sec 3.4] for further details.

For improved accuracy, we can use variational inference, in which we assume $q(\mathbf{f}_X) = \mathcal{N}(\mathbf{f}_X|\mathbf{m}, \mathbf{S})$; we then optimize \mathbf{m} and \mathbf{S} using (stochastic) gradient descent, rather than assuming \mathbf{S} is the Hessian at the mode. See Section 18.5.4 for the details.

Once we have a Gaussian posterior $q(\mathbf{f}_X|\mathcal{D})$, we can then use standard GP prediction to compute $q(f_*|\mathbf{x}_*, \mathcal{D})$. Finally, we can approximate the posterior predictive distribution over binary labels using

$$\pi_* = p(y_* = 1|\mathbf{x}_*, \mathcal{D}) = \int p(y_* = 1|f_*) q(f_*|\mathbf{x}_*, \mathcal{D}) df_* \quad (18.106)$$

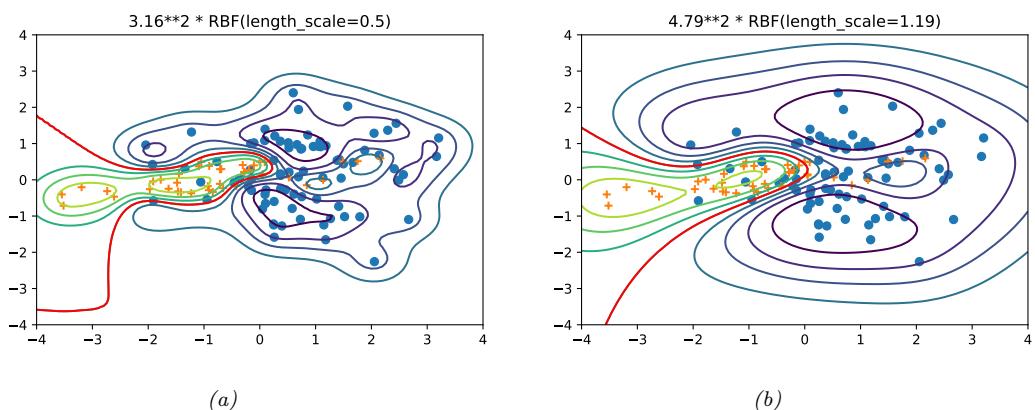


Figure 18.9: Contours of the posterior predictive probability for a binary classifier generated by a GP with an SE kernel. (a) Manual kernel parameters: short length scale, $\ell = 0.5$, variance $3.16^2 \approx 9.98$. (b) Learned kernel parameters: long length scale, $\ell = 1.19$, variance $4.79^2 \approx 22.9$. Generated by `gpc demo 2d.ipynb`.

This 1d integral can be computed using the probit approximation from Section 15.3.5. In this case we have $\pi_* \approx \sigma(\kappa(v)\mathbb{E}[f_*])$, where $v = \mathbb{V}[f_*]$ and $\kappa^2(v) = (1 + \pi v / 8)^{-1}$.

In Figure 18.9, we show a synthetic binary classification problem in 2d. We use an SE kernel. On the left, we show predictions using hyper-parameters set by hand; we use a short length scale, hence the very sharp turns in the decision boundary. On the right, we show the predictions using the learned hyper-parameters; the model favors more parsimonious explanation of the data.

18.4.2 Multi-class classification

The multi-class case is somewhat harder, since the function now needs to return a vector of C logits to get $p(y_n|\mathbf{x}_n) = \text{Cat}(y_n|\text{softmax}(\mathbf{f}_n))$, where $\mathbf{f}_n = (f_n^1, \dots, f_n^C)$. It is standard to assume that $f^c \sim \text{GP}(0, \mathcal{K}_c)$. Thus we have one latent function per class, which are a priori independent, and which may use different kernels.

We can derive a Laplace approximation for this model as discussed in [RW06, Sec 3.5]. Alternatively, we can use a variational approach, using the local variational bound to the multinomial softmax in [Cha12]. An alternative variational method, based on data augmentation with auxiliary variables, is described in [Wen+19b; Liu+19a; GFWQ20].

18.4.3 GPs for Poisson regression (Cox process)

In this section, we illustrate Poisson regression where the underlying log rate function is modeled by a GP. This is known as a **Cox process**. We can perform approximate posterior inference in this model using Laplace, MCMC or SVI (stochastic variational inference). In Figure 18.10 we give a 1d example, where we use a Matern $\frac{5}{2}$ kernel. We apply MCMC and SVI. In the VI case, we additionally have to specify the form of the posterior; we use a Gaussian approximation for the variational GP posterior $p(\mathbf{f}|\mathbf{X}, \mathbf{y})$, and a point estimate for the kernel parameters.

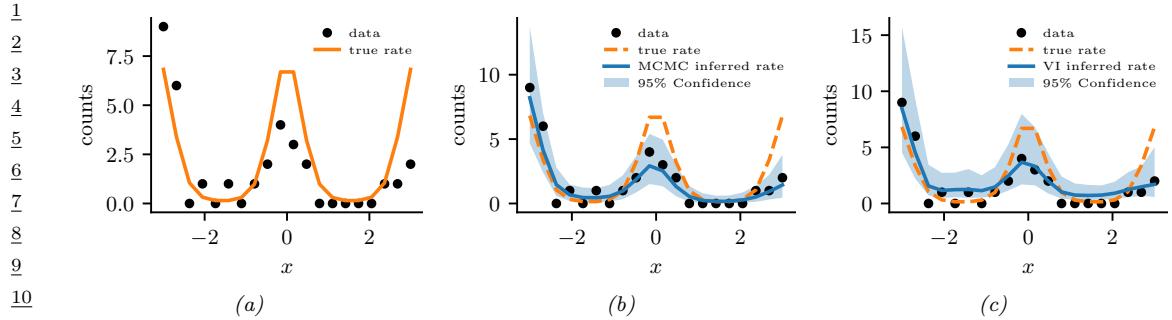


Figure 18.10: Poisson regression with a GP. (a) Observed data (black dots) and true log rate function (yellow line). (b) Posterior predictive distribution (shading shows 1 and 2 σ bands) from MCMC. (c) Posterior predictive distribution from SVI. Generated by [gp_poisson_1d.ipynb](#).

An interesting application of this is to spatial **disease mapping**. For example, [VPV10] discuss the problem of modeling the relative risk of heart attack in different regions in Finland. The data consists of the heart attacks in Finland from 1996-2000 aggregated into 20km x 20km lattice cells. The likelihood has the following form: $y_n \sim \text{Poi}(e_n r_n)$, where e_n is the known expected number of deaths (related to the population of cell n and the overall death rate), and r_n is the **relative risk** of cell n which we want to infer. Since the data counts are small, we regularize the problem by sharing information with spatial neighbors. Hence we assume $f \triangleq \log(r) \sim \text{GP}(0, K)$. We use a Matern kernel (Section 18.2.1.1) with $\nu = 3/2$, and a length scale and magnitude that are estimated from data.

Figure 18.11 gives an example of this method in action (using Laplace approximation). On the left we plot the posterior mean relative risk (RR), and on the right, the posterior variance. We see that the RR is higher in Eastern Finland, which is consistent with other studies. We also see that the variance in the North is higher, since there are fewer people living there.

18.4.4 Other likelihoods

Many other likelihoods are possible. For example, [VJV09] uses a Student- t likelihood in order to perform robust regression. A general method for performing approximate variational inference in GPs with such non-conjugate likelihoods is discussed in [WSS21].

18.5 Scaling GP inference to large datasets

In Section 18.3.6, we saw that the best way to perform GP inference and training is to compute a Cholesky decomposition of the $N \times N$ Gram matrix. Unfortunately, this takes $O(N^3)$ time. In this section, we discuss methods to scale up GPs to handle large N . See Table 18.3 for a summary, and [Liu+20c] for more details.¹

¹ 1. We focus on efficient methods for evaluating the marginal likelihood and the posterior predictive distribution. For an efficient method for sampling a function from the posterior, see [Wil+20a].

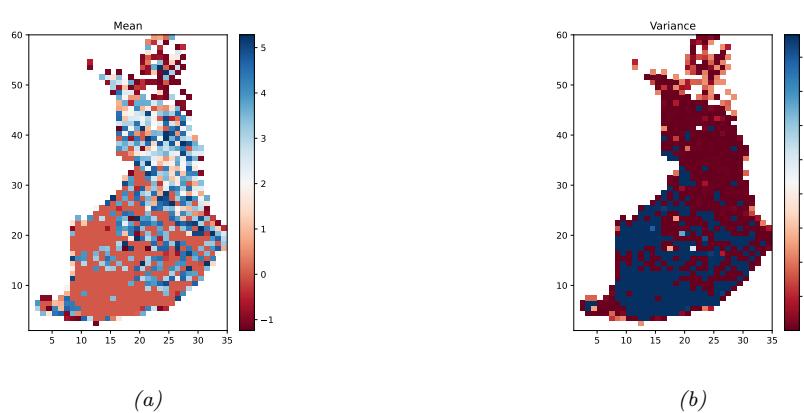


Figure 18.11: We show the relative risk of heart disease in Finland using a Poisson GP fit to 911 data points. Left: posterior mean. Right: posterior variance. Generated by `gp_spatial_demo.ipynb`.

Method	Cost	Section
Cholesky	$O(N^3)$	Section 18.3.6
Conj. Grad.	$O(CN^2)$	Section 18.5.5
Inducing	$O(NM^2 + M^3 + DNM)$	Section 18.5.3
Variational	$O(NM^2 + M^3 + DNM)$	Section 18.5.4
SVGP	$O(BM^2 + M^3 + DNM)$	Section 18.5.4.3
KISS-GP	$O(CN + CDM^D \log M)$	Section 18.5.5.3
SKIP	$O(DLN + DLM \log M + L^3 N \log D + CL^2 N)$	Section 18.5.5.3

Table 18.3: Summary of time to compute the log marginal likelihood of a GP regression model. Notation: N is number of training examples, M is number of inducing points, B is size of minibatch, D is dimensionality of input vectors (assuming $\mathcal{X} = \mathbb{R}^D$), C is number of conjugate gradient iterations. L is number of Lanczos iterations. Based on Table 2 of [Gar+18a].

18.5.1 Subset of data

The simplest approach to speeding up GP inference is to throw away some of the data. Suppose we keep a subset of M examples. In this case, exact inference will take $O(M^3)$ time. This is called the **subset-of-data** approach.

The key question is: how should we choose the subset? The simplest approach is to pick random examples (this method was recently analysed in [HIY19]). However, intuitively it makes more sense to try to pick a subset that in some sense “covers” the original data, so it contains approximately the same information (up to some tolerance) without the redundancy. Clustering algorithms are one heuristic approach, but we can also use coresets methods, which can provably find such an information-preserving subset (see e.g., [Hug+19] for an application of this idea to GPs).

1 **18.5.1.1 Informative vector machine**

3 Clustering and coresets methods are unsupervised, in that they only look at the features \mathbf{x}_i and
4 not the labels y_i , which can be suboptimal. The **informative vector machine** [HLS03] uses a
5 greedy strategy to iteratively add the labeled example (\mathbf{x}_j, y_j) that maximally reduces the entropy
6 of the function's posterior, $\Delta_j = \mathbb{H}(p(f_j)) - \mathbb{H}(p^{\text{new}}(f_j))$, where $p^{\text{new}}(f_j)$ is the posterior of f
7 at \mathbf{x}_j after conditioning on y_j . (This is very similar to active learning.) To compute Δ_j , let
8 $p(f_j) = \mathcal{N}(\mu_j, v_j)$, and $p(f_j|y_j) \propto p(f_j)\mathcal{N}(y_j|f_j, \sigma^2) = \mathcal{N}(f_j|\mu_j^{\text{new}}, v_j^{\text{new}})$, where $(v_j^{\text{new}})^{-1} = v_j^{-1} + \sigma^{-2}$.
9 Since $\mathbb{H}(\mathcal{N}(\mu, v)) = \log(2\pi ev)/2$, we have $\Delta_j = 0.5 \log(1 + v_j/\sigma^2)$. Since this is a monotonic function
10 of v_j , we can maximize it by choosing the site with the largest variance. (In fact, entropy is a
11 submodular function, so we can use submodular optimization algorithms to improve on the IVM, as
12 shown in [Kra+08].)
13

14
15 **18.5.1.2 Discussion**

16 The main problem with the subset of data approach is that it ignores some of the data, which can
17 reduce predictive accuracy and increase uncertainty about the true function. Fortunately there
18 are other scalable methods that avoid this problem, essentially by approximately representing (or
19 compressing) the training data, as we discuss below.
20

21
22 **18.5.2 Nyström approximation**

23 Suppose we had a rank M approximation to the $N \times N$ matrix gram matrix of the following form:
24

$$\mathbf{K}_{X,X} \approx \mathbf{U}\Lambda\mathbf{U}^\top \quad (18.107)$$

27 where Λ is a diagonal matrix of the M leading eigenvalues, and \mathbf{U} is the matrix of the corresponding
28 M eigenvectors, each of size N . In this case, we can use the matrix inversion lemma to write
29

$$\mathbf{K}_\sigma^{-1} = (\mathbf{K}_{X,X} + \sigma^2 \mathbf{I}_N)^{-1} \approx \sigma^{-2} \mathbf{I}_N + \sigma^{-2} \mathbf{U}(\sigma^2 \Lambda^{-1} + \mathbf{U}^\top \mathbf{U})^{-1} \mathbf{U}^\top \quad (18.108)$$

32 which takes $O(NM^2)$ time. Similarly, one can show (using the Sylvester determinant lemma) that
33

$$|\mathbf{K}_\sigma| \approx |\Lambda| |\sigma^2 \Lambda^{-1} + \mathbf{U}^\top \mathbf{U}| \quad (18.109)$$

36 which also takes $O(NM^2)$ time.

37 Unfortunately, directly computing such an eigendecomposition takes $O(N^3)$ time, which does not
38 help. However, suppose we pick a subset Z of $M < N$ points. We can partition the Gram matrix as
39 follows (where we assume the chosen points come first, and then the remaining points):
40

$$\mathbf{K}_{X,X} = \begin{pmatrix} \mathbf{K}_{Z,Z} & \mathbf{K}_{Z,X-Z} \\ \mathbf{K}_{X-Z,Z} & \mathbf{K}_{X-Z,X-Z} \end{pmatrix} \quad (18.110)$$

44 Let $\mathbf{K}_{Z,X}$ denote the top $M \times N$ block, and $\mathbf{K}_{X,Z}$ denote its transpose. We now compute an
45 eigendecomposition of $\mathbf{K}_{Z,Z}$ to get the eigenvalues $\{\lambda_i\}_{i=1}^M$ and eigenvectors $\{\mathbf{u}_i\}_{i=1}^M$. We now use
46 these to approximate the full matrix as shown below, where the scaling constants are chosen so that
47

1 $\|\tilde{\mathbf{u}}_i\| \approx 1$:

2

$$\tilde{\lambda}_i \triangleq \frac{N}{M} \lambda_i \quad (18.111)$$

3

$$\tilde{\mathbf{u}} \triangleq \sqrt{\frac{M}{N}} \frac{1}{\lambda_i} \mathbf{K}_{X,Z} \mathbf{u}_i \quad (18.112)$$

4

$$\mathbf{K}_{X,X} \approx \sum_{i=1}^M \tilde{\lambda}_i \tilde{\mathbf{u}}_i \tilde{\mathbf{u}}_i^\top \quad (18.113)$$

5

$$= \sum_{i=1}^M \frac{N}{M} \lambda_i \sqrt{\frac{M}{N}} \frac{1}{\lambda_i} \mathbf{K}_{X,Z} \mathbf{u}_i \sqrt{\frac{M}{N}} \frac{1}{\lambda_i} \mathbf{u}_i^\top \mathbf{K}_{X,Z}^\top \quad (18.114)$$

6

$$= \mathbf{K}_{X,Z} \left(\sum_{i=1}^M \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^\top \right) \mathbf{K}_{Z,X} \quad (18.115)$$

7

$$= \mathbf{K}_{X,Z} \mathbf{K}_{Z,Z}^{-1} \mathbf{K}_{Z,X} \quad (18.116)$$

8 This is known as the **Nyström approximation** [WS01]. If we define

9

$$\mathbf{Q}_{A,B} \triangleq \mathbf{K}_{A,Z} \mathbf{K}_{Z,Z}^{-1} \mathbf{K}_{Z,B} \quad (18.117)$$

10 then we can write the approximate Gram matrix as $\mathbf{Q}_{X,X}$. We can then replace \mathbf{K}_σ with $\hat{\mathbf{Q}}_{X,X} =$
11 $\mathbf{Q}_{X,X} + \sigma^2 \mathbf{I}_N$. Computing the eigendecomposition takes $O(M^3)$ time, and computing $\hat{\mathbf{Q}}_{X,X}^{-1}$ takes
12 $O(NM^2)$ time. Thus complexity is now linear in N instead of cubic.

13 If we are approximating *only* $\hat{\mathbf{K}}_{X,X}$ in $\mu_{*|X}$ in Equation (18.53) and $\Sigma_{*|X}$ in Equation (18.54)
14 $\hat{\mathbf{Q}}_{X,X}$, then this is inconsistent with the other un-approximated kernel function evaluations in these
15 formulae, and can result in the predictive variance being negative. One solution to this is to use the
16 same \mathbf{Q} approximation for all terms.

17 18.5.3 Inducing point methods

18 In this section, we discuss an approximation method based on **inducing points**, also called **pseudo**
19 **inputs**, which are like a learned summary of the training data that we can condition on, rather than
20 conditioning on all of it.

21 Let \mathbf{X} be the observed inputs, and $\mathbf{f}_X = f(\mathbf{X})$ be the unknown vector of function values (for which
22 we have noisy observations \mathbf{y}). Let \mathbf{f}_* be the unknown function values at one or more test points
23 \mathbf{X}_* . Finally, let us assume we have M additional inputs, \mathbf{Z} , with unknown function values \mathbf{f}_Z (often
24 denoted by \mathbf{u}). The exact joint prior has the form

25

$$p(\mathbf{f}_X, \mathbf{f}_*) = \int p(\mathbf{f}_*, \mathbf{f}_X, \mathbf{f}_Z) d\mathbf{f}_Z = \int p(\mathbf{f}_*, \mathbf{f}_X | \mathbf{f}_Z) p(\mathbf{f}_Z) d\mathbf{f}_Z = \mathcal{N}\left(\mathbf{0}, \begin{pmatrix} \mathbf{K}_{X,X} & \mathbf{K}_{X,*} \\ \mathbf{K}_{*,X} & \mathbf{K}_{*,*} \end{pmatrix}\right) \quad (18.118)$$

26 (We write $p(\mathbf{f}_X, \mathbf{f}_*)$ instead of $p(\mathbf{f}_X, \mathbf{f}_* | \mathbf{X}, \mathbf{X}_*)$, since the inputs can be thought of as just indices
27 into the random function f .)

28 We will choose \mathbf{f}_Z in such a way that it acts as a sufficient statistic for the data, so that we can
29 predict \mathbf{f}_* just using \mathbf{f}_Z instead of \mathbf{f}_X , i.e., we assume $\mathbf{f}_* \perp \mathbf{f}_X | \mathbf{f}_Z$. Thus we approximate the prior

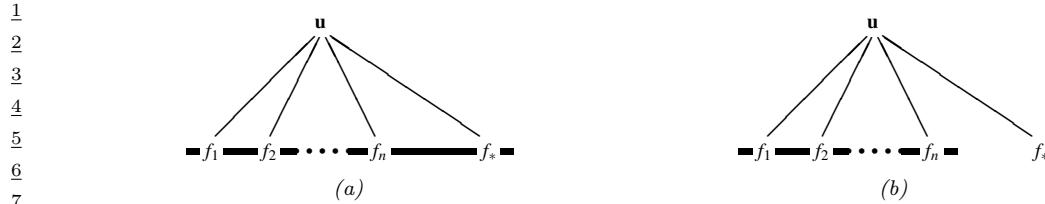


Figure 18.12: Illustration of the graphical model for a GP on n observations, $\mathbf{f}_{1:n}$, and one test case, f_* , with inducing variables \mathbf{u} . The thick lines indicate that all variables are fully interconnected. The observations y_i (not shown) are locally connected to each f_i . (a) no approximations are made. (b) we assume f_* is conditionally independent of \mathbf{f}_X given \mathbf{u} . From Figure 1 of [QCR05]. Used with kind permission of Joaquin Quinonero-Candela.

13

14

15 as follows:

16

$$17 \quad p(\mathbf{f}_*, \mathbf{f}_X, \mathbf{f}_Z) = p(\mathbf{f}_* | \mathbf{f}_X, \mathbf{f}_Z) p(\mathbf{f}_X | \mathbf{f}_Z) p(\mathbf{f}_Z) \approx p(\mathbf{f}_* | \mathbf{f}_Z) p(\mathbf{f}_X | \mathbf{f}_Z) p(\mathbf{f}_Z) \quad (18.119)$$

18

19 See Figure 18.12 for an illustration of this assumption, and Section 18.5.3.4 for details on how to
20 choose the inducing set Z . (Note that this method is often called a “sparse GP”, because it makes
21 predictions for \mathbf{f}_* using a subset of the training data, namely \mathbf{f}_Z , instead of all of it, \mathbf{f}_X .)

22 From this, we can derive the following train and test conditionals

$$23 \quad p(\mathbf{f}_X | \mathbf{f}_Z) = \mathcal{N}(\mathbf{f}_X | \mathbf{K}_{X,Z} \mathbf{K}_{Z,Z}^{-1} \mathbf{f}_Z, \mathbf{K}_{X,X} - \mathbf{Q}_{X,X}) \quad (18.120)$$

$$25 \quad p(\mathbf{f}_* | \mathbf{f}_Z) = \mathcal{N}(\mathbf{f}_* | \mathbf{K}_{*,Z} \mathbf{K}_{Z,Z}^{-1} \mathbf{f}_Z, \mathbf{K}_{*,*} - \mathbf{Q}_{*,*}) \quad (18.121)$$

27 The above equations can be seen as exact inference on noise-free observations \mathbf{f}_Z . To gain
28 computational speedups, we will make further approximations of the form $\tilde{\mathbf{Q}}_{X,X} \approx \mathbf{K}_{X,X} - \mathbf{Q}_{X,X}$
29 and $\tilde{\mathbf{Q}}_{*,*} \approx \mathbf{K}_{*,*} - \mathbf{Q}_{*,*}$. We consider various choices for these values below. All of these choices
30 result in an initial training cost of $O(M^3 + NM^2)$, and then take $O(M)$ time for the predictive mean
31 for each test case, and $O(M^2)$ time for the predictive variance. (Compare this to $O(N^3)$ training
32 time and $O(N)$ and $O(N^2)$ testing time for exact inference.)

33

34 18.5.3.1 SOR/ DIC

35

36 Suppose we assume $\tilde{\mathbf{Q}}_{X,X} = \mathbf{0}$ and $\tilde{\mathbf{Q}}_{*,*} = \mathbf{0}$, so the conditionals are deterministic. This is called the
37 **deterministic inducing conditional (DIC)** approximation [QCR05], or the **subset of regressors**
38 (**SOR**) approximation [Sil85; SB01]. The corresponding joint prior has the form

$$39 \quad q_{\text{SOR}}(\mathbf{f}_X, \mathbf{f}_*) = \mathcal{N}(\mathbf{0}, \begin{pmatrix} \mathbf{Q}_{X,X} & \mathbf{Q}_{X,*} \\ \mathbf{Q}_{*,X} & \mathbf{Q}_{*,*} \end{pmatrix}) \quad (18.122)$$

42

Consequently the predictive distribution is

$$44 \quad q_{\text{SOR}}(\mathbf{f}_* | \mathbf{y}) = \mathcal{N}(\mathbf{f}_* | \mathbf{Q}_{*,X} \hat{\mathbf{Q}}_{X,X}^{-1} \mathbf{y}, \mathbf{Q}_{*,*} - \mathbf{Q}_{*,X} \hat{\mathbf{Q}}_{X,X}^{-1} \mathbf{Q}_{X,*}) \quad (18.123)$$

$$45 \quad = \mathcal{N}(\mathbf{f}_* | \sigma^{-2} \mathbf{K}_{*,Z} \Sigma \mathbf{K}_{Z,X} \mathbf{y}, \mathbf{K}_{*,Z} \Sigma \mathbf{K}_{Z,*}) \quad (18.124)$$

47

where we have defined $\hat{\mathbf{Q}}_{X,X} = \mathbf{Q}_{X,X} + \sigma^2 \mathbf{I}_N$, and $\Sigma = (\sigma^{-2} \mathbf{K}_{Z,X} \mathbf{K}_{X,Z} + \mathbf{K}_{Z,Z})^{-1}$.

This predictive distribution is equivalent to the usual one for GPs except we have replaced $\mathbf{K}_{X,X}$ by $\mathbf{Q}_{X,X}$. This is equivalent to performing GP inference with the following kernel function

$$\mathcal{K}_{\text{SOR}}(\mathbf{x}_i, \mathbf{x}_j) = \mathcal{K}(\mathbf{x}_i, \mathbf{Z}) \mathbf{K}_{Z,Z}^{-1} \mathcal{K}(\mathbf{Z}, \mathbf{x}_j) \quad (18.125)$$

The kernel matrix has rank M , so the GP is degenerate. Furthermore, the kernel will be near 0 when \mathbf{x}_i or \mathbf{x}_j is far from one of the chosen points \mathbf{Z} , which can result in an underestimate of the predictive variance.

18.5.3.2 DTC

One way to overcome the overconfidence of DIC is to only assume $\hat{\mathbf{Q}}_{X,X} = \mathbf{0}$, but let $\tilde{\mathbf{Q}}_{*,*} = \mathbf{K}_{*,*} - \mathbf{Q}_{*,*}$ be exact. This is called the **deterministic training conditional** or **DTC** method [SWL03].

The corresponding joint prior has the form

$$q_{\text{dtc}}(\mathbf{f}_X, \mathbf{f}_*) = \mathcal{N}(\mathbf{0}, \begin{pmatrix} \mathbf{Q}_{X,X} & \mathbf{Q}_{X,*} \\ \mathbf{Q}_{*,X} & \mathbf{K}_{*,*} \end{pmatrix}) \quad (18.126)$$

Hence the predictive distribution becomes

$$q_{\text{dtc}}(\mathbf{f}_* | \mathbf{y}) = \mathcal{N}(\mathbf{f}_* | \mathbf{Q}_{*,X} \hat{\mathbf{Q}}_{X,X}^{-1} \mathbf{y}, \mathbf{K}_{*,*} - \mathbf{Q}_{*,X} \hat{\mathbf{Q}}_{X,X}^{-1} \mathbf{Q}_{X,*}) \quad (18.127)$$

$$= \mathcal{N}(\mathbf{f}_* | \sigma^{-2} \mathbf{K}_{*,Z} \Sigma \mathbf{K}_{Z,X} \mathbf{y}, \mathbf{K}_{*,*} - \mathbf{Q}_{*,*} + \mathbf{K}_{*,Z} \Sigma \mathbf{K}_{Z,*}) \quad (18.128)$$

The predictive mean is the same as in SOR, but the variance is larger (since $\mathbf{K}_{*,*} - \mathbf{Q}_{*,*}$ is positive definite) due to the uncertainty of \mathbf{f}_* given \mathbf{f}_Z .

18.5.3.3 FITC

A widely used approximation assumes $q(\mathbf{f}_X | \mathbf{f}_Z)$ is fully factorized, i.e,

$$q(\mathbf{f}_X | \mathbf{f}_Z) = \prod_{n=1}^N p(f_n | \mathbf{f}_Z) = \mathcal{N}(\mathbf{f}_X | \mathbf{K}_{X,Z} \mathbf{K}_{Z,Z}^{-1} \mathbf{f}_Z, \text{diag}(\mathbf{K}_{X,X} - \mathbf{Q}_{X,X})) \quad (18.129)$$

This is called the **fully independent training conditional** or **FITC** assumption, and was first proposed in [SG06a]. This throws away less uncertainty than the SOR and DTC methods, since it does not make any deterministic assumptions about the relationship between \mathbf{f}_X and \mathbf{f}_Z .

The joint prior has the form

$$q_{\text{fitc}}(\mathbf{f}_X, \mathbf{f}_*) = \mathcal{N}(\mathbf{0}, \begin{pmatrix} \mathbf{Q}_{X,X} - \text{diag}(\mathbf{Q}_{X,X} - \mathbf{K}_{X,X}) & \mathbf{Q}_{X,*} \\ \mathbf{Q}_{*,X} & \mathbf{K}_{*,*} \end{pmatrix}) \quad (18.130)$$

The predictive distribution for a single test case is given by

$$q_{\text{fitc}}(f_* | \mathbf{y}) = \mathcal{N}(f_* | \mathbf{k}_{*,Z} \Sigma \mathbf{K}_{Z,X} \Lambda^{-1} \mathbf{y}, k_{**} - q_{**} + \mathbf{k}_{*,Z} \Sigma \mathbf{k}_{Z,*}) \quad (18.131)$$

where $\Lambda \triangleq \text{diag}(\mathbf{K}_{X,X} - \mathbf{Q}_{X,X} + \sigma^2 \mathbf{I}_N)$, and $\Sigma \triangleq (\mathbf{K}_{Z,Z} + \mathbf{K}_{Z,X} \Lambda^{-1} \mathbf{K}_{X,Z})^{-1}$. If we have a batch of test cases, we can assume they are conditionally independent (an approach known as **fully independent conditional** or **FIC**), and multiply the above equation.

The computational cost is the same as for SOR and DTC, but the approach avoids some of the pathologies due to a non-degenerate kernel. In particular, one can show that the FIC method is equivalent to exact GP inference with the following non-degenerate kernel:

$$\mathcal{K}_{\text{fic}}(\mathbf{x}_i, \mathbf{x}_j) = \begin{cases} \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) & \text{if } i = j \\ \mathcal{K}_{\text{SOR}}(\mathbf{x}_i, \mathbf{x}_j) & \text{if } i \neq j \end{cases} \quad (18.132)$$

18.5.3.4 Learning the inducing points

So far, we have not specified how to choose the inducing points or pseudo inputs \mathbf{Z} . We can treat these like kernel hyperparameters, and choose them so as to maximize the log marginal likelihood, given by

$$\log q(\mathbf{y}|\mathbf{X}, \mathbf{Z}) = \log \int \int p(\mathbf{y}|\mathbf{f}_X) q(\mathbf{f}_X|\mathbf{X}, \mathbf{f}_Z) p(\mathbf{f}_Z|\mathbf{Z}) d\mathbf{f}_Z d\mathbf{f} \quad (18.133)$$

$$= \log \int p(\mathbf{y}|\mathbf{f}_X) q(\mathbf{f}_X|\mathbf{X}, \mathbf{Z}) d\mathbf{f}_X \quad (18.134)$$

$$= -\frac{1}{2} \log |\mathbf{Q}_{X,X} + \Lambda| - \frac{1}{2} \mathbf{y}^\top (\mathbf{Q}_{X,X} + \Lambda)^{-1} \mathbf{y} - \frac{n}{2} \log(2\pi) \quad (18.135)$$

where the definition of Λ depends on the method, namely $\Lambda_{\text{SOR}} = \Lambda_{\text{dtc}} = \sigma^2 \mathbf{I}_N$, and $\Lambda_{\text{fitc}} = \text{diag}(\mathbf{K}_{X,X} - \mathbf{Q}_{X,X}) + \sigma^2 \mathbf{I}_N$.

If the input domain is \mathbb{R}^d , we can optimize $\mathbf{Z} \in \mathbb{R}^{Md}$ using gradient methods. However, one of the appeals of kernel methods is that they can handle structured inputs, such as strings and graphs (see Section 18.2.1.4). In this case, we cannot use gradient methods to select the inducing points. A simple approach is to select the inducing points from the training set, as in the subset of data approach in Section 18.5.1, or using the efficient selection mechanism in [Cao+15]. However, we can also use discrete optimization methods, such as simulated annealing (Section 12.9.1), as discussed in [For+18a]. See Figure 18.13 for an illustration.

18.5.4 Sparse variational methods

In this section, we discuss a variational approach to GP inference that is similar to the inducing point methods in Section 18.5.3, but which generalizes it to also handle non-conjugate likelihoods.

This is called the **sparse variational GP** or **SVGP** approximation. For more details, see [Lei+20].

See also [WKS21] for connections between SVGP and the Nyström method.)

To explain the idea behind SVGP, let us assume, for simplicity, that the function f is defined over a finite set \mathcal{X} of possible inputs, which we partition into three subsets: the training set \mathbf{X} , a set of inducing points \mathbf{Z} , and all other points (which we can think of as the test set), \mathbf{X}_* . (We assume these sets are disjoint.) Let \mathbf{f}_X , \mathbf{f}_Z and \mathbf{f}_* represent the corresponding unknown function values on these points, and let $\mathbf{f} = [\mathbf{f}_X, \mathbf{f}_Z, \mathbf{f}_*]$ be all the unknowns. (Here we work with a fixed-length vector \mathbf{f} , but the result generalizes to Gaussian processes, as explained in [Mat+16].) We assume the function is sampled from a GP, so $p(\mathbf{f}) = \mathcal{N}(m(\mathcal{X}), \mathcal{K}(\mathcal{X}, \mathcal{X}))$.

47

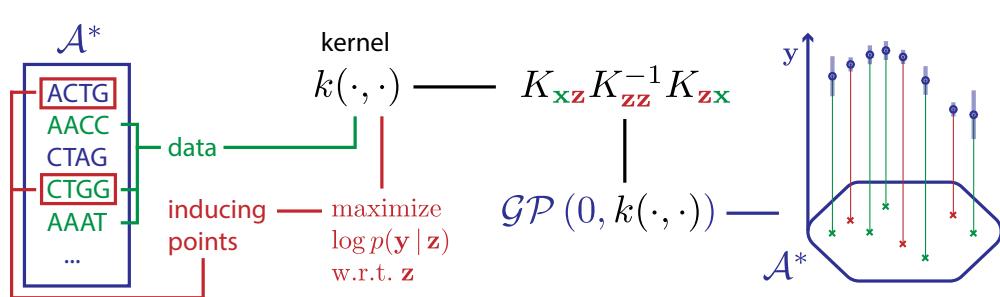


Figure 18.13: Illustration of how to choose inducing points from a discrete input domain (here DNA sequences of length 4) to maximize the log marginal likelihood. From Figure 1 of [For+18a]. Used with kind permission of Vincent Fortuin.

The inducing point methods in Section 18.5.3 approximates the GP prior by assuming $p(\mathbf{f}_*, \mathbf{f}_X, \mathbf{f}_Z) \approx p(\mathbf{f}_* | \mathbf{f}_Z)p(\mathbf{f}_X | \mathbf{f}_Z)p(\mathbf{f}_Z)$. The inducing points \mathbf{f}_Z are chosen to maximize the likelihood of the observed data. We then perform exact inference in this approximate model. By contrast, in this section, we will keep the model unchanged, but we will instead approximate the posterior $p(\mathbf{f} | \mathbf{y})$ using variational inference. This approach is known as the **variational free energy (VFE)** method for GPs [Tit09; Mat+16]; it is also called SVGP.

In the VFE view, the inducing points \mathbf{Z} and inducing variables \mathbf{f}_Z are variational parameters, rather than model parameters, which avoids the risk of overfitting. Furthermore, one can show that as the number of inducing points m increases, the quality of the posterior consistently improves, eventually recovering exact inference. By contrast, in the classical inducing point method, increasing m does not always result in better performance [BWR16].

In more detail, the VFE approach tries to find an approximate posterior $q(\mathbf{f})$ to minimize $D_{\text{KL}}(q(\mathbf{f}) \| p(\mathbf{f} | \mathbf{y}))$. The key assumption is that $q(\mathbf{f}) = q(\mathbf{f}_*, \mathbf{f}_X, \mathbf{f}_Z) = p(\mathbf{f}_*, \mathbf{f}_X | \mathbf{f}_Z)q(\mathbf{f}_Z)$, where $p(\mathbf{f}_*, \mathbf{f}_X | \mathbf{f}_Z)$ is computed exactly using the GP prior, and $q(\mathbf{f}_Z)$ is learned, by minimizing $\mathcal{K}(q) = D_{\text{KL}}(q(\mathbf{f}) \| p(\mathbf{f} | \mathbf{y}))$.² Intuitively, $q(\mathbf{f}_Z)$ acts as a “bottleneck” which “absorbs” all the observations from \mathbf{y} ; posterior predictions for elements of \mathbf{f}_X or \mathbf{f}_* are then made via their dependence on \mathbf{f}_Z , rather than their dependence on each other.

² One can show that $D_{\text{KL}}(q(\mathbf{f}) \| p(\mathbf{f} | \mathbf{y})) = D_{\text{KL}}(q(\mathbf{f}_X, \mathbf{f}_Z) \| p(\mathbf{f}_X, \mathbf{f}_Z | \mathbf{y}))$, which is the original objective from [Tit09].

We can derive the form of the loss, which is used to compute the posterior $q(\mathbf{f}_Z)$, as follows:

$$\mathcal{K}(q) = D_{\text{KL}}(q(\mathbf{f}_*, \mathbf{f}_X, \mathbf{f}_Z) \parallel p(\mathbf{f}_*, \mathbf{f}_X, \mathbf{f}_Z | \mathbf{y})) \quad (18.136)$$

$$= \int q(\mathbf{f}_*, \mathbf{f}_X, \mathbf{f}_Z) \log \frac{q(\mathbf{f}_*, \mathbf{f}_X, \mathbf{f}_Z)}{p(\mathbf{f}_*, \mathbf{f}_X, \mathbf{f}_Z | \mathbf{y})} d\mathbf{f}_* d\mathbf{f}_X d\mathbf{f}_Z \quad (18.137)$$

$$= \int p(\mathbf{f}_*, \mathbf{f}_X | \mathbf{f}_Z) q(\mathbf{f}_Z) \log \frac{\cancel{p(\mathbf{f}_* | \mathbf{f}_X, \mathbf{f}_Z)} \cancel{p(\mathbf{f}_X | \mathbf{f}_Z)} q(\mathbf{f}_Z) p(\mathbf{y})}{\cancel{p(\mathbf{f}_* | \mathbf{f}_X, \mathbf{f}_Z)} \cancel{p(\mathbf{f}_X | \mathbf{f}_Z)} \cancel{p(\mathbf{f}_Z | p(\mathbf{y} | \mathbf{f}_X))}} d\mathbf{f}_* d\mathbf{f}_X d\mathbf{f}_Z \quad (18.138)$$

$$= \int p(\mathbf{f}_*, \mathbf{f}_X | \mathbf{f}_Z) q(\mathbf{f}_Z) \log \frac{q(\mathbf{f}_Z) p(\mathbf{y})}{p(\mathbf{f}_Z) p(\mathbf{y} | \mathbf{f}_X)} d\mathbf{f}_* d\mathbf{f}_X d\mathbf{f}_Z \quad (18.139)$$

$$= \int q(\mathbf{f}_Z) \log \frac{q(\mathbf{f}_Z)}{p(\mathbf{f}_Z)} d\mathbf{f}_Z - \int p(\mathbf{f}_X | \mathbf{f}_Z) q(\mathbf{f}_Z) \log p(\mathbf{y} | \mathbf{f}_X) d\mathbf{f}_X d\mathbf{f}_Z + C \quad (18.140)$$

$$= D_{\text{KL}}(q(\mathbf{f}_Z) \parallel p(\mathbf{f}_Z)) - \mathbb{E}_{q(\mathbf{f}_X)} [\log p(\mathbf{y} | \mathbf{f}_X)] + C \quad (18.141)$$

where $C = \log p(\mathbf{y})$ is an irrelevant constant.

We can alternatively write the objective as an evidence lower bound that we want to maximize:

$$\log p(\mathbf{y}) = \mathcal{K}(q) + \mathbb{E}_{q(\mathbf{f}_X)} [\log p(\mathbf{y} | \mathbf{f}_X)] - D_{\text{KL}}(q(\mathbf{f}_Z) \parallel p(\mathbf{f}_Z)) \quad (18.142)$$

$$\geq \mathbb{E}_{q(\mathbf{f}_X)} [\log p(\mathbf{y} | \mathbf{f}_X)] - D_{\text{KL}}(q(\mathbf{f}_Z) \parallel p(\mathbf{f}_Z)) \triangleq \mathcal{L}(q) \quad (18.143)$$

Now suppose we choose a Gaussian posterior approximation, $q(\mathbf{f}_Z) = \mathcal{N}(\mathbf{f}_Z | \mathbf{m}, \mathbf{S})$. Since $p(\mathbf{f}_Z) = \mathcal{N}(\mathbf{f}_Z | \mathbf{0}, \mathcal{K}(\mathbf{Z}, \mathbf{Z}))$, we can compute the KL term in closed form using the formula for KL divergence between Gaussians (Equation (5.80)).

As for the expected log-likelihood term, we need to compute $q(\mathbf{f}_X)$. Since $q(\mathbf{f}_Z)$ is Gaussian, this can be done in closed form as follows:

$$q(\mathbf{f}_X | \mathbf{m}, \mathbf{S}) = \int p(\mathbf{f}_X | \mathbf{f}_Z, \mathbf{X}, \mathbf{Z}) q(\mathbf{f}_Z | \mathbf{m}, \mathbf{S}) d\mathbf{f}_Z = \mathcal{N}(\mathbf{f}_X | \tilde{\mu}, \tilde{\Sigma}) \quad (18.144)$$

$$\tilde{\mu}_i = m(\mathbf{x}_i) + \boldsymbol{\alpha}(\mathbf{x}_i)^T (\mathbf{m} - m(\mathbf{Z})) \quad (18.145)$$

$$\tilde{\Sigma}_{ij} = \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) - \boldsymbol{\alpha}(\mathbf{x}_i)^T (\mathcal{K}(\mathbf{Z}, \mathbf{Z}) - \mathbf{S}) \boldsymbol{\alpha}(\mathbf{x}_j) \quad (18.146)$$

$$\boldsymbol{\alpha}(\mathbf{x}_i) = \mathcal{K}(\mathbf{Z}, \mathbf{Z})^{-1} \mathcal{K}(\mathbf{Z}, \mathbf{x}_i) \quad (18.147)$$

Hence $q(f_n) = \mathcal{N}(f_n | \tilde{\mu}_n, \tilde{\Sigma}_{nn})$, which we can use to compute the expected log likelihood:

$$\mathbb{E}_{q(\mathbf{f}_X)} [\log p(\mathbf{y} | \mathbf{f}_X)] = \sum_{n=1}^N \mathbb{E}_{q(f_n)} [\log p(y_n | f_n)] \quad (18.148)$$

We discuss how to compute the expected loglikelihood below.

18.5.4.1 Gaussian likelihood

If we have a Gaussian observation model, we can compute the expected log likelihood in closed form.

In particular, if we assume $m(\mathbf{x}) = \mathbf{0}$, we have

$$\mathbb{E}_{q(f_n)} [\log \mathcal{N}(y_n | f_n, \beta^{-1})] = \log \mathcal{N}(y_n | \mathbf{k}_n^T \mathbf{K}_{Z,Z}^{-1} \mathbf{m}, \beta^{-1}) - \frac{1}{2} \beta \tilde{k}_{nn} - \frac{1}{2} \text{tr}(\mathbf{S} \boldsymbol{\Lambda}_n) \quad (18.149)$$

where $\tilde{k}_{nn} = k_{nn} - \mathbf{k}_n^\top \mathbf{K}_{Z,Z}^{-1} \mathbf{k}_n$, \mathbf{k}_n is the n 'th column of $\mathbf{K}_{Z,X}$ and $\Lambda_n = \beta \mathbf{K}_{Z,Z}^{-1} \mathbf{k}_n \mathbf{k}_n^\top \mathbf{K}_{Z,Z}^{-1}$. Hence the overall ELBO has the form

$$\mathcal{L}(q) = \log \mathcal{N}(\mathbf{y} | \mathbf{K}_{X,Z} \mathbf{K}_{Z,Z}^{-1} \mathbf{m}, \beta^{-1} \mathbf{I}_N) - \frac{1}{2} \beta \text{tr}(\mathbf{K}_{X,Z} \mathbf{K}_{Z,Z}^{-1} \mathbf{S} \mathbf{K}_{Z,Z}^{-1} \mathbf{K}_{Z,X}) \quad (18.150)$$

$$- \frac{1}{2} \beta \text{tr}(\mathbf{K}_{X,X} - \mathbf{Q}_{X,X}) - D_{\text{KL}}(q(\mathbf{f}_Z) \| p(\mathbf{f}_Z)) \quad (18.151)$$

To compute the gradients of this, we leverage the following result [OA09]:

$$\frac{\partial}{\partial \mu} \mathbb{E}_{\mathcal{N}(x|\mu, \sigma^2)} [h(x)] = \mathbb{E}_{\mathcal{N}(x|\mu, \sigma^2)} \left[\frac{\partial}{\partial x} h(x) \right] \quad (18.152)$$

$$\frac{\partial}{\partial \sigma^2} \mathbb{E}_{\mathcal{N}(x|\mu, \sigma^2)} [h(x)] = \frac{1}{2} \mathbb{E}_{\mathcal{N}(x|\mu, \sigma^2)} \left[\frac{\partial^2}{\partial x^2} h(x) \right] \quad (18.153)$$

We then substitute $h(x)$ with $\log p(y_n|f_n)$. Using this, one can show

$$\nabla_{\mathbf{m}} \mathcal{L}(q) = \beta \mathbf{K}_{Z,Z}^{-1} \mathbf{K}_{Z,X} \mathbf{y} - \Lambda \mathbf{m} \quad (18.154)$$

$$\nabla_{\mathbf{S}} \mathcal{L}(q) = \frac{1}{2} \mathbf{S}^{-1} - \frac{1}{2} \Lambda \quad (18.155)$$

Setting the derivatives to zero gives the optimal solution:

$$\mathbf{S} = \Lambda^{-1} \quad (18.156)$$

$$\Lambda = \beta \mathbf{K}_{Z,Z}^{-1} \mathbf{K}_{Z,X} \mathbf{K}_{X,Z} \mathbf{K}_{Z,Z}^{-1} + \mathbf{K}_{Z,Z}^{-1} \quad (18.157)$$

$$\mathbf{m} = \beta \Lambda^{-1} \mathbf{K}_{Z,Z}^{-1} \mathbf{K}_{Z,X} \mathbf{y} \quad (18.158)$$

This is called **sparse GP regression** or **SGPR** [Tit09].

With these parameters, the lower bound on the log marginal likelihood is given by

$$\log p(\mathbf{y}) \geq \log \mathcal{N}(\mathbf{y} | \mathbf{0}, \mathbf{K}_{X,Z} \mathbf{K}_{Z,Z}^{-1} \mathbf{K}_{Z,X} + \beta^{-1} \mathbf{I}) - \frac{1}{2} \beta \text{tr}(\mathbf{K}_{X,X} - \mathbf{Q}_{X,X}) \quad (18.159)$$

where $\mathbf{Q}_{X,X} = \mathbf{K}_{X,Z} \mathbf{K}_{Z,Z}^{-1} \mathbf{K}_{Z,X}$. (This is called the “collapsed” lower bound, since we have marginalized out \mathbf{f}_Z .) If $Z = X$, then $\mathbf{K}_{Z,Z} = \mathbf{K}_{Z,X} = \mathbf{K}_{X,X}$, so the bound becomes tight, and we have $\log p(\mathbf{y}) = \log \mathcal{N}(\mathbf{y} | \mathbf{0}, \mathbf{K}_{X,X} + \beta^{-1} \mathbf{I})$.

Equation (18.159) is almost the same as the log marginal likelihood for the DTC model in Equation (18.135), except for the trace term; it is this latter term that prevents overfitting, due to the fact that we treat \mathbf{f}_Z as variational parameters of the posterior rather than model parameters of the prior.

18.5.4.2 Non-Gaussian likelihood

In this section, we briefly consider the case of non-Gaussian likelihoods, which arise when using GPs for classification or for count data (see Section 18.4). We can compute the gradients of the expected log likelihood by defining $h(f_n) = \log p(y_n|f_n)$ and then using a Monte Carlo approximation to Equation (18.152) and Equation (18.153). In the case of a binary classifier, we can use the results in Table 18.2 to compute the inner $\frac{\partial}{\partial f_n} h(f_n)$ and $\frac{\partial^2}{\partial f_n^2} h(f_n)$ terms.

1 2 **18.5.4.3 Minibatch SVI**

3 Computing the optimal variational solution in Section 18.5.4.1 requires solving a batch optimization
4 problem, which takes $O(M^3 + NM^2)$ time. This may still be too slow if N is large, unless M is
5 small, which compromises accuracy.

6 An alternative approach is to perform stochastic optimization of the VFE objective, instead of
7 batch optimization. This is known as stochastic variational inference (see Section 10.3.1). The
8 key observation is that the log likelihood in Equation (18.148) is a sum of N terms, which we can
9 approximate with minibatch sampling to compute noisy estimates of the gradient, as proposed in
10 [HFL13].

11 In more detail, the objective becomes

$$\mathcal{L}(q) = \left[\frac{N}{B} \sum_{b=1}^B \frac{1}{|\mathcal{B}_b|} \sum_{n \in \mathcal{B}_b} \mathbb{E}_{q(f_n)} [\log p(y_n | f_n)] \right] - D_{\text{KL}}(q(\mathbf{f}_Z) \| (p(\mathbf{f}_Z))) \quad (18.160)$$

16 where \mathcal{B}_b is the b 'th batch, and B is the number of batches. Since the GP model (with Gaussian
17 likelihoods) is in the exponential family, we can efficiently compute the natural gradient (Section 6.4)
18 of Equation (18.160) wrt the canonical parameters of $q(\mathbf{f}_Z)$; this converges much faster than following
19 the standard gradient. See [HFL13] for details.
20

21 22 **18.5.5 Exploiting parallelization and structure via kernel matrix multiplies**

23 It takes $O(N^3)$ time to compute the Cholesky decomposition of $\mathbf{K}_{X,X}$, which is needed to solve the
24 linear system $\mathbf{K}_\sigma \boldsymbol{\alpha} = \mathbf{y}$ and to compute $|\mathbf{K}_{X,X}|$. An alternative to Cholesky decomposition is to
25 use linear algebra methods, often called **Krylov subspace methods** based just on **matrix vector**
26 **multiplication** or **MVM**. These approaches are often much faster.

27 In short, if the kernel matrix $\mathbf{K}_{X,X}$ has special algebraic structure, which is often the case through
28 either the choice of kernel or the structure of the inputs, then it is typically easier to exploit this
29 structure in performing fast matrix multiplies. Moreover, even if the kernel matrix **does not** have
30 special structure, matrix multiplies are trivial to parallelize, and can thus be greatly accelerated by
31 GPUs, unlike Cholesky based methods which are largely sequential. Algorithms based on matrix
32 multiplies are in harmony with modern hardware advances, which enable significant parallelization.
33

34 35 **18.5.5.1 Using conjugate gradient and Lanczos methods**

36 We can solve the linear system $\mathbf{K}_\sigma \boldsymbol{\alpha} = \mathbf{y}$ using conjugate gradients (CG). The key computational
37 step in CG is the ability to perform MVMs. Let $\tau(\mathbf{K}_\sigma)$ be the time complexity of a single MVM
38 with \mathbf{K}_σ . For a dense $n \times n$ matrix, we have $\tau(\mathbf{K}_\sigma) = n^2$; however, we can speed this up if \mathbf{K}_σ is
39 sparse or structured, as we discuss below.

40 Even if \mathbf{K}_σ is dense, we may still be able to save time by solving the linear system approximately.
41 In particular, if we perform C iterations, CG will take $O(C\tau(\mathbf{K}_\sigma))$ time. If we run for $C = n$, and
42 $\tau(\mathbf{K}_\sigma) = n^2$, it gives the exact solution in $O(n^3)$ time. However, often we can use fewer iterations
43 and still get good accuracy, depending on the condition number of \mathbf{K}_σ .

44 We can compute the log determinant of a matrix using the MVM primitive with a similar iterative
45 method known as **stochastic Lanczos quadrature** [UCS17; Don+17a]. This takes $O(L\tau(\mathbf{K}_\sigma))$
46 time for L iterations.

47

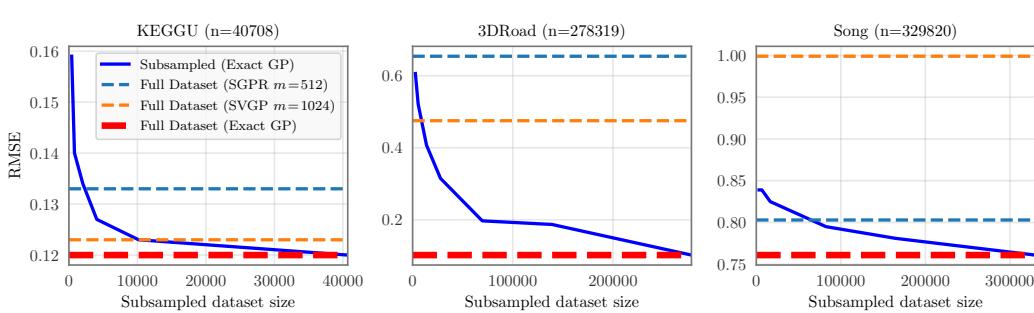


Figure 18.14: RMSE on test set as a function of training set size using a GP with Matern 3/2 kernel with shared lengthscale across all dimensions. Solid lines: exact inference. Dashed blue: SGPR method (closed-form batch solution to the Gaussian variational approximation) of Section 18.5.4.1 with $M = 512$ inducing points. Dashed orange: SVGP method (SGD on Gaussian variational approximation) of Section 18.5.4.3 with $M = 1024$ inducing points. Number of input dimensions: KEGGU $D = 27$, 3DRoad $D = 3$, Song $D = 90$. From Figure 4 of [Wan+19a]. Used with kind permission of Andrew Wilson.

These methods have been used in the **blackbox matrix-matrix multiplication** (BBMM) inference procedure of [Gar+18a], which formulates a batch approach to CG that can be effectively parallelized on GPUs. Using 8 GPUs, this enabled the authors of [Wan+19a] to perform exact inference for a GP regression model on $N \sim 10^4$ datapoints in seconds, $N \sim 10^5$ datapoints in minutes, and $N \sim 10^6$ datapoints in hours.

Interestingly, Figure 18.14 shows that exact GP inference on a subset of the data can often outperform approximate inference on the full data. We also see that performance of exact GPs continues to significantly improve as we increase the size of the data, suggesting that GPs are not only useful in the small-sample setting. In particular, the BBMM is an exact method, and so will preserve the non-parametric representation of a GP with a non-degenerate kernel. By contrast, standard scalable approximations typically operate by replacing the exact kernel with an approximation that corresponds to a parametric model. The non-parametric GPs are able to grow their capacity with more data, benefiting more significantly from the structure present in large datasets.

18.5.5.2 Kernels with compact support

Suppose we use a kernel with **compact support**, where $\mathcal{K}(\mathbf{x}, \mathbf{x}') = 0$ if $\|\mathbf{x} - \mathbf{x}'\| > \epsilon$ for some threshold ϵ (see e.g., [MR09]), then \mathbf{K}_σ will be sparse, so $\tau(\mathbf{K}_\sigma)$ will be $O(N)$. We can also induce sparsity and structure in other ways, as we discuss in Section 18.5.5.3.

18.5.5.3 KISS

One way to ensure that MVMs are fast is to force the kernel matrix to have structure. The **structured kernel interpolation** (SKI) method of [WN15] does this as follows. First it assumes we have a set of inducing points, with Gram matrix $\mathbf{K}_{Z,Z}$. It then interpolates these values to predict the entries of the full kernel matrix using

$$\mathbf{K}_{X,X} \approx \mathbf{W}_X \mathbf{K}_{Z,Z} \mathbf{W}_X^\top \quad (18.161)$$

1 where \mathbf{W}_X is a sparse matrix containing interpolation weights. If we use cubic interpolation, each
 2 row only has 4 nonzeros. Thus we can compute $(\mathbf{W}_X \mathbf{K}_{Z,Z} \mathbf{W}_X^\top) \mathbf{v}$ for any vector \mathbf{v} in $O(N + M^2)$
 3 time.
 4

5 Note that the **SKI** approach generalizes all inducing point methods. For example, we can recover the
 6 subset of regressors method (SOR) method by setting the interpolation weights to $\mathbf{W} = \mathbf{K}_{X,Z} \mathbf{K}_{Z,Z}^{-1}$.
 7 We can identify this procedure as performing a global Gaussian process interpolation strategy on the
 8 user specified kernel. See [WN15] and [WDN15] for more details.

9 In 1d, we can further reduce the running time by choosing the inducing points to be on a regular
 10 grid, so that $\mathbf{K}_{Z,Z}$ is a Toeplitz matrix. In higher dimensions, we need to use a multidimensional grid
 11 of points, resulting in $\mathbf{K}_{Z,Z}$ being a Kronecker product of Toeplitz matrices. This enables matrix
 12 vector multiplication in $O(N + M \log M)$ time and $O(N + M)$ space. The resulting method is called
 13 **KISS-GP** [WN15], which stands for “kernel interpolation for scalable, structured GPs”.

14 Unfortunately, the KISS method can take exponential time in the input dimensions D when
 15 exploiting Kronecker structure in $\mathbf{K}_{Z,Z}$, due to the need to create a fully connected multidimensional
 16 lattice. In [Gar+18b], they propose a method called **SKIP**, which stands for “SKI for products”.
 17 The idea is to leverage the fact that many kernels (including ARD) can be written as a product of 1d
 18 kernels: $\mathcal{K}(\mathbf{x}, \mathbf{x}') = \prod_{d=1}^D \mathcal{K}^d(\mathbf{x}, \mathbf{x}')$. This can be combined with the 1d SKI method to enable fast
 19 MVMs. The overall running time to compute the log marginal likelihood (which is the bottleneck
 20 for kernel learning) using C iterations of CG and a Lanczos decomposition of rank L , becomes
 21 $O(DL(N + M \log M) + L^3 N \log D + CL^2 N)$. Typical values are $L \sim 10^1$ and $C \sim 10^2$.

22

23 18.5.5.4 Tensor train methods

24 Consider the Gaussian VFE approach in Section 18.5.4. We have to estimate the covariance \mathbf{S} and
 25 the mean \mathbf{m} . We can represent \mathbf{S} efficiently using Kronecker structure, as used by KISS. Additionally,
 26 we can represent \mathbf{m} efficiently using the **tensor train decomposition** [Ose11] in combination with
 27 **SKI** [WN15]. The resulting **TT-GP** method can scale efficiently to billions of inducing points, as
 28 explained in [INK18].
 29

30

31 18.5.6 Converting a GP to a SSM

32 Consider a function defined on a 1d scalar input, such as a time index. For many kernels, the
 33 corresponding GP can be modeled using a stochastic differential equation. This induces a block
 34 tri-diagonal precision matrix for the posterior $p(\mathbf{f}_{1:T} | \mathbf{y}_{1:T})$. We can therefore convert this to a
 35 linear-Gaussian state space model (Section 29.1), and perform exact inference in $O(T)$ time using
 36 Kalman smoothing, as explained in [SSH13; Ada+20]. This conversion can be done exactly for
 37 Matern kernels and approximately for Gaussian (RBF) kernels (see [SS19, Ch. 12]). In [SGF21], they
 38 describe how to reduce the linear dependence on T to $\log(T)$ time using a **parallel prefix scan**
 39 operator, that can be run efficiently on GPUs.
 40

41

42 18.6 Learning the kernel

43

44 In [Mac98], David MacKay asked: “How can Gaussian processes replace neural networks? Have we
 45 thrown the baby out with the bathwater?” This remark was made in the late 1990s, at the end of
 46 the second wave of neural networks. Researchers and practitioners had grown weary of the design
 47

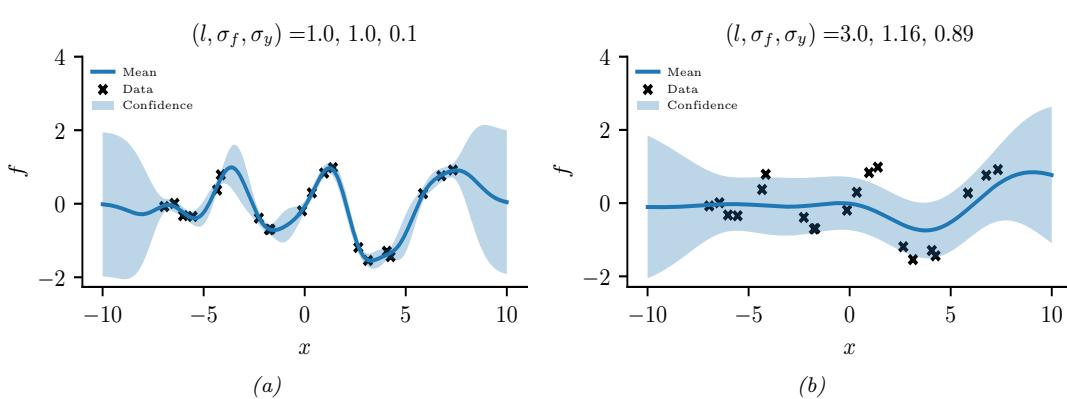


Figure 18.15: Some 1d GPs with RBF kernels but different hyper-parameters fit to 20 noisy observations. The hyper-parameters $(\ell, \sigma_f, \sigma_y)$ are as follows: (a) $(1, 1, 0.1)$ (b) $(3.0, 1.16, 0.89)$. Adapted from Figure 2.5 of [RW06]. Generated by [gpr_demo_change_hparams.ipynb](#).

decisions associated with neural networks — such as activation functions, optimization procedures, architecture design — and the lack of a principled framework to make these decisions. Gaussian processes, by contrast, were perceived as flexible and principled probabilistic models, which naturally followed from Radford Neal’s results on infinite neural networks [Nea96], which we discuss in more depth in Section 18.7.

However, MacKay [Mac98] noted that neural networks could discover rich representations of data through adaptive hidden basis functions, while Gaussian processes with standard kernel functions, such as the RBF kernel, are essentially just smoothing devices. Indeed, the generalization properties of Gaussian processes hinge on the suitability of the kernel function. *Learning* the kernel is how we do representation learning with Gaussian processes, and in many cases will be crucial for good performance — especially when we wish to perform extrapolation, making predictions far away from the data [WA13; Wil+14].

As we will see, learning a kernel is in many ways analogous to training a neural network. Moreover, neural networks and Gaussian processes can be synergistically combined through approaches such as deep kernel learning (see Section 18.6.6) and NN-GPs (Section 18.7.2).

18.6.1 Empirical Bayes for the kernel parameters

Suppose, as in Section 18.3.2, we are performing 1d regression using a GP with an RBF kernel. Since the data has observation noise, the kernel has the following form:

$$\mathcal{K}_y(x_p, x_q) = \sigma_f^2 \exp\left(-\frac{1}{2\ell^2}(x_p - x_q)^2\right) + \sigma_y^2 \delta_{pq} \quad (18.162)$$

Here ℓ is the horizontal scale over which the function changes, σ_f^2 controls the vertical scale of the function, and σ_y^2 is the noise variance. Figure 18.15 illustrates the effects of changing these parameters. We sampled 20 noisy data points from the SE kernel using $(\ell, \sigma_f, \sigma_y) = (1, 1, 0.1)$, and then made predictions various parameters, conditional on the data. In Figure 18.15(a), we use

$(\ell, \sigma_f, \sigma_y) = (1, 1, 0.1)$, and the result is a good fit. In Figure 18.15(b), we increase the length scale to $\ell = 3$; now the function looks smoother, but we are arguably underfitting.

To estimate the kernel parameters θ (sometimes called hyperparameters), we could use exhaustive search over a discrete grid of values, with validation loss as an objective, but this can be quite slow. (This is the approach used by nonprobabilistic methods, such as SVMs, to tune kernels.) Here we consider an empirical Bayes approach, which will allow us to use continuous optimization methods, which are much faster. In particular, we will maximize the marginal likelihood

$$p(\mathbf{y}|\mathbf{X}, \theta) = \int p(\mathbf{y}|\mathbf{f}, \mathbf{X})p(\mathbf{f}|\mathbf{X}, \theta)d\mathbf{f} \quad (18.163)$$

(The reason it is called the marginal likelihood, rather than just likelihood, is because we have marginalized out the latent Gaussian vector \mathbf{f} .) Since $p(\mathbf{f}|\mathbf{X}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K})$, and $p(\mathbf{y}|\mathbf{f}) = \prod_{n=1}^N \mathcal{N}(y_n|f_n, \sigma_y^2)$, the marginal likelihood is given by

$$\log p(\mathbf{y}|\mathbf{X}, \theta) = \log \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K}_\sigma) = -\frac{1}{2}\mathbf{y}^\top \mathbf{K}_\sigma^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K}_\sigma| - \frac{N_D}{2} \log(2\pi) \quad (18.164)$$

where the dependence of \mathbf{K}_σ on θ is implicit. The first term is a data fit term, the second term is a model complexity term, and the third term is just a constant. To understand the tradeoff between the first two terms, consider a SE kernel in 1D, as we vary the length scale ℓ and hold σ_y^2 fixed. Let $J(\ell) = -\log p(\mathbf{y}|\mathbf{X}, \ell)$. For short length scales, the fit will be good, so $\mathbf{y}^\top \mathbf{K}_\sigma^{-1} \mathbf{y}$ will be small. However, the model complexity will be high: \mathbf{K} will be almost diagonal, since most points will not be considered “near” any others, so the $\log |\mathbf{K}_\sigma|$ will be large. For long length scales, the fit will be poor but the model complexity will be low: \mathbf{K} will be almost all 1’s, so $\log |\mathbf{K}_\sigma|$ will be small.

We now discuss how to maximize the marginal likelihood. One can show that

$$\frac{\partial}{\partial \theta_j} \log p(\mathbf{y}|\mathbf{X}, \theta) = \frac{1}{2} \mathbf{y}^\top \mathbf{K}_\sigma^{-1} \frac{\partial \mathbf{K}_\sigma}{\partial \theta_j} \mathbf{K}_\sigma^{-1} \mathbf{y} - \frac{1}{2} \text{tr}(\mathbf{K}_\sigma^{-1} \frac{\partial \mathbf{K}_\sigma}{\partial \theta_j}) \quad (18.165)$$

$$= \frac{1}{2} \text{tr} \left((\boldsymbol{\alpha} \boldsymbol{\alpha}^\top - \mathbf{K}_\sigma^{-1}) \frac{\partial \mathbf{K}_\sigma}{\partial \theta_j} \right) \quad (18.166)$$

where $\boldsymbol{\alpha} = \mathbf{K}_\sigma^{-1} \mathbf{y}$. It takes $O(N^3)$ time to compute \mathbf{K}_σ^{-1} , and then $O(N^2)$ time per hyper-parameter to compute the gradient.

The form of $\frac{\partial \mathbf{K}_\sigma}{\partial \theta_j}$ depends on the form of the kernel, and which parameter we are taking derivatives with respect to. Often we have constraints on the hyper-parameters, such as $\sigma_y^2 \geq 0$. In this case, we can define $\theta = \log(\sigma_y^2)$, and then use the chain rule.

Given an expression for the log marginal likelihood and its derivative, we can estimate the kernel parameters using any standard gradient-based optimizer. However, since the objective is not convex, local minima can be a problem, as we illustrate below, so we may need to use multiple restarts.

41

42 18.6.1.1 Example

43

44 Consider Figure 18.16. We use the SE kernel in Equation (18.162) with $\sigma_f^2 = 1$, and plot 45 $\log p(\mathbf{y}|\mathbf{X}, \ell, \sigma_y^2)$ (where \mathbf{X} and \mathbf{y} are the 7 data points shown in panels b and c as we vary ℓ 46 and σ_y^2). The two local optima are indicated by + in panel (a). The bottom left optimum corresponds 47

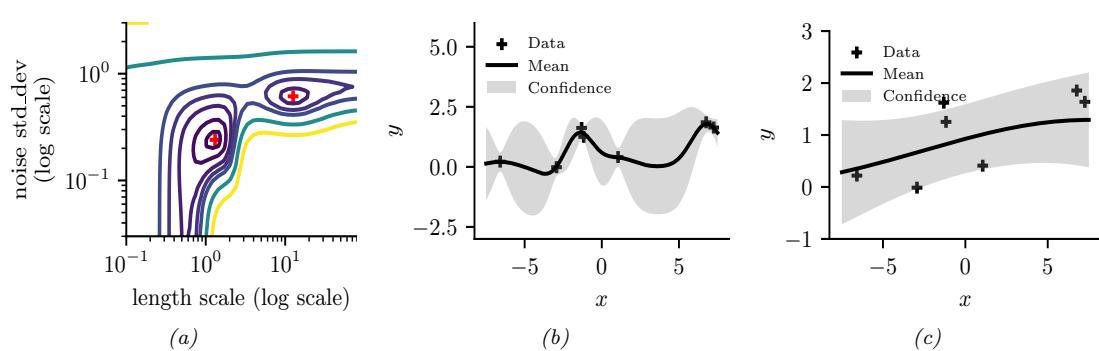


Figure 18.16: Illustration of local minima in the marginal likelihood surface. (a) We plot the log marginal likelihood vs σ_y^2 and ℓ , for fixed $\sigma_f^2 = 1$, using the 7 data points shown in panels b and c. (b) The function corresponding to the lower left local minimum, $(\ell, \sigma_n^2) \approx (1, 0.2)$. This is quite “wiggly” and has low noise. (c) The function corresponding to the top right local minimum, $(\ell, \sigma_n^2) \approx (10, 0.8)$. This is quite smooth and has high noise. The data was generated using $(\ell, \sigma_n^2) = (1, 0.1)$. Adapted from Figure 5.5 of [RW06]. Generated by [gpr_demo_marglik.ipynb](#).

to a low-noise, short-length scale solution (shown in panel b). The top right optimum corresponds to a high-noise, long-length scale solution (shown in panel c). With only 7 data points, there is not enough evidence to confidently decide which is more reasonable, although the more complex model (panel b) has a marginal likelihood that is about 60% higher than the simpler model (panel c). With more data, the more complex model would become even more preferred.

Figure 18.16 illustrates some other interesting (and typical) features. The region where $\sigma_y^2 \approx 1$ (top of panel a) corresponds to the case where the noise is very high; in this regime, the marginal likelihood is insensitive to the length scale (indicated by the horizontal contours), since all the data is explained as noise. The region where $\ell \approx 0.5$ (left hand side of panel a) corresponds to the case where the length scale is very short; in this regime, the marginal likelihood is insensitive to the noise level (indicated by the vertical contours), since the data is perfectly interpolated. Neither of these regions would be chosen by a good optimizer.

18.6.2 Bayesian inference for the kernel parameters

When we have a small number of datapoints (e.g., when using GPs for blackbox optimization, as we discuss in Section 6.8), using a point estimate of the kernel parameters can give poor results [Bul11; WF14]. As a simple example, if the function values that have been observed so far are all very similar, then we may estimate $\hat{\sigma} \approx 0$, which will result in overly confident predictions.³

To overcome such overconfidence, we can compute a posterior over the kernel parameters. If the dimensionality of θ is small, we can compute a discrete grid of possible values, centered on the MAP

³ In [WSN00; BBV11b], they show how we can put a conjugate prior on σ^2 and integrate it out, to generate a Student version of the GP, which is more robust.

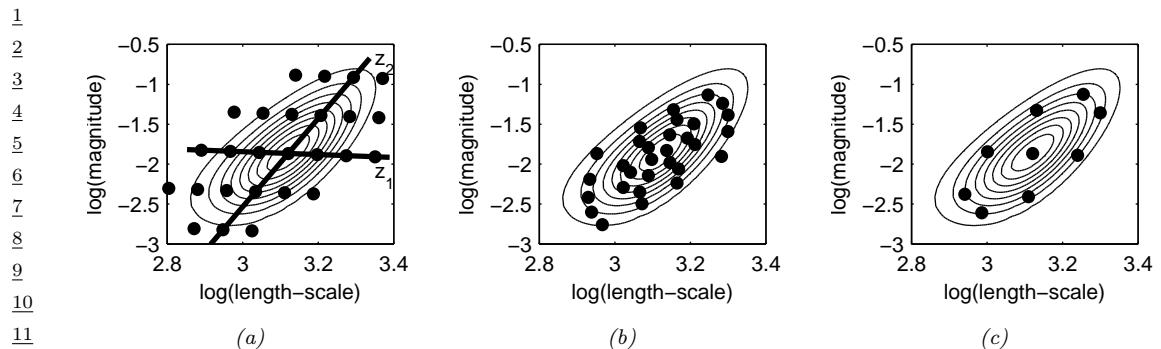


Figure 18.17: Three different approximations to the posterior over hyper-parameters: grid-based, Monte Carlo, and central composite design. From Figure 3.2 of [Van10]. Used with kind permission of Jarno Vanhatalo.

estimate $\hat{\theta}$ (computed as above). We can then approximate the posterior using

$$p(\mathbf{f}|\mathcal{D}) = \sum_{s=1}^S p(\mathbf{f}|\mathcal{D}, \boldsymbol{\theta}_s) p(\boldsymbol{\theta}_s|\mathcal{D}) w_s \quad (18.167)$$

where w_s denotes the weight for grid point s .

In higher dimensions, a regular grid suffers from the curse of dimensionality. One alternative is to place grid points at the mode, and at a distance $\pm 1\text{sd}$ from the mode along each dimension, for a total of $2|\theta| + 1$ points. This is called a **central composite design** [RMC09]. See Figure 18.17 for an illustration.

In higher dimensions, we can use Monte Carlo inference for the kernel parameters when computing Equation (18.167). For example, [MA10] shows how to use slice sampling (Section 12.4.1) for this task, [Hen+15] shows how to use HMC (Section 12.5), and [BBV11a] shows how to use SMC (Chapter 13).

In Figure 18.18, we illustrate the difference between kernel optimization vs kernel inference. We fit a 1d dataset using a kernel of the form

$$\frac{32}{32} \quad \mathcal{K}(r) = \sigma_1^2 \mathcal{K}_{\text{SE}}(r; \tau) \mathcal{K}_{\cos}(r; \rho_1) + \sigma_2^s \mathcal{K}_{32}(r; \rho_2) \quad (18.168)$$

³⁴ where $\mathcal{K}_{\text{SE}}(r; \ell)$ is the squared exponential kernel (Equation (18.13)), $\mathcal{K}_{\cos}(r; \rho_1)$ is the cosine kernel (Equation (18.20)), and $\mathcal{K}_{32}(r; \rho_2)$ is the Matern $\frac{3}{2}$ kernel (Equation (18.16)). We then compute a point-estimate of the kernel parameters using empirical Bayes, and posterior samples using HMC. We then predicting the posterior mean of f on a 1d test set by plugging in the MLE or averaging over samples. We see that the latter captures more uncertainty (beyond the uncertainty captured by the Gaussian itself).

18.6.3 Multiple kernel learning for additive kernels

⁴³ A special case of kernel learning arises when the kernel is a sum of B base kernels

$$\frac{44}{45} \quad \mathcal{K}(\mathbf{x}, \mathbf{x}') = \sum_{b=1}^B w_b \mathcal{K}_b(\mathbf{x}, \mathbf{x}') \quad (18.169)$$

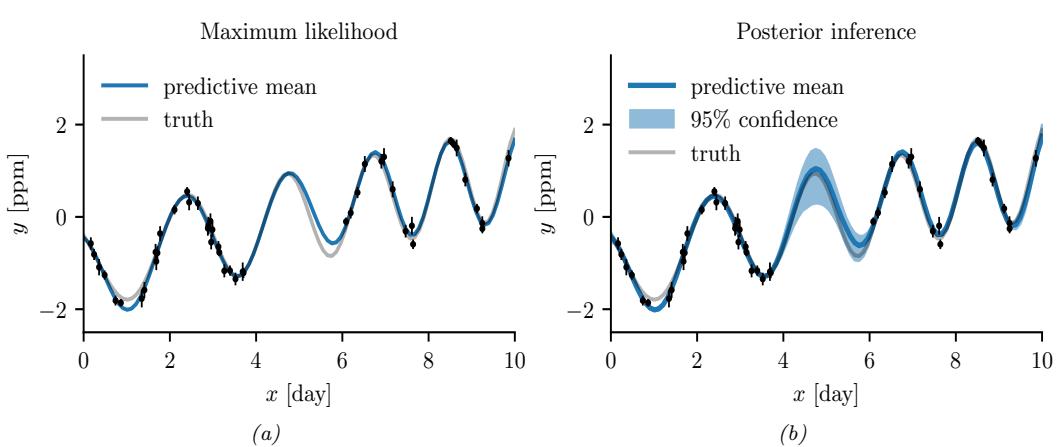


Figure 18.18: Difference between estimation and inference for kernel hyper-parameters. (a) Empirical Bayes approach based on optimization. We plot the posterior predicted mean given a plug-in estimate, $\mathbb{E}[f(x)|\mathcal{D}, \hat{\theta}]$. (b) Bayesian approach based on HMC. We plot the posterior predicted mean, marginalizing over hyper-parameters, $\mathbb{E}[f(x)|\mathcal{D}]$. Generated by [gp_kernel_opt.ipynb](#).

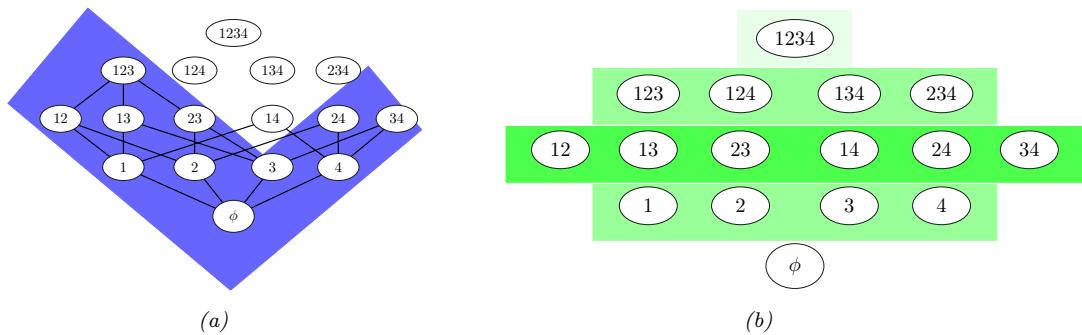


Figure 18.19: Comparison of different additive model classes for a 4d function. Circles represent different interaction terms, ranging from first-order to fourth-order. Left: hierarchical kernel learning uses a nested hierarchy of terms. Right: additive GPs use a weighted sum of additive kernels of different orders. Color shades represent different weighting terms. Adapted from Figure 6.2 of [Duv14].

Optimizing the weights $w_b > 0$ using structural risk minimization is known as **multiple kernel learning**; see e.g., [Rak+08] for details.

Now suppose we constrain the base kernels to depend on a subset of the variables. Furthermore, suppose we enforce a hierarchical inclusion property (e.g., including the kernel k_{123} means we must also include k_{12} , k_{13} and k_{23}), as illustrated in Figure 18.19(left). This is called **hierarchical kernel learning**. We can find a good subset from this model class using convex optimization [Bac09]; however, this requires the use of cross validation to estimate the weights. A more efficient approach is to use the empirical Bayes approach described in [DNR11].

In many cases, it is common to restrict attention to first order additive kernels, i.e., $\mathcal{K}(\mathbf{x}, \mathbf{x}') =$

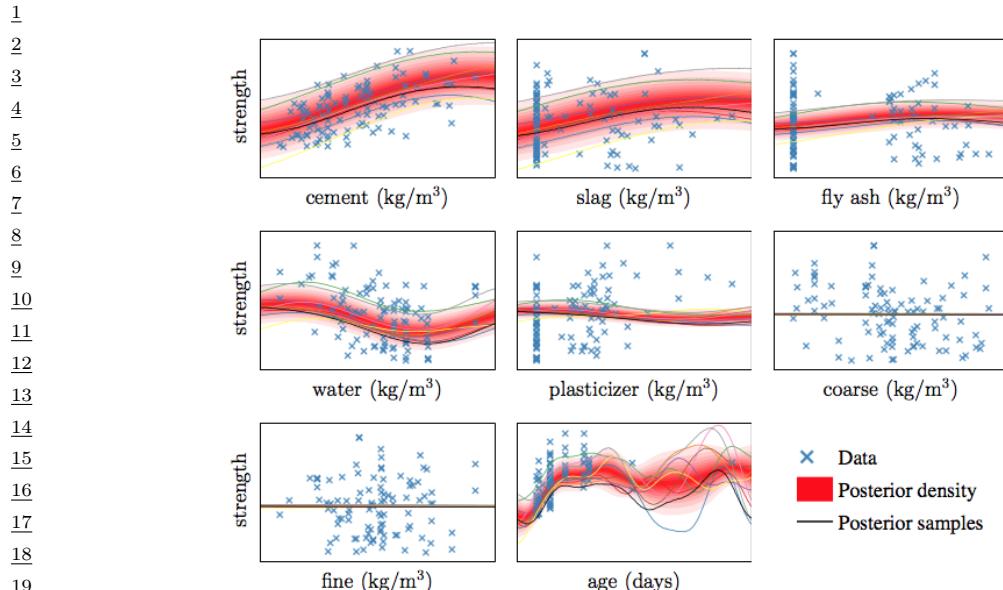


Figure 18.20: Predictive distribution of each term in a GP-GAM model applied to a dataset with 8 continuous inputs and 1 continuous output, representing the strength of some concrete. From Figure 2.7 of [Duv14]. Used with kind permission of David Duvenaud.

$\sum_{d=1}^D \mathcal{K}_d(x_d, x'_d)$. The resulting function has the form

$$f(\mathbf{x}) = f_1(x_1) + \dots + f_D(x_D) \quad (18.170)$$

This is called a **generalized additive model** or **GAM**.

Figure 18.20 shows an example of this, where each base kernel has the form $\mathcal{K}_d(x_d, x'_d) = \sigma_d^2 \text{SE}(x_d, x'_d | \ell_d)$. In Figure 18.20, we see that the σ_d^2 terms for the coarse and fine features are set to zero, indicating that these inputs have no impact on the response variable.

[DBW20] considers additive kernels operating on different linear projections of the inputs:

$$\mathcal{K}(\mathbf{x}, \mathbf{x}') = \sum_{b=1}^B w_b \mathcal{K}_b(\mathbf{P}_b \mathbf{x}, \mathbf{P}_b \mathbf{x}') \quad (18.171)$$

Surprisingly, they show that these models can match or exceed the performance of kernels operating on the original space, even when the projections are into a **single** dimension, and not learned. In other words, it is possible to reduce many regression problems to a single dimension without loss in performance. This finding is particularly promising for scalable inference, such as KISS (see Section 18.5.5.3), and active learning, which are greatly simplified in a low dimensional setting.

More recently, [LBH22] has proposed the **orthogonal additive kernel** (OAK), which imposes an orthogonality constraint on the additive functions. This ensures an identifiable, low-dimensional representation of the functional relationship, and results in improved performance.

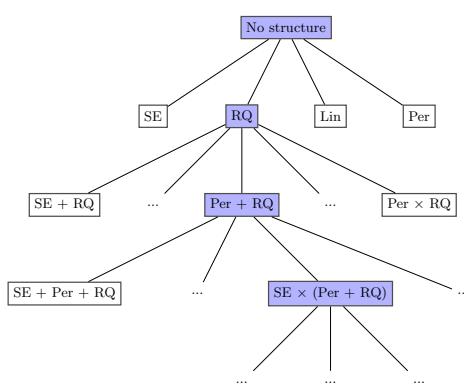


Figure 18.21: Example of a search tree over kernel expressions. Adapted from Figure 3.2 of [Duv14].

18.6.4 Automatic search for compositional kernels

Although the above methods can estimate the hyperparameters of a specified set of kernels, they do not choose the kernels themselves (other than the special case of selecting a subset of kernels from a set). In this section, we describe a method, based on [Duv+13], for sequentially searching through the space of increasingly complex GP models so as to find a parsimonious description of the data.

We start with a simple kernel, such as the white noise kernel, and then consider replacing it with a set of possible alternative kernels, such as an SE kernel, RQ kernel, etc. We use the BIC score (Section 3.9.5.2) to evaluate each candidate model (choice of kernel) m . This has the form $BIC(m) = \log p(\mathcal{D}|m) - \frac{1}{2}|m|\log N$, where $p(\mathcal{D}|m)$ is the marginal likelihood, and $|m|$ is the number of parameters. The first term measures fit to the data, and the second term is a complexity penalty. We can also consider replacing a kernel by the addition of two kernels, $k \rightarrow (k + k')$, or the multiplication of two kernels, $k \rightarrow (k \times k')$. See Figure 18.21 for an illustration of the search space.

Searching through this space is similar to what a human expert would do. In particular, if we find structure in the residuals, such as periodicity, we can propose a certain “move” through the space. We can also start with some structure that is assumed to hold globally, such as linearity, but if we find this only holds locally, we can multiply the kernel by an SE kernel. We can also add input dimensions incrementally, to capture higher order interactions.

Figure 18.22 shows the output of this process applied to a dataset of monthly totals of international airline passengers. The input to the GP is the set of time stamps, $\mathbf{x} = 1 : t$; there are no other features.

The observed data lies in between the dotted vertical lines; curves outside of this region are extrapolations. We see that the system has discovered a fairly interpretable set of patterns in the data. Indeed, it is possible to devise an algorithm to automatically convert the output of this search process to a natural language summary, as shown in [Llo+14]. In this example, it summarizes the data as being generated by the addition of 4 underlying trends: a linearly increasing function; an approximately periodic function with a period of 1.0 years, and with linearly increasing amplitude; a smooth function; and uncorrelated noise with linearly increasing standard deviation.

Recently, [Sun+18] showed how to create a DNN which learns the kernel given two input vectors.

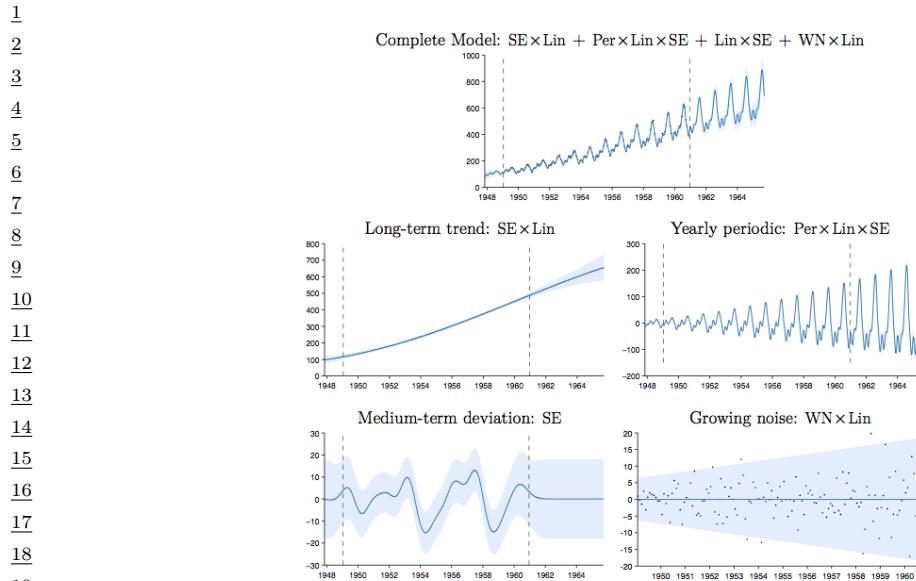


Figure 18.22: Top row: airline dataset and posterior distribution of the model discovered after a search of depth 10. Subsequent rows: predictions of the individual components. From Figure 3.5 of [Duv14], based on [Llo+14]. Used with kind permission of David Duvenaud.

The hidden units are defined as sums and products of elementary kernels, as in the above search based approach. However, the DNN can be trained in a differentiable way, so is much faster.

18.6.5 Spectral mixture kernel learning

Any shift-invariant (stationary) kernel can be converted via the Fourier transform to its dual form, known as its **spectral density**. This means that learning the spectral density is equivalent to learning any shift-invariant kernel. For example, if we take the Fourier transform of an RBF kernel, we get a Gaussian spectral density centered at the origin. If we take the Fourier transform of a Matern kernel, we get a Student- t spectral density centred at the origin. Thus standard approaches to multiple kernel learning, which typically involve additive compositions of RBF and Matern kernels with different length-scale parameters, amount to density estimation with a scale mixture of Gaussian or Student- t distributions at the origin. Such models are very inflexible for density estimation, and thus also very limited in being able to perform kernel learning.

On the other hand, *scale-location* mixture of Gaussians can model any density to arbitrary precision. Moreover, with even a small number of components these mixtures of Gaussians are highly flexible. Thus a spectral density corresponding to a scale-location mixture of Gaussians forms an expressive basis for all shift-invariant kernels. One can evaluate the inverse Fourier transform for a Gaussian mixture analytically, to derive the **spectral mixture kernel** [WA13], which we can express for

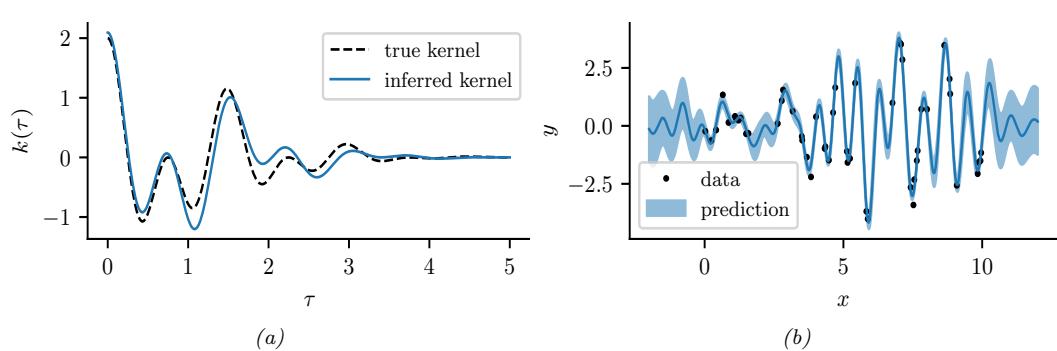


Figure 18.23: Illustration of a GP with a spectral mixture kernel in 1d. (a) Learned vs true kernel. (b) Predictions using learned kernel. Generated by [gp_spectral_mixture.ipynb](#).

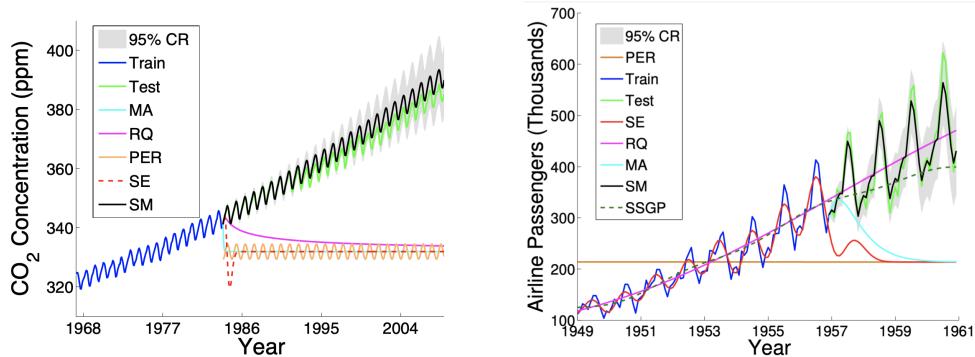


Figure 18.24: Extrapolations (point predictions and 95% credible set) on CO₂ and airline datasets using Gaussian processes with Matern, rational quadratic, periodic, RBF (SE), and spectral mixture kernels, each with hyperparameters learned using empirical Bayes. From [Wil14].

one-dimensional inputs x as:

$$\mathcal{K}(x, x') = \sum_i w_i \cos((x - x')(2\pi\mu_i)) \exp(-2\pi^2(x - x')^2 v_i) \quad (18.172)$$

The mixture weights w_i , as well as the means μ_i and variances v_i of the Gaussians in the spectral density, can be learned by empirical Bayes optimization (Section 18.6.1) or in a fully-Bayesian procedure (Section 18.6.2) [Jan+17]. We illustrate the former approach in Figure 18.23.

By learning the parameters of the spectral mixture kernel, we can discover representations that enable extrapolation — to make reasonable predictions far away from the data. For example, in Section 18.8.1, compositions of kernels are carefully hand-crafted to extrapolate CO₂ concentrations. But in this instance, the human statistician is doing all of the interesting representation learning. Figure 18.24 shows Gaussian processes with learned spectral mixture kernels instead automatically extrapolating on CO₂ and airline passenger problems.

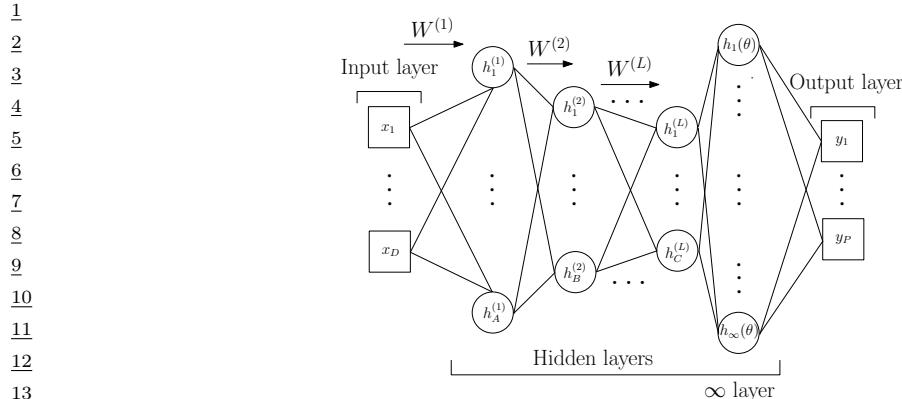


Figure 18.25: Deep Kernel Learning: A Gaussian process with a deep kernel maps D dimensional inputs \mathbf{x} through L parametric hidden layers followed by a hidden layer with an infinite number of basis functions, with base kernel hyperparameters $\boldsymbol{\theta}$. Overall, a Gaussian process with a deep kernel produces a probabilistic mapping with an infinite number of adaptive basis functions parametrized by $\gamma = \{\mathbf{w}, \boldsymbol{\theta}\}$. All parameters γ are learned through the marginal likelihood of the Gaussian process. From Figure 1 of [Wil+16].

19

20

21 These kernels can also be used to extrapolate higher dimensional large-scale spatio-temporal
22 patterns. Large datasets can provide relatively more information for expressive kernel learning.
23 However, scaling an expressive kernel learning approach poses different challenges than scaling a
24 standard Gaussian process model. One faces additional computational constraints, and the need
25 to retain significant model structure for expressing the rich information available in a large dataset.
26 Indeed, in Figure 18.24 we can separately understand the effects of the kernel learning approach and
27 scalable inference procedure, in being able to discover structure necessary to extrapolate textures.
28 An expressive kernel model and a scalable inference approach that preserves a *non-parametric*
29 representation are needed for good performance.

30 Structure exploiting inference procedures, such as Kronecker methods, as well as KISS-GP and
31 conjugate gradient based approaches, are appropriate for these tasks — since they generally preserve
32 or exploit existing structure, rather than introducing approximations that corrupt the structure.
33 Spectral mixture kernels combined with these scalable inference techniques have been used to great
34 effect for spatiotemporal extrapolation problems, including land-surface temperature forecasting,
35 epidemiological modeling, and policy-relevant applications.

36

37 18.6.6 Deep kernel learning

38

39 **Deep kernel learning** [SH07; Wil+16] combines the structural properties of neural networks with
40 the non-parametric flexibility and uncertainty representation provided by Gaussian processes. For
41 example, we can define a “deep RBF kernel” as follows:

$$42 \quad \mathcal{K}_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{x}') = \exp \left[-\frac{1}{2\sigma^2} \|\mathbf{h}_{\boldsymbol{\theta}}^L(\mathbf{x}) - \mathbf{h}_{\boldsymbol{\theta}}^L(\mathbf{x}')\|^2 \right] \quad (18.173)$$

43 where $\mathbf{h}_{\boldsymbol{\theta}}^L(\mathbf{x})$ are the outputs of layer L from a DNN. We can then learn the parameters $\boldsymbol{\theta}$ by
44 maximizing the marginal likelihood of the Gaussian processes.

45

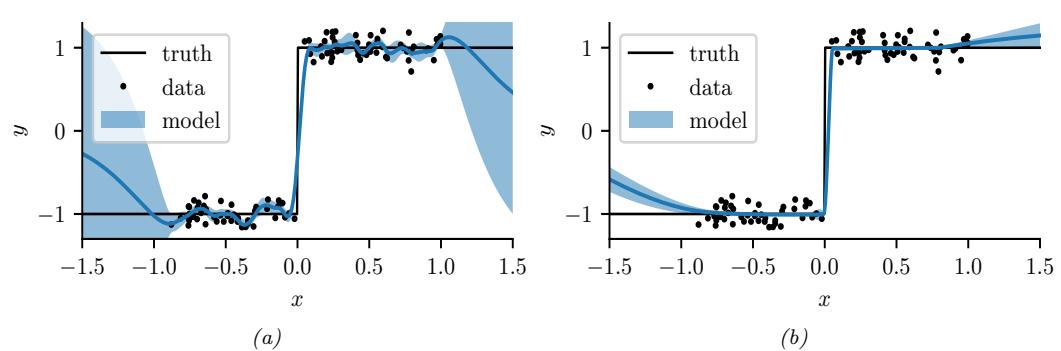


Figure 18.26: Modeling a discontinuous function with (a) a GP with a “shallow” Matern $\frac{3}{2}$ kernel, and (b) a GP with a “deep” MLP + Matern kernel. Generated by [gp_deep_kernel_learning.ipynb](#).

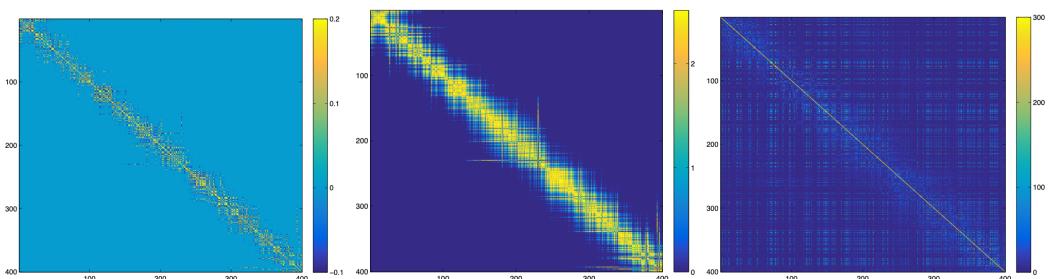


Figure 18.27: Left: The learned covariance matrix of a deep kernel with spectral mixture base kernel on a set of test cases for the Olivetti faces dataset, where the test samples are ordered according to the orientations of the input faces. **Middle:** The respective covariance matrix using a deep kernel with RBF base kernel. **Right:** The respective covariance matrix using a standard RBF kernel. From Figure 5 of [Wil+16].

This framework is illustrated in Figure 18.25. We can understand the neural network features as inputs into a base kernel. The neural network can either be (i) pre-trained, (ii) learned jointly with the base kernel parameters, or (iii) pre-trained and then fine-tuned through the marginal likelihood. This approach can be viewed as a “last-layer” Bayesian model, where a Gaussian process is applied to the final layer of a neural network. The base kernel often provides a good measure of distance in feature space, desirably encouraging predictions to have high uncertainty as we move far away from the data.

We can use deep kernel learning to help the GP learn discontinuous functions, as illustrated in Figure 18.26. On the left we show the results of a GP with a standard Matern $\frac{3}{2}$ kernel. It is clear that the out-of-sample predictions are poor. On the right we show the results of the same model where we first transform the input through a learned 2 layer MLP (with 15 and 10 hidden units). It is clear that the model is working much better.

As a more complex example, we consider a regression problem where we wish to map faces (vectors of pixel intensities) to a continuous valued orientation angle. In Figure 18.27, we evaluate the deep kernel matrix (with RBF and spectral mixture base kernels, discussed in Section 18.6.5) on data

ordered by orientation angle. We can see that the learned deep kernels, in the left two panels, have a pronounced diagonal band, meaning that they have *discovered* that faces with similar orientation angles are correlated. On the other hand, in the right panel we see that the entries even for a learned RBF kernel are highly diffuse. Since the RBF kernel essentially uses Euclidean distance as a metric for similarity, it is unable to learn a representation that effectively solves this problem. In this case, one must do highly non-Euclidean metric learning.

However, [ORW21] show that the approach to DKL based on maximizing the marginal likelihood can result in overfitting that is worse than standard DNN learning. They propose a fully Bayesian approach, in which they use SGLD (Section 12.7.1) to sample the DNN weights as well as the GP hyperparameters.

12

13 18.7 GPs and DNNs

14

In Section 18.6.6, we showed how we can combine the structural properties of neural networks with GPs. In Section 18.7.1 we show that, in the limit of infinitely wide networks, a neural network defines a GP with a certain kernel. These kernels are fixed, so the method is not performing representation learning, as a standard neural network would (see e.g., [COB18; Woo+19]). Nonetheless, these kernels are interesting in their own right, for example in modelling non-stationary covariance structure. In Section 18.7.2, we discuss the connection between SGD training of DNNs and GPs. And in Section 18.7.3, we discuss deep GPs, which are similar to DNNs in that they consist of many layers of functions which are composed together, but each layer is a nonparametric function.

23

24 18.7.1 Kernels derived from infinitely wide DNNs (NN-GP)

25

In this section, we show that an MLP with one hidden layer, whose width goes to infinity, and which has a Gaussian prior on all the parameters, converges to a Gaussian process with a well-defined kernel.⁴ This result was first shown for in [Nea96; Wil98], and was later extended to deep MLPs in [DFS16; Lee+18], to CNNs in [Nov+19], and to general DNNs in [Yan19]. The resulting kernel is called the NN-GP kernel [Lee+18].

We will consider the following model:

$$32 \quad f_k(\mathbf{x}) = b_k + \sum_{j=1}^H v_{jk} h_j(\mathbf{x}), \quad h_j(\mathbf{x}) = \varphi(u_{0j} + \mathbf{x}^\top \mathbf{u}_j) \quad (18.174)$$

where H is the number of hidden units, and $\varphi()$ is some nonlinear activation function, such as ReLU. We will assume Gaussian priors on the parameters:

$$38 \quad b_k \sim \mathcal{N}(0, \sigma_b), v_{jk} \sim \mathcal{N}(0, \sigma_v), u_{0j} \sim \mathcal{N}(0, \sigma_0), \mathbf{u}_j \sim \mathcal{N}(0, \Sigma) \quad (18.175)$$

Let $\boldsymbol{\theta} = \{b_k, v_{jk}, u_{0j}, \mathbf{u}_j\}$ be all the parameters. The expected output from unit k when applied to one input vector is given by

$$41 \quad \mathbb{E}_{\boldsymbol{\theta}} [f_k(\mathbf{x})] = \mathbb{E}_{\boldsymbol{\theta}} \left[b_k + \sum_{j=1}^H v_{jk} h_j(\mathbf{x}) \right] = \underbrace{\mathbb{E}_{\boldsymbol{\theta}} [b_k]}_{=0} + \sum_{j=1}^H \underbrace{\mathbb{E}_{\boldsymbol{\theta}} [v_{jk}]}_{=0} \mathbb{E}_{\boldsymbol{\theta}} [h_j(\mathbf{x})] = 0 \quad (18.176)$$

45 4. Our presentation is based on http://cbl.eng.cam.ac.uk/pub/Intranet/MLG/ReadingGroup/presentation_matthias.pdf.

47

The covariance in the output for unit k when the function is applied to two different inputs is given by the following.⁵

$$\mathbb{E}_{\boldsymbol{\theta}} [f_k(\mathbf{x})f_k(\mathbf{x}')] = \mathbb{E}_{\boldsymbol{\theta}} \left[\left(b_k + \sum_{j=1}^H v_{jk} h_j(\mathbf{x}) \right) \left(b_k + \sum_{j=1}^H v_{jk} h_j(\mathbf{x}') \right) \right] \quad (18.177)$$

$$= \sigma_b^2 + \sum_{j=1}^H \mathbb{E}_{\boldsymbol{\theta}} [v_{jk}^2] \mathbb{E}_{\boldsymbol{u}} [h_j(\mathbf{x})h_j(\mathbf{x}')] = \sigma_b^2 + \sigma_v^2 H \mathbb{E}_{\boldsymbol{u}} [h_j(\mathbf{x})h_j(\mathbf{x}')] \quad (18.178)$$

Now consider the limit $H \rightarrow \infty$. We scale the magnitude of the output by defining $\sigma_v^2 = \omega/H$. Since the input to k 'th output unit is an infinite sum of random variables (from the hidden units $h_j(\mathbf{x})$), we can use the **central limit theorem** to conclude that the output converges to a Gaussian with mean and variance given by

$$\mathbb{E}[f_k(\mathbf{x})] = 0, \quad \mathbb{V}[f_k(\mathbf{x})] = \sigma_b^2 + \omega \mathbb{E}_{\boldsymbol{u}} [h(\mathbf{x})^2] \quad (18.179)$$

Furthermore, the joint distribution over $\{f_k(\mathbf{x}_n) : n = 1 : N\}$ for any $N \geq 2$ converges to a multivariate Gaussian with covariance given by

$$\mathbb{E}[f_k(\mathbf{x})f_k(\mathbf{x}')] = \sigma_b^2 + \omega \mathbb{E}_{\boldsymbol{u}} [h(\mathbf{x})h(\mathbf{x}')] \triangleq \mathcal{K}(\mathbf{x}, \mathbf{x}') \quad (18.180)$$

Thus the MLP converges to a GP. To compute the kernel function, we need to evaluate

$$C(\mathbf{x}, \mathbf{x}') = \mathbb{E}_{\boldsymbol{u}} [h(u_0 + \mathbf{u}^\top \mathbf{x})h(u_0 + \mathbf{u}^\top \mathbf{x}')] = \mathbb{E}_{\boldsymbol{u}} [h(\tilde{\mathbf{u}}^\top \tilde{\mathbf{x}})h(\tilde{\mathbf{u}}^\top \tilde{\mathbf{x}}')] \quad (18.181)$$

where we have defined $\tilde{\mathbf{x}} = (1, \mathbf{x})$ and $\tilde{\mathbf{u}} = (u_0, \mathbf{u})$. Let us define

$$\tilde{\Sigma} = \begin{pmatrix} \sigma_0^2 & 0 \\ 0 & \Sigma \end{pmatrix} \quad (18.182)$$

Then we have

$$C(\mathbf{x}, \mathbf{x}') = \int h(\tilde{\mathbf{u}}^\top \tilde{\mathbf{x}})h(\tilde{\mathbf{u}}^\top \tilde{\mathbf{x}}') \mathcal{N}(\tilde{\mathbf{u}} | \mathbf{0}, \tilde{\Sigma}) d\tilde{\mathbf{u}} \quad (18.183)$$

This can be computed in closed form for certain activation functions, as shown in Table 18.4.

This is sometimes called the **neural net kernel**. Note that this is a **nonstationary kernel** (i.e., not a function of $\|\mathbf{x} - \mathbf{x}'\|$), and sample paths from it are nearly discontinuous and tend to constant values for large positive or negative inputs, as illustrated in Figure 18.28.

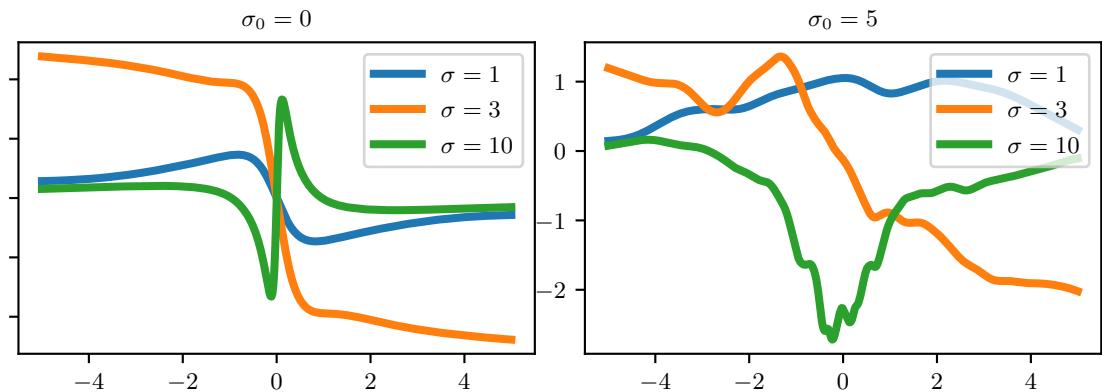
18.7.2 Neural tangent kernel (NTK)

In Section 18.7.1 we derived the NN-GP kernel, under the assumption that all the weights are random. A natural question is: can we derive a kernel from a DNN after it has been trained, or more generally, while it is being trained. It turns out that this can be done, as we show below.

⁵ We are using the fact that $u \sim \mathcal{N}(0, \sigma^2)$ implies $\mathbb{E}[u^2] = \mathbb{V}[u] = \sigma^2$.

$\frac{1}{2}$	$h(\tilde{\mathbf{x}})$	$C(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}')$
$\frac{3}{2}$	$\text{erf}(\tilde{\mathbf{x}}^\top \tilde{\mathbf{u}})$	$\frac{2}{\pi} \arcsin(f_1(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}'))$
$\frac{4}{2}$	$\mathbb{I}(\tilde{\mathbf{x}}^\top \tilde{\mathbf{u}} \geq 0)$	$\pi - \theta(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}')$
$\frac{5}{2}$	$\text{ReLU}(\tilde{\mathbf{x}}^\top \tilde{\mathbf{u}})$	$\frac{f_2(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}')}{\pi} \sin(\theta(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}')) + \frac{\pi - \theta(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}')}{\pi} \tilde{\mathbf{x}}^\top \tilde{\Sigma} \tilde{\mathbf{x}}$

$\frac{7}{8}$ Table 18.4: Some neural net GP kernels. Here we define $f_1(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}') = \frac{2\tilde{\mathbf{x}}^\top \tilde{\Sigma} \tilde{\mathbf{x}}'}{\sqrt{(1+2\tilde{\mathbf{x}}^\top \tilde{\Sigma} \tilde{\mathbf{x}})(1+2(\tilde{\mathbf{x}}')^\top \tilde{\Sigma} \tilde{\mathbf{x}}')}}$, $f_2(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}') =$
 $\frac{1}{9} \|\tilde{\Sigma}^{\frac{1}{2}} \tilde{\mathbf{x}}\| \|\tilde{\Sigma}^{\frac{1}{2}} \tilde{\mathbf{x}}'\|$, $f_3(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}') = \sqrt{(\tilde{\mathbf{x}}^\top \tilde{\Sigma} \tilde{\mathbf{x}})((\tilde{\mathbf{x}}')^\top \tilde{\Sigma} \tilde{\mathbf{x}}')}$, and $\theta(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}') = \arccos(f_3(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}'))$. Results are derived in
 $\frac{10}{11}$ [Wil98; CS09].



$\frac{24}{25}$ Figure 18.28: Sample output from a GP with an NNGP kernel derived from an infinitely wide one layer MLP with activation function of the form $h(x) = \text{erf}(x \cdot u + u_0)$ where $u \sim \mathcal{N}(0, \sigma)$ and $u_0 \sim \mathcal{N}(0, \sigma_0)$. Generated by [nngp_1d.ipynb](#). Used with kind permission of Matthias Bauer.

$\frac{29}{30}$ Let $\mathbf{f} = [f(\mathbf{x}_n; \theta)]_{n=1}^N$ be the $N \times 1$ prediction vector, let $\nabla_{\mathbf{f}} \mathcal{L} = [\frac{\partial \mathcal{L}}{\partial f(\mathbf{x}_n)}]_{n=1}^N$ be the $N \times 1$ loss gradient vector, let $\boldsymbol{\theta} = [\theta_p]_{p=1}^P$ be the $P \times 1$ vector of parameters, and let $\nabla_{\boldsymbol{\theta}} \mathbf{f} = [\frac{\partial f(\mathbf{x}_n)}{\partial \theta_p}]$ be the $P \times N$ matrix of partials. Suppose we perform continuous time gradient descent with fixed learning rate η . The parameters evolve over time as follows:

$$\frac{34}{35} \quad \partial_t \boldsymbol{\theta}_t = -\eta \nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{f}_t) = -\eta \nabla_{\boldsymbol{\theta}} \mathbf{f}_t \cdot \nabla_{\mathbf{f}} \mathcal{L}(\mathbf{f}_t) \quad (18.184)$$

$\frac{36}{37}$ Thus the function evolves over time as follows:

$$\frac{38}{39} \quad \partial_t \mathbf{f}_t = \nabla_{\boldsymbol{\theta}} \mathbf{f}_t^\top \partial_t \boldsymbol{\theta}_t = -\eta \nabla_{\boldsymbol{\theta}} \mathbf{f}_t^\top \nabla_{\boldsymbol{\theta}} \mathbf{f}_t \cdot \nabla_{\mathbf{f}} \mathcal{L}(\mathbf{f}_t) = -\eta \mathcal{T}_t \cdot \nabla_{\mathbf{f}} \mathcal{L}(\mathbf{f}_t) \quad (18.185)$$

$\frac{40}{41}$ where \mathcal{T}_t is the $N \times N$ kernel matrix

$$\frac{42}{43} \quad \mathcal{T}_t(\mathbf{x}, \mathbf{x}') \triangleq \nabla_{\boldsymbol{\theta}} f_t(\mathbf{x}) \cdot \nabla_{\boldsymbol{\theta}} f_t(\mathbf{x}') = \sum_{p=1}^P \frac{\partial f(\mathbf{x}; \boldsymbol{\theta})}{\partial \theta_p} \Big|_{\boldsymbol{\theta}_t} \frac{\partial f(\mathbf{x}'; \boldsymbol{\theta})}{\partial \theta_p} \Big|_{\boldsymbol{\theta}_t} \quad (18.186)$$

$\frac{45}{46}$ If we let the learning rate η become infinitesimally small, and the widths go to infinity, one can show that this kernel converges to a constant matrix, this is known as the **neural tangent kernel** or

1 NTK [JGH18]:

2

$$\mathcal{T}(\mathbf{x}, \mathbf{x}') \triangleq \nabla_{\boldsymbol{\theta}} f(\mathbf{x}; \boldsymbol{\theta}_{\infty}) \cdot \nabla_{\boldsymbol{\theta}} f(\mathbf{x}'; \boldsymbol{\theta}_{\infty}) \quad (18.187)$$

3

4 Details on how to compute this kernel for various models, such as CNNs, graph neural nets, and
5 general neural nets, can be found in [Aro+19; Du+19; Yan19]. A software library to compute the
6 NN-GP kernel and NTK is available in [Ano19].

7 The assumptions behind the NTK results in the parameters barely changing from their initial
8 values (which is why a linear approximation around the starting parameters is valid). This can
9 still lead to a change in the final predictions (and zero final training error), because the final layer
10 weights can learn to use the random features just like in kernel regression. However, this phenomenon
11 — which has been called “**lazy training**” [COB18] — is not representative of DNN behavior in
12 practice [Woo+19], where parameters often change a lot. Fortunately it is possible to use a different
13 parameterization which does result in feature learning in the infinite width limit [YH21].

14 dev

15 18.7.3 Deep GPs

16 A **deep Gaussian process** or **DGP** is a composition of GPs [DL13]. More formally, a DGP of L
17 layers is a hierarchical model of the form

18

$$\text{DGP}(\mathbf{x}) = f_L \circ \cdots \circ f_1(\mathbf{x}), \quad f_i(\cdot) = [f_i^{(1)}(\cdot), \dots, f_i^{(H_i)}(\cdot)], \quad f_i^{(j)} \sim \text{GP}(0, \mathcal{K}_i(\cdot, \cdot)) \quad (18.188)$$

19

20 This is similar to a deep neural network, except the hidden nodes are now hidden functions.

21 A natural question is: what is gained by this approach compared to a standard GP? Although
22 conventional single-layer GPs are nonparametric, and can model any function (assuming the use of a
23 non-degenerate kernel) with enough data, in practice their performance is limited by the choice of
24 kernel. It is tempting to think that deep kernel learning (Section 18.6.6) can solve this problem, but
25 in theory a GP on top of a DNN is still just a GP. However, one can show that a composition of GPs
26 is strictly more general. Unfortunately, inference in deep GPs is rather complicated. so we leave the
27 details to [Supplementary](#) Section 18.1. See also [Jak21] for a recent survey on this topic.

28 18.8 Gaussian processes for timeseries forecasting

29 It is possible to use Gaussian processes to perform timeseries forecasting (see e.g., [Rob+13]). The
30 basic idea is to model the unknown output as a function of time, $f(t)$, and to represent a prior about
31 the form of f as a GP; we then update this prior given the observed evidence, and forecast into the
32 future. Naively this would take $O(T^3)$ time. However, for certain stationary kernels, it is possible to
33 reformulate the problem as a linear-Gaussian state space model, and then use the Kalman smoother
34 to perform inference in $O(T)$ time, as explained in [SSH13; SS19; Ada+20]. This conversion can be
35 done exactly for Matern kernels and approximately for Gaussian (RBF) kernels (see [SS19, Ch. 12]).
36 In [SGF21], they describe how to reduce the linear dependence on T to $\log(T)$ time using a parallel
37 prefix scan operator, that can be run efficiently on GPUs (see Section 8.3.3.4).

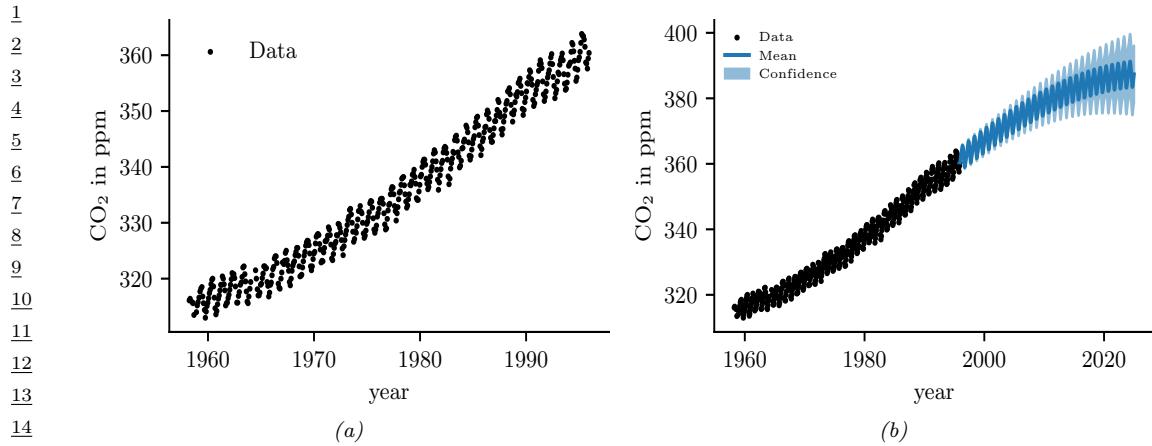


Figure 18.29: (a) The observed Mauna Loa CO₂ time series. (b) Forecasts from a GP. Generated by [gp_mauna_loa.ipynb](#).

18.8.1 Example: Mauna Loa

In this section, we use the Mauna Loa CO₂ dataset from Section 29.12.3. We show the raw data in Figure 18.29(a). We see that there is periodic (or quasi-periodic) signal with a year-long period superimposed on a long term trend. Following [RW06, Sec 5.4.3], we will model this with a composition of kernels:

$$\mathcal{K}(r) = \mathcal{K}_1(r) + \mathcal{K}_2(r) + \mathcal{K}_3(r) + \mathcal{K}_4(r) \quad (18.189)$$

where $\mathcal{K}_i(t, t') = \mathcal{K}_i(t - t')$ for the i 'th kernel.

To capture the long term smooth rising trend, we let \mathcal{K}_1 be a squared exponential (SE) kernel, where θ_0 is the amplitude and θ_1 is the length scale:

$$\mathcal{K}_1(r) = \theta_0^2 \exp\left(-\frac{r^2}{2\theta_1^2}\right) \quad (18.190)$$

To model the periodicity, we can use periodic or exp-sine-squared kernel from Equation (18.19) with a period of 1 year. However, since it is not clear if the seasonal trend is exactly periodic, we multiply this periodic kernel with another SE kernel to allow for a decay away from periodicity; the result is \mathcal{K}_2 , where θ_2 is the magnitude, θ_3 is the decay time for the periodic component, $\theta_4 = 1$ is the period, and θ_5 is the smoothness of the periodic component.

$$\mathcal{K}_2(r) = \theta_2^2 \exp\left(-\frac{r^2}{2\theta_3^2} - \theta_5 \sin^2\left(\frac{\pi r}{\theta_4}\right)\right) \quad (18.191)$$

To model the (small) medium term irregularities, we use a rational quadratic kernel (Equation (18.21)):

$$\mathcal{K}_3(r) = \theta_6^2 \left[1 + \frac{r^2}{2\theta_7^2 \theta_8}\right]^{-\theta_8} \quad (18.192)$$

1 where θ_6 is the magnitude, θ_7 is the typical length scale, and θ_8 is the shape parameter.
2

3 The magnitude of the independent noise can be incorporated into the observation noise of the
4 likelihood function. For the correlated noise, we use another SE kernel:
5

6
$$\mathcal{K}_4(r) = \theta_9^2 \exp\left(-\frac{r^2}{2\theta_{10}^2}\right) \quad (18.193)$$

7

8 where θ_9 is the magnitude of the correlated noise, and θ_{10} is the length scale. (Note that the
9 combination of \mathcal{K}_1 and \mathcal{K}_4 is non-identifiable, but this does not affect predictions.)
10

11 We can fit this model by optimizing the marginal likelihood wrt $\boldsymbol{\theta}$ (see Section 18.6.1). The
12 resulting forecast is shown in Figure 18.29(b).
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

19 Beyond the iid assumption

19.1 Introduction

The standard approach to supervised ML assumes the training and test sets both contain independent and identically distributed (iid) samples from the same distribution. However, there are many settings in which the test distribution may be different from the training distribution; this is known as **distribution shift**, as we discuss in Section 19.2.

In some cases, we may have data from multiple related distributions, not just train and test, as we discuss in Section 19.6. We may also encounter data in a streaming setting, where the data distribution may be changing continuously, or in a piecewise constant fashion, as we discuss in Section 19.7. Finally, in Section 19.8, we discuss settings in which the test distribution is chosen by an adversary to minimize performance of a prediction system.

19.2 Distribution shift

Suppose we have a labeled training set from a **source distribution** $p(\mathbf{x}, \mathbf{y})$ which we use to fit a predictive model $p(\mathbf{y}|\mathbf{x})$. At test time we encounter data from the **target distribution** $q(\mathbf{x}, \mathbf{y})$. If $p \neq q$, we say that there has been a **distribution shift** or **dataset shift** [QC+08; BD+10]. This can adversely affect the performance of predictive models, as we illustrate in Section 19.2.1. In Section 19.2.2 we give a taxonomy of some kinds of distribution shift using the language of causal graphical models. We then proceed to discuss a variety of strategies that can be adopted to ameliorate the harm caused by distribution shift. In particular, in Section 19.3, we discuss techniques for detecting shifts, so that we can abstain from giving an incorrect prediction if the model is not confident. In Section 19.4, we discuss techniques to improve robustness to shifts; in particular, given labeled data from $p(\mathbf{x}, \mathbf{y})$, we aim to create a model that approximates $q(\mathbf{y}|\mathbf{x})$. In Section 19.5, we discuss techniques to adapt the model to the target distribution given some labeled or unlabeled data from the target.

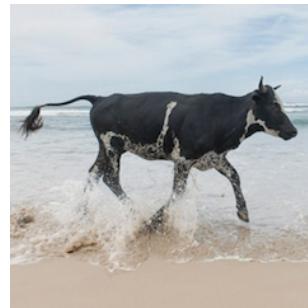
19.2.1 Motivating examples

Figure 19.1 shows how shifting the test distribution slightly, by adding a small amount of Gaussian noise, can hurt performance of an otherwise high accuracy image classifier. Similar effects occur with other kinds of **common corruptions**, such as image blurring [HD19]. Analogous problems can also occur in the text domain [Ryc+19], and the speech domain (c.f. male vs female speakers in

10 Figure 19.1: Effect of Gaussian noise of increasing magnitude on an image classifier. The model is a
11 ResNet-50 CNN trained on ImageNet. From Figure 23 of [For+19]. Used with kind permission of Justin
12 Gilmer.



(A) Cow: 0.99, Pasture: 0.99,
Grass: 0.99, No Person: 0.98,
Mammal: 0.98



(B) No Person: 0.99, Water: 0.98,
Beach: 0.97, Outdoors: 0.97,
Seashore: 0.97



(C) No Person: 0.97, **Mammal:**
0.96, Water: 0.94, Beach: 0.94, Two:
0.94

27 Figure 19.2: Illustration of how image classifiers generalize poorly to new environments. (a) In the training
28 data, most cows occur on grassy backgrounds. (b-c) In these test image, the cow occurs “out of context”, namely
29 on a beach. The background is considered a “spurious correlation”. In (b), the cow is not detected. In (c), it is
30 classified with a generic “mammal” label. Top five labels and their confidences are produced by ClarifAI.com,
31 which is a state of the art commerical vision system. From Figure 1 of [BVHP18]. Used with kind permission
 of Sara Beery.

³⁴ Figure 34.1). These examples illustrate that high performing predictive models can be very sensitive to small changes in the input distribution.

36 Performance can also drop on “clean” images, but which exhibit other kinds of shift. Figure 19.2
 37 gives an amusing example of this. In particular, it illustrates how the performance of a CNN image
 38 classifier can be very accurate on **in-domain** data, but can be very inaccurate on **out-of-domain**
 39 data, such as images with a different background, or taken at a different time or location (see e.g.,
 40 [Koh+20b]) or from a novel viewing angle (see e.g., [KH22]).

41 The root cause of many of these problems is the fact that discriminative models often leverage
42 features that are predictive of the output *in the training set*, but which are not reliable in general.
43 For example, in an image classification dataset, we may find that green grass in the background
44 is very predictive of the class label “cow”, but this is not a feature that is stable across different
45 distributions; these are called **spurious correlations** or **shortcut features**. Unfortunately, such
46 features are often easier for models to learn, for reasons explained in [Gei+20a; Xia+21; Sha+20;
47

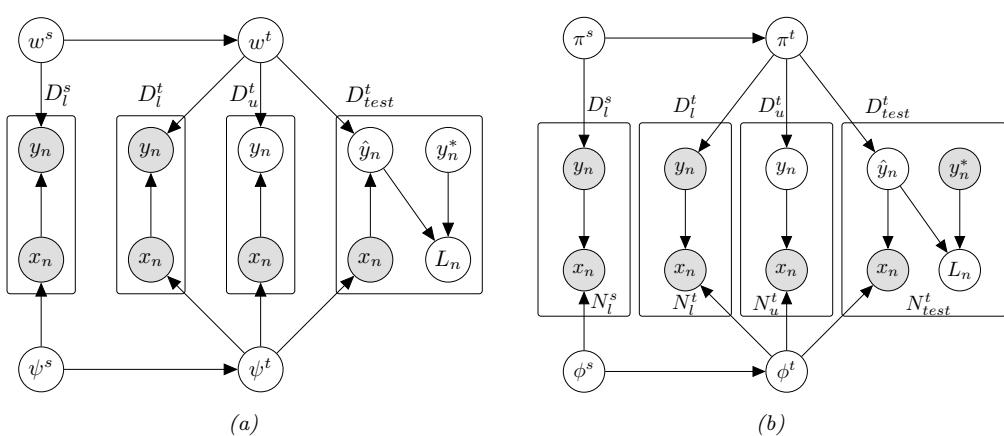


Figure 19.3: Models for distribution shift from source s to target t . Here \mathcal{D}_s^s is the labeled training set from the source, \mathcal{D}_L^t is an optional labeled training set from the target, \mathcal{D}_U^t is an optional unlabeled training set from the target, and \mathcal{D}_{test}^t is a labeled test set from the target. In the latter case, \hat{y}_n is the prediction on the n 'th test case (generated by the model), y_n^* is the true value, and $\ell_n = \ell(y_n^*, \hat{y}_n)$ is the corresponding loss. (Note that we don't evaluate the loss on the source distribution.) (a) Discriminative (causal) model. (b) Generative (anti-causal).

Pez+21].

Relying on these shortcuts can have serious real-world consequences. For example, [Zec+18a] found that a CNN trained to recognize pneumonia was relying on hospital-specific metal tokens in the chest X-ray scans, rather than focusing on the lungs themselves, and thus the model did not generalize to new hospitals.

Analogous problems arise with other kinds of ML models, as well as other data types, such as text (e.g., changing “he” to “she” can flip the output of a sentiment analysis system), audio (e.g., adding background noise can easily confuse speech recognition systems), and medical records [Ros22]. Furthermore, the changes to the input needed to change the output can often be imperceptible, as we discuss in the section on adversarial robustness (Section 19.8).

19.2.2 A causal view of distribution shift

In the sections below, we briefly summarize some canonical kinds of distribution shift. We adopt a causal view of the problem, following [Sch+12a; Zha+13b; BP16; Mei18a; CWG20; Bud+21; SCS22]).¹ (See Section 4.7 for a brief discussion of causal DAGs, and Chapter 36 for more details.)

We assume the inputs to the model (the covariates) are X and the outputs to be predicted (the labels) are Y . If we believe that X causes Y , denoted $X \rightarrow Y$, we call it **causal prediction**. If we believe that Y causes X , denoted $Y \rightarrow X$, we call it **anti-causal prediction** [Sch+12a].

The decision about which model to use depends on our assumptions about the underlying **data**.

1. In the causality literature, the question of whether a model can generalize to a new distribution is called the question of **external validity**. If a model is externally valid, we say that it is **transportable** from one distribution to another [BP16].

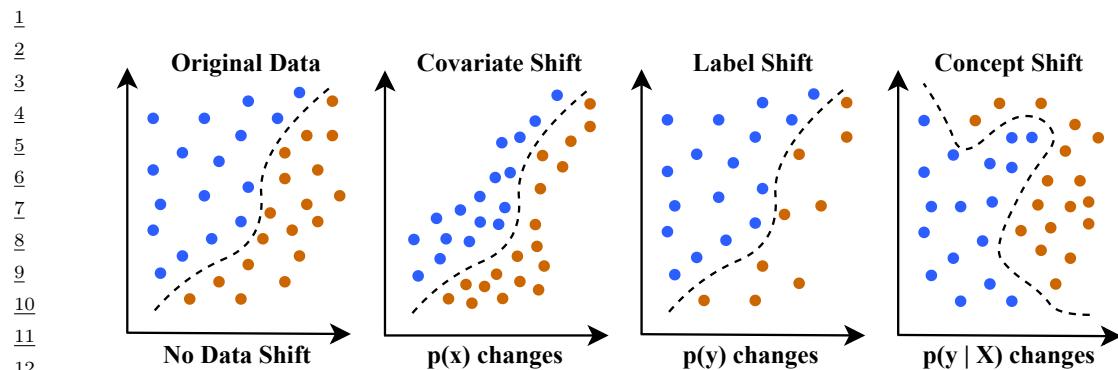


Figure 19.4: Illustration of the 4 main kinds of distribution shift for a 2d binary classification problem.
Adapted from Figure 1 of [al21].

generating process. For example, suppose X is a medical image, and Y is an image segmentation created by a human expert or an algorithm. If we change the image, we will change the annotation, and hence $X \rightarrow Y$. Now suppose X is a medical image and Y is the ground truth disease state of the patient, as estimated by some other means (e.g., a lab test). In this case, we have $Y \rightarrow X$, since changing the disease state will change the appearance of the image. As another example, suppose X is a text review of a movie, and Y is a measure of how informative the review is. Clearly we have $X \rightarrow Y$. Now suppose Y is the star rating of the movie, representing the degree to which the user liked it; this will affect the words that they write, and hence $Y \rightarrow X$.

Based on the above discussion, we can factor the joint distribution in two possible ways. One way is to define a discriminative (causal) model:

$$p_{\theta}(x, y) = p_{\psi}(x)p_w(y|x) \quad (19.1)$$

See Figure 19.3a. Alternatively we can define a generative (anti-causal) model:

$$p_{\theta}(x, y) = p_{\pi}(y)p_{\phi}(x|y) \quad (19.2)$$

See Figure 19.3b. For each of these 2 models model types, different parts of the distribution may change from source to target. This gives rise to 4 canonical type of shift, as we discuss in Section 19.2.3.

19.2.3 The four main types of distribution shift

The four main types of distribution shift are summarized in Section 19.2 and are illustrated in Figure 19.4. We give more details below.

41

19.2.3.1 Covariate shift

44 In a causal (discriminative) model, if $p_{\psi}(x)$ changes (so $\psi^s \neq \psi^t$), we call it **covariate shift**, also
45 called **domain shift**. For example, the training distribution may be clean images of coffee pots, and
46 the test distribution may be images of coffee pots with Gaussian noise, as shown in Figure 19.1; or the

47

Name	Source	Target	Causal
Covariate / domain shift	$p(X)p(Y X)$	$q(X)p(Y X)$	Causal
Concept shift	$p(X)p(Y X)$	$p(X)q(Y X)$	Causal
Label (prior) shift	$p(Y)p(X Y)$	$q(Y)p(X Y)$	Anti-causal
Manifestation shift	$p(Y)p(X Y)$	$p(Y)q(X Y)$	Anti-causal

Table 19.1: The 4 main types of distribution shift.

training distribution may be photos of objects in a catalog, with uncluttered white backgrounds, and the test distribution may be photos of the same kinds of objects collected “in the wild”; or the training data may be synthetically generated images, and the test distribution may be real images. Similar shifts can occur in the text domain; for example, the training distribution may be movie reviews written in English, and the test distribution may be translations of these reviews into Spanish.

Some standard strategies to combat covariate shift include importance weighting (Section 19.5.2) and domain adaptation (Section 19.5.3).

19.2.3.2 Concept shift

In a causal (discriminative) model, if $p_w(\mathbf{y}|\mathbf{x})$ changes (so $w^s \neq w^t$), we call it **concept shift**, also called **annotation shift**. For example, consider the medical imaging context: the conventions for annotating images might be different between the training distribution and test distribution. Another example of concept shift occurs when a new label can occur in the target distribution that was not part of the source distribution. This is related to open world recognition, discussed in Section 19.3.4.

Since concept shift is a change in what we “mean” by a label, it is impossible to fix this problem without seeing labeled examples from the target distribution, which defines each label by means of examples.

19.2.3.3 Label / prior shift

In an anti-causal (generative) model, if $p_\pi(\mathbf{y})$ changes (i.e., $\pi^s \neq \pi^t$), we call it **label shift**, also called **prior shift** or **prevalence shift**. For example, consider the medical imaging context, where $Y = 1$ if the patient has some disease and $Y = 0$ otherwise. If the training distribution is an urban hospital and the test distribution is a rural hospital, then the prevalence of the disease, represented by $p(Y = 1)$, might very well be different.

Some standard strategies to combat label shift are to reweight the output of a discriminative classifier using an estimate of the new label distribution, as we discuss in Section 19.5.4.

19.2.3.4 Manifestation shift

In an anti-causal (generative) model, if $p_\phi(\mathbf{x}|\mathbf{y})$ changes (i.e., $\phi^s \neq \phi^t$), we call **manifestation shift** [CWG20], or **Conditional shift** [Zha+13b]. This is, in some sense, the inverse of concept shift. For example, consider the medical imaging context: the way that the same disease Y manifests itself in the shape of a tumor X might be different. This is usually due to the presence of a hidden confounding factor that has changed between source and target (e.g., different age of the patients).

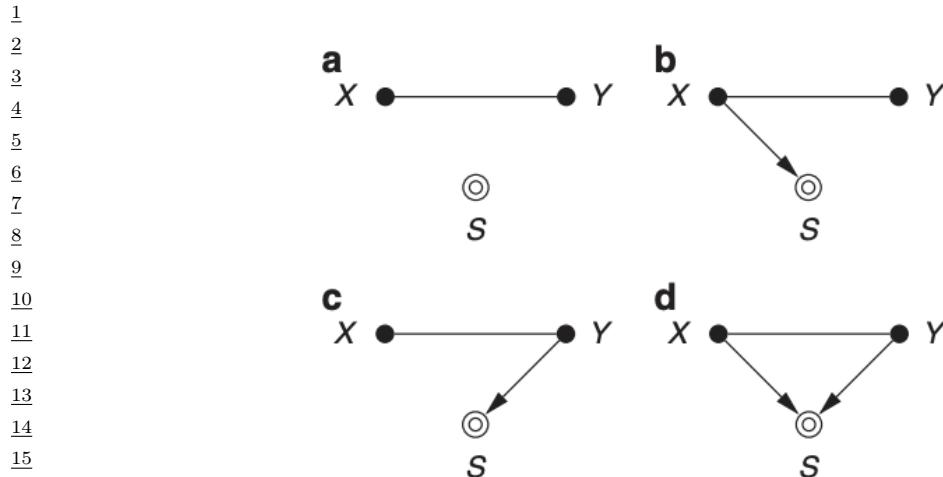


Figure 19.5: Causal diagrams for different sample selection strategies. Undirected edges can be oriented in either direction. The selection variable S is set to 1 if its parent nodes match the desired criterion; only these samples are included in the dataset. (a) No selection. (b) Selection on X . (c) Selection on Y . (d) Selection on X and Y . Adapted from Figure 4 of [CWG20].

19.2.4 Selection bias

In some cases, we may induce a shift in the distribution just due to the way the data is collected, without any changes to the underlying distributions. In particular, let $S = 1$ if a sample from the population is included in the training set, and $S = 0$ otherwise. Thus the source distribution is $p(X, Y) = p(X, Y|S = 1)$ but the target distribution is $q(X, Y) = p(X, Y|S \in \{0, 1\}) = p(X, Y)$, so there is no selection.

In Figure 19.5 we visualize the four kinds of selection. For example, suppose we select based on X meeting certain criteria, e.g., images of a certain quality, or exhibiting a certain pattern; this can induce domain shift or covariate shift. Now suppose we select based on Y meeting certain criteria, e.g., we are more likely to select rare examples where $Y = 1$, in order to **balance the dataset** (for reasons of computational efficiency); this can induce label shift. Finally, suppose we select based on both X and Y ; this can induce non-causal dependencies between X and Y , a phenomenon known as **selection bias** (see Section 4.2.4.2 for details).

19.3 Detecting distribution shifts

In general it will not be possible to make a model robust to all of the ways a distribution can shift at test time, nor will we always have access to test samples at training time. As an alternative, it may be sufficient for the model to *detect* that a shift has happened, and then to respond in the appropriate way. There are several ways of detecting distribution shift, some of which we summarize below. (See also Section 29.5.6, where we discuss changepoint detection in time series data.) The main distinction between methods is based on whether we have a set of samples from the target

1 distribution, or just a single sample, and whether the test samples are labeled or unlabeled. We
2 discuss these different scenarios below.

5 19.3.1 Detecting shifts using two-sample testing

7 Suppose we collect a set of samples from the source and target distribution. We can then use standard
8 techniques for **two-sample testing** to estimate if the null hypothesis, $p(\mathbf{x}, y) = q(\mathbf{x}, y)$, is true or
9 not. (If we have unlabeled samples, we just test if $p(\mathbf{x}) = q(\mathbf{x})$.) For example, we can use MMD
10 (Section 2.7.3) to measure the distance between the set of input samples (see e.g., [Liu+20a]). Or we
11 can measure (Euclidean) distances in the embedding space of a classifier trained on the source (see
12 e.g., [KM22]).

13 In some cases it may be possible to just test if the distribution of the labels $p(y)$ has changed, which
14 is an easier problem than testing for changes in the distribution of inputs $p(\mathbf{x})$. In particular, if the
15 label shift assumption (Section 19.2.3.3) holds (i.e., $q(\mathbf{x}|y) = p(\mathbf{x}|y)$), plus some other assumptions,
16 then we can use the blackbox shift estimation technique from Section 19.5.4 to estimate $q(y)$. If
17 we find that $q(y) = p(y)$, then we can conclude that $q(\mathbf{x}, y) = p(\mathbf{x}, y)$. In [RGL19], they showed
18 experimentally that this method worked well for detecting distribution shifts even when the label
19 shift assumption does not hold.

20 It is also possible to use conformal prediction (Section 14.3) to develop “distribution free” methods
21 for detecting covariate shift, given only access to a calibration set and some conformity scoring
22 function [HL20].

24 19.3.2 Detecting single out-of-distribution (OOD) inputs

26 Now suppose we just have *one* unlabeled sample from the target distribution, $\mathbf{x} \sim q$, and we want to
27 know if \mathbf{x} is in-distribution (**ID**) or out-of-distribution (**OOD**). We will call this problem **out-of-**
28 **distribution detection**, although it is also called **anomaly detection**, and **novelty detection**.²

29 The OOD detection problem requires making a binary decision about whether the test sample is ID
30 or OOD. If it is ID, we may optionally require that we return its class label, as shown in Figure 19.6.
31 In the sections below, we give a brief overview of techniques that have been proposed for tackling
32 this problem, but for more details, see e.g., [Pan+21; Ruf+21; Bul+20; Yan+21; Sal+21; Hen+19b].

34 19.3.2.1 Supervised ID/OOD methods (outlier exposure)

36 The simplest method for OOD detection assumes we have access to labeled ID and OOD samples at
37 training time. Then we just fit a binary classifier to distinguish the OOD or background class (called
38 “**known unknowns**”) from the ID class (called “**known knowns**”). This technique is called **outlier**
39 **exposure** (see e.g., [HMD19; Thu+21; Bit+21]) and can work well. However, in most cases we will
40 not have enough examples from the OOD distribution, since the OOD set is basically the set of all
41 possible inputs except for the ones of interest.

42 43 2. The task of **outlier detection** is somewhat different from anomaly or OOD detection, despite the similar name.
44 In the outlier detection literature, the assumption is that there is a single unlabeled dataset, and the goal is to identify
45 samples which are “untypical” compared to the majority. This is often used for **data cleaning**. (Note that this is a
46 **transductive learning** task, where the model is trained and evaluated on the same data. We focus on inductive
47 tasks, where we train a model on one dataset, and then test it on another.)

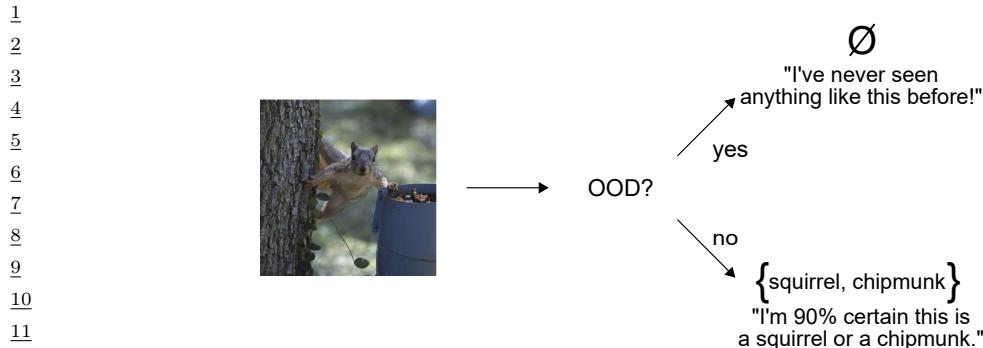


Figure 19.6: Illustration of a two-stage decision problem. First we must decide if the input image is out-of-distribution (OOD) or not. If it is not, we must return the set of class labels that have high probability. From [AB21]. Used with kind permission of Anastasios Angelopoulos.

Instead of trying to solve the binary ID/OOD classification problem, we can directly try to predict the class of the input. Let the probabilities over the C labels be $p_c = p(y = c|\mathbf{x})$, and let the logits be $\ell_c = \log p_c$. We can derive a **confidence score** or **uncertainty metric** in a variety of ways from these quantities, e.g., the max probability $s = \max_c p_c$, the margin $s = \max_c \ell_c - \max_c^2 \ell_c$ (where \max^2 means the second largest element), the entropy $s = H(p_{1:C})^3$, the “**energy score**” $s = \sum_c \ell_c$ [Liu+21b], etc. In [Mil+21; Vaz+22] they show that the simple max probability baseline performs very well in practice.

We want to solve the two-stage decision problems illustrated in Figure 19.6. We define the prediction set as follows:

$$\mathcal{T}_\lambda(\mathbf{x}) = \begin{cases} \emptyset & \text{if } \text{OOD}(\mathbf{x}) > \lambda_1 \\ \text{APS}(\mathbf{x}) & \text{otherwise} \end{cases} \quad (19.3)$$

where $\text{OOD}(\mathbf{x})$ is some heuristic OOD score (such as max class probability), and $\text{APS}(\mathbf{x})$ is the adaptive prediction set method of Section 14.3.1, which returns the set of the top K class labels, such that the sum of their probabilities exceeds threshold λ_2 . (Formally, $\text{APS}(\mathbf{x}) = \{\pi_1, \dots, \pi_K\}$ where π sorts $f(\mathbf{x})_{1:C}$ in descending order, and $K = \min\{K' : \sum_{c=1}^{K'} f(\mathbf{x})_c > \lambda_2\}$.)

³ [Kir+21] argues against using entropy, since it confuses uncertainty about which of the C labels to use with uncertainty about whether any of the labels is suitable, compared to a “none-of-the-above” option.

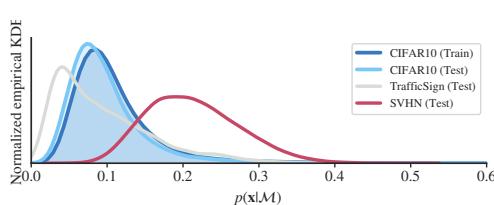


Figure 19.7: Likelihoods from a Glow normalizing flow model (Section 23.2.1) trained on CIFAR10 and evaluated on different test sets. The SVHN street sign dataset has lower visual complexity, and hence higher likelihood. Qualitatively similar results are obtained for other generative models and data set. From Figure 1 of [Ser+20]. Used with kind permission of Joan Serra.

method (see [Ang+21]). The resulting thresholds will jointly minimize the following risks:

$$R_1(\boldsymbol{\lambda}) = p(\mathcal{T}_{\boldsymbol{\lambda}}(\mathbf{x}) = \emptyset) \quad (19.4)$$

$$R_2(\boldsymbol{\lambda}) = p(\mathbf{y} \notin \mathcal{T}_{\boldsymbol{\lambda}}(\mathbf{x}) | \mathcal{T}_{\boldsymbol{\lambda}}(\mathbf{x}) \neq \emptyset) \quad (19.5)$$

where $p(\mathbf{x}, \mathbf{y})$ is the true but unknown source distribution (of ID samples, no OOD samples required), R_1 is the chance that an ID sample will be incorrectly rejected as OOD (type-I error), and R_2 is the chance (conditional on the decision to classify) that the true label is not in the predicted set. The goal is to set λ_1 as large as possible (so we can detect OOD examples when they arise) while controlling the type-I error (e.g., we may want to ensure that we falsely flag (as OOD) no more than 10% of in-distribution samples). We then set λ_2 in the usual way for the APS method in Section 14.3.1.

19.3.2.4 Unsupervised methods

If we don't have labeled examples, a natural approach to OOD detection is to fit an unconditional density model (such as a VAE) to the ID samples, and then to evaluate the likelihood $p(\mathbf{x})$ and compare this to some threshold value. Unfortunately for many kinds of deep model and datasets, we sometimes find that $p(\mathbf{x})$ is lower for samples that are from the source distribution than from a novel target distribution. For example, if we train a pixel-CNN model (Section 22.3.2) or a normalizing-flow model (Chapter 23) on Fashion-MNIST and evaluate it on MNIST, we find it gives higher likelihood to the MNIST samples [Nal+19a; Ren+19; KIW20; ZGR21]. This phenomenon occurs for several other models and datasets (see Figure 19.7).

One solution to this is to use log a **likelihood ratio** relative to a baseline density model, $R(\mathbf{x}) = \log p(\mathbf{x})/q(\mathbf{x})$, as opposed to the raw log likelihood, $L(\mathbf{x}) = \log p(\mathbf{x})$. (This technique was explored in [Ren+19], amongst other papers.) An important advantage of this is that the ratio is invariant to transformations of the data. To see this, let $\mathbf{x}' = \phi(\mathbf{x})$ be some invertible, but possibly nonlinear, transformation. By the change of variables, we have $p(\mathbf{x}') = p(\mathbf{x})|\det \text{Jac}(\phi^{-1})(\mathbf{x})|$. Thus $L(\mathbf{x}')$ will differ from $L(\mathbf{x})$ in a way that depends on the transformation. By contrast, we have $R(\mathbf{x}) = R(\mathbf{x}')$, regardless of ϕ , since

$$R(\mathbf{x}') = \log p(\mathbf{x}') - \log q(\mathbf{x}') = \log p(\mathbf{x}) + \log |\det \text{Jac}(\phi^{-1})(\mathbf{x})| - \log q(\mathbf{x}) - \log |\det \text{Jac}(\phi^{-1})(\mathbf{x})| \quad (19.6)$$

1 Various other strategies have been proposed, such as computing the log-likelihood adjusted by a
 2 measure of the complexity (coding length computed by a lossless compression algorithm) of the input
 3 [Ser+20], computing the likelihood of model features instead of inputs [Mor+21a], etc.
 4

5 A closely related technique relies on **reconstruction error**. The idea is to fit an autoencoder or
 6 VAE (Section 21.2) to the ID samples, and then measure the reconstruction error of the input: a
 7 sample that is OOD is likely to incur larger error (see e.g. [Pol+19]). However, this suffers from the
 8 same problems as density estimation methods.

9 An alternative to trying to estimate the likelihood, or reconstruct the output, is to use a GAN
 10 (Chapter 26) that is trained to discriminate “real” from “fake” data. This has been extended to the
 11 open set recognition setting in the **OpenGAN** method of [KR21b].
 12

13 19.3.3 Selective prediction

14 Suppose the system has a confidence level of p that an input is OOD (see Section 19.3.4 for a
 15 discussion of some ways to compute such confidence scores). If p is below some threshold, the system
 16 may choose to **abstain** from classifying it with a specific label. By varying the threshold, we can
 17 control the tradeoff between accuracy and abstention rate. This is called **selective prediction** (see
 18 e.g., [EW10; GEY19; Ziy+19; JKG18]), and is useful for applications where an error can be more
 19 costly than asking a human expert for help (e.g., medical image classification).
 20

21

22 19.3.3.1 Example: SGLD vs SGD for MLPs

23 One way to improve performance of OOD detection is to “be Bayesian” about the parameters of the
 24 model, so that the uncertainty in their values is reflected in the posterior predictive distribution.
 25 This can result in better performance in selective prediction tasks.

26 In this section, we give a simple example of this, where we fit a shallow MLP to the MNIST dataset
 27 using either standard SGD (specifically RMSprop) or Stochastic Gradient Langevin Dynamics (see
 28 Section 12.7.1), which is a form of MCMC inference. We use 6,000 training steps, where each step
 29 uses a minibatch of size 1,000. After fitting the model to the training set, we evaluate its predictions
 30 on the test set. To assess how well calibrated the model is, we select a subset of predictions whose
 31 confidence is above a threshold t . (The confidence value is just the probability assigned to the MAP
 32 class.) As we increase the threshold t from 0 to 1, we make predictions on fewer examples, but the
 33 accuracy should increase. This is shown in Figure 19.8: the green curve is the fraction of the test set
 34 for which we make a prediction, and the blue curve is the accuracy. On the left we show SGD, and
 35 on the right we show SGLD. In this case, performance is quite similar, although SGD has slightly
 36 higher accuracy. However, the story changes somewhat when there is distribution shift.

37 To study the effects under distribution shift, we apply both models to FashionMNIST data. We
 38 show the results in Figure 19.9. The accuracy of both models is very low (less than the chance level
 39 of 10%). but SGD remains quite confident in many more of its predictions than SGLD, which is
 40 more conservative. To see this, consider a confidence threshold of 0.5: the SGD approach predicts on
 41 about 97% of the examples (recall that the green curve corresponds to the right hand axis), whereas
 42 the SGLD only predicts on about 70% of the examples.

43 More details on the behavior of Bayesian neural networks under distribution shift can be found in
 44 Section 17.4.6.2.
 45

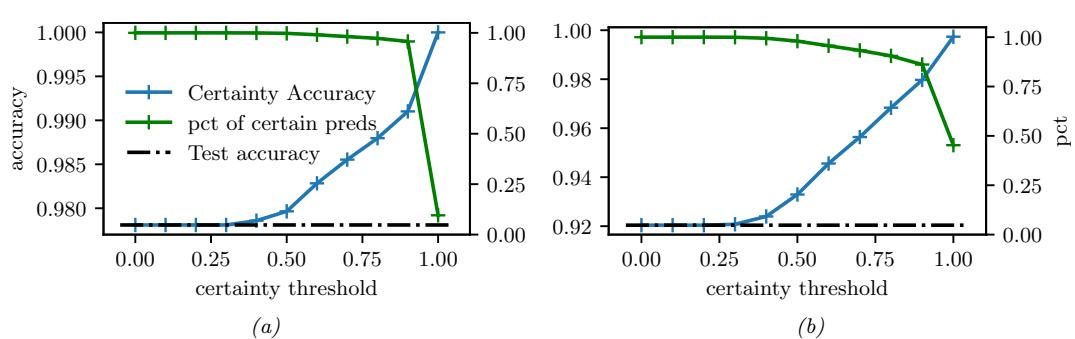


Figure 19.8: Accuracy vs confidence plots for an MLP fit to the MNIST training set, and then evaluated on one batch from the MNIST test set. Scale for blue accuracy curve is on the left, scale for green percentage predicted curve is on the right. (a) Plugin approach, computed using SGD. (b) Bayesian approach, computed using 10 samples from SGLD. Generated by `bnn_mnist_sgld.ipynb`.

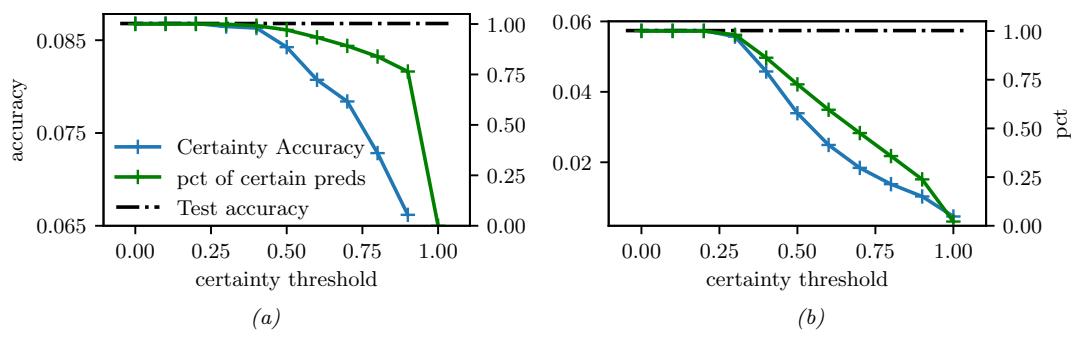


Figure 19.9: Similar to Figure 19.8, except that performance is evaluated on the Fashion MNIST dataset. (a) SGD. (b) SGLD. Generated by `bnn_mnist_sgld.ipynb`.

19.3.4 Open world recognition

In Section 19.3.3, we discussed methods that “refuse to classify” if the input is suspected to be OOD, i.e., is not one of the existing classes. An alternative approach is to treat the input as an example from a new class. If the set of classes is allowed to grow over time in this way, the problem is called **open world classification** [BB15a].⁴ Note that open world classification is most naturally tackled in the context of a continual learning system, which we discuss in Section 19.7.3.

19.4 Robustness to distribution shifts

In this section, we discuss techniques to improve the **robustness** of a model to distribution shifts. In particular, given labeled data from $p(\mathbf{x}, \mathbf{y})$, we aim to create a model that approximates $q(\mathbf{y}|\mathbf{x})$.

⁴ This is not to be confused with **open set recognition**, which refers to classification problems in which the system should label samples from unknown classes as OOD, rather than trying to incrementally learn about new classes. See e.g., [GHC20] for a review of open set recognition methods.

1 **19.4.1 Data augmentation**

3 A simple approach to potentially increasing the robustness of a predictive model to distribution shifts
4 is to simulate samples from the target distribution by modifying the source data. This is called **data**
5 **augmentation**, and is widely used in the deep learning community. For example, it is standard to
6 apply small perturbations to images (e.g., shifting them or rotating them), while keeping the label
7 the same (assuming that the label should be invariant to such changes); see e.g., [SK19; Hen+20] for
8 details. Similarly, in NLP (natural language processing), it is standard to change words that should
9 not affect the label (e.g., replacing “he” with “she” in a sentiment analysis system), or to use **back**
10 **translation** (from a source language to a target language and back) to generate paraphrases; see
11 e.g., [Fen+21] for a review of such techniques. For a causal perspective on data augmentation, see
12 e.g., [Kau+21].
13

14 **19.4.2 Distributionally robust optimization**

16 We can make a discriminative model that is robust to (some forms of) covariate shift by solving the
17 following **distributionally robust optimization** (DRO) problem:

$$\min_{f \in \mathcal{F}} \max_{\mathbf{w} \in \mathcal{W}} \frac{1}{N} \sum_{n=1}^N w_n \ell(f(\mathbf{x}_n), \mathbf{y}_n) \quad (19.7)$$

21 where the samples are from the source distribution, $(\mathbf{x}_n, \mathbf{y}_n) \sim p$. This is an example of a **min-max**
22 **optimization problem**, in which we want to minimize the worst case risk. The specification of the
23 robustness set, \mathcal{W} , is a key factor that determines how well the method works, and how difficult the
24 optimization problem is. Typically it is specified in terms of an ℓ_2 ball around the inputs, but this
25 could also be defined in a feature (embedding space). It is also possible to define the robustness set
26 in terms of local changes to a structural causal model [Mei18a]. For more details on DRO, see e.g.,
27 [WYG14; Hu+18; CP20a; LFG21; Sag+20].
28

29 **19.5 Adapting to distribution shifts**

31 In this section, we discuss techniques to **adapt** the model to the target distribution. If we have some
32 labeled data from the target distribution, we can use transfer learning, as we discuss in Section 19.5.1.
33 However, getting labeled data from the target distribution is often not an option. Therefore, in the
34 other sections, we discuss techniques that just rely on *unlabeled* data from the target distribution.
35

36 **19.5.1 Supervised adaptation using transfer learning**

38 Suppose we have labeled training data from a source distribution, $\mathcal{D}^s = \{(\mathbf{x}_n, \mathbf{y}_n) \sim p : n = 1 : N_s\}$,
39 and also some some labeled data from the target distribution, $\mathcal{D}^t = \{(\mathbf{x}_n, \mathbf{y}_n) \sim q : n = 1 : N_t\}$. Our
40 goal is to minimize the risk on the target distibution q , which can be computed using
41

$$R(f, q) = \mathbb{E}_{q(\mathbf{x}, \mathbf{y})} [\ell(\mathbf{y}, f(\mathbf{x}))] \quad (19.8)$$

43 We can approximate the risk empirically using

$$\hat{R}(f, \mathcal{D}^t) = \frac{1}{|\mathcal{D}^t|} \sum_{(\mathbf{x}_n, \mathbf{y}_n) \in \mathcal{D}^t} \ell(\mathbf{y}_n, f(\mathbf{x}_n)) \quad (19.9)$$

If \mathcal{D}^t is large enough, we can directly optimize this using standard empirical risk minimization (ERM). However, if \mathcal{D}^t is small, we might want to use \mathcal{D}^s somehow as a regularizer. This is called **transfer learning**, since we hope to “transfer knowledge” from p to q . There are many approaches to transfer learning (see e.g., [Zhu+21] for a review). We briefly mention a few below.

19.5.1.1 Pre-train and fine-tune

The simplest and most widely used approach to transfer learning is the **pre-train and fine-tune** approach. We first fit a model to the source distribution by computing $f^s = \operatorname{argmin}_f \hat{R}(f, \mathcal{D}^s)$. (Note that the source data may be unlabeled, in which case we can use self-supervised learning methods.) We then adapt the model to work on the target distribution by computing

$$f^t = \operatorname{argmin}_f \hat{R}(f, \mathcal{D}^t) + \lambda \|f - f^s\| \quad (19.10)$$

where $\|f - f^s\|$ is some distance between the functions, and $\lambda \geq 0$ controls the degree of regularization.

Since we assume that we have very few samples from the target distribution, we typically “freeze” most of the parameters of the source model. (This makes an implicit assumption that the features that are useful for the source distribution also work well for the target.) We can then solve Equation (19.10) by “chopping off the head” from f^s and replacing it with a new linear layer, to map to the new set of labels for the target distribution, and then compute a new MAP estimate for the parameters on the target distribution. (We can also compute a prior for the parameters of the source model, and use it to compute a posterior for the parameters of the target model, as discussed in Section 17.2.3.)

This approach is very widely used in practice, since it is simple and effective. In particular, it is common to take a large pre-trained model, such as a transformer, that has been trained (often using self supervised learning, Section 32.3.3) on a lot of data, such as the entire web, and then to use this model as a feature extractor (see e.g., [Kol+20]). The features are fed to the downstream model, which may be a linear classifier or a shallow MLP, which is trained on the target distribution.

19.5.1.2 Prompt tuning (in-context learning)

Recently another approach to transfer learning has been developed, that leverages large models, such as transformers (Section 22.4), which are trained on massive web datasets, usually in an unsupervised way, and then adapted to a small, task-specific target distribution. The interesting thing about this approach is the parameters of the original model are not changed; instead, the model is simply “conditioned” on new training data, usually in the form of a text **prompt** \mathbf{z} . That is, we compute

$$f^t(\mathbf{x}) = f^s(\mathbf{x} \cup \mathbf{z}) \quad (19.11)$$

where we (manually or automatically) optimize \mathbf{z} while keeping f^s frozen. This approach is called **prompt tuning** or **in-context learning** (see e.g., [Liu+21a]), and is an instance of **few-shot learning** (see Figure 22.5 for an example).

Here \mathbf{z} acts like a small training dataset, and f^s uses attention (Section 16.2.7) to “look at” all its inputs, comparing \mathbf{x} with the examples in \mathbf{z} , and uses this to make a prediction. This works because the text training data often has a similar hierarchical structure (see [Xie+22] for a Bayesian interpretation).

1 **19.5.2 Weighted ERM for covariate shift**

3 In this section we reconsider the risk minimization objective in Equation (19.8), but leverage unlabeled
4 data from the target distribution to estimate it. If we make the covariate shift assumption (i.e.,
5 $q(\mathbf{x}, \mathbf{y}) = q(\mathbf{x})p(\mathbf{y}|\mathbf{x})$), then we have

7
$$R(f, q) = \int q(\mathbf{x})q(\mathbf{y}|\mathbf{x})\ell(\mathbf{y}, f(\mathbf{x}))d\mathbf{x}d\mathbf{y} \quad (19.12)$$

9
$$= \int q(\mathbf{x})p(\mathbf{y}|\mathbf{x})\ell(\mathbf{y}, f(\mathbf{x}))d\mathbf{x}d\mathbf{y} \quad (19.13)$$

11
$$= \int \frac{q(\mathbf{x})}{p(\mathbf{x})}p(\mathbf{x})p(\mathbf{y}|\mathbf{x})\ell(\mathbf{y}, f(\mathbf{x}))d\mathbf{x}d\mathbf{y} \quad (19.14)$$

14
$$\approx \frac{1}{N} \sum_{(\mathbf{x}_n, \mathbf{y}_n) \in \mathcal{D}_L^s} w_n \ell(\mathbf{y}_n, f(\mathbf{x}_n)) \quad (19.15)$$

17 where the weights are given by the ratio

19
$$w_n = w(\mathbf{x}_n) = \frac{q(\mathbf{x}_n)}{p(\mathbf{x}_n)} \quad (19.16)$$

21 Thus we can solve the covariate shift problem by using **weighted ERM**.

23 However, this raises two questions. First, why do we need to use this technique, since a discriminative
24 model $p(\mathbf{y}|\mathbf{x})$ should work for any input \mathbf{x} , regardless of which distribution it comes from. Second,
25 given that we do need to use this method, in practice how should we estimate the weights $w_n =$
26 $w(\mathbf{x}_n) = \frac{q(\mathbf{x}_n)}{p(\mathbf{x}_n)}$. We discuss these issues below.

28 **19.5.2.1 Why is covariate shift a problem for discriminative models?**

29 For a discriminative model of the form $p(\mathbf{y}|\mathbf{x})$, it might seem that such a change in $p(\mathbf{x})$ will not
30 affect the predictions. If the predictor $p(\mathbf{y}|\mathbf{x})$ is the correct model for all parts of the input space \mathbf{x} ,
31 then this conclusion is warranted. However, most models will only be accurate in certain parts of the
32 input space. This is illustrated in Figure 19.10b, where we show that a linear model fit to the source
33 distribution may perform much worse on the target distribution than a model that weights target
34 points more heavily during training.

36 **19.5.2.2 How should we estimating the ERM weights?**

38 One approach to estimating the ERM weights $w_n = w(\mathbf{x}_n) = \frac{q(\mathbf{x}_n)}{p(\mathbf{x}_n)}$ is to learn a density model for
39 the source and target. However, density estimation is difficult for high dimensional features. An
40 alternative approach is to try to approximate the density ratio, by fitting a binary classifier to
41 distinguish the two distributions, as discussed in Section 2.7.5. In particular, suppose we have an
42 equal number of samples from $p(\mathbf{x})$ and $q(\mathbf{x})$. Let us label the first set with $c = -1$ and the second
43 set with $c = 1$. Then we have

45
$$p(c = 1|\mathbf{x}) = \frac{q(\mathbf{x})}{q(\mathbf{x}) + p(\mathbf{x})} \quad (19.17)$$

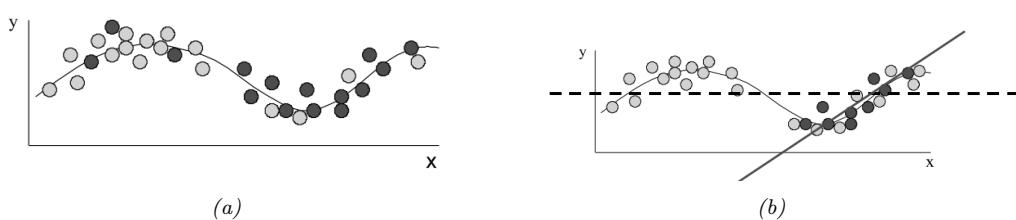


Figure 19.10: (a) Illustration of covariate shift. Light gray represents training distribution, dark gray represents test distribution. We see the test distribution has shifted to the right but the underlying input-output function is constant. (b) Dashed line: fitting a linear model across the full support of X . Solid black line: fitting the same model only on parts of input space that have high likelihood under the test distribution. From Figures 1–2 of [Sto09]. Used with kind permission of Amos Storkey.

and hence $\frac{p(c=1|\mathbf{x})}{p(c=-1|\mathbf{x})} = \frac{q(\mathbf{x})}{p(\mathbf{x})}$. If the classifier has the form $f(\mathbf{x}) = p(c=1|\mathbf{x}) = \sigma(h(\mathbf{x})) = \frac{1}{1+\exp(-h(\mathbf{x}))}$, where $h(\mathbf{x})$ is the prediction function that returns the logits, then the importance weights are given by

$$w_n = \frac{1/(1 + \exp(-h(\mathbf{x}_n)))}{\exp(-h(\mathbf{x}_n))/(1 + \exp(-h(\mathbf{x}_n)))} = \exp(h(\mathbf{x}_n)) \quad (19.18)$$

Of course this method requires that \mathbf{x} values that may occur in the test distribution should also be possible in the training distribution, i.e. $q(\mathbf{x}) > 0 \implies p(\mathbf{x}) > 0$. Hence there are no guarantees about this method being able to interpolate beyond the training distribution.

19.5.3 Unsupervised domain adaptation for covariate shift

We now turn to methods that only need access to unlabeled examples from the target distribution.

The technique of **unsupervised domain adaptation** or **UDA** assumes access to a labeled dataset from the source distribution, $\mathcal{D}_1 = \mathcal{D}_L^s \sim p(\mathbf{x}, \mathbf{y})$ and an unlabeled dataset from the target distribution, $\mathcal{D}_2 = \mathcal{D}_U^t \sim q(\mathbf{x})$. It then uses the unlabeled target data to improve robustness or invariance of the predictor, rather than using a weighted ERM method.

There are many forms of UDA (see e.g., [KL21; CB20] for reviews). Here we just focus on one method, called **domain adversarial learning** [Gan+16a]. Let $f_\alpha : \mathcal{X}_1 \cup \mathcal{X}_2 \rightarrow \mathcal{H}$ be a feature extractor defined on the two input domains, let $c_\beta : \mathcal{H} \rightarrow \{1, 2\}$ be a classifier that maps from the feature space to the domain from which the input was taken, either domain 1 or 2 (source or target), and let $g_\gamma : \mathcal{H} \rightarrow \mathcal{Y}$ be a classifier that maps from the feature space to the label space. We want to train the feature extractor so that it cannot distinguish whether the input is coming from the source or target distribution; in this case, it will only be able to use features that are common to both domains. Hence we optimize

$$\min_{\gamma} \max_{\alpha, \beta} \frac{1}{N_1 + N_2} \sum_{\mathbf{x}_n \in \mathcal{D}_1, \mathcal{D}_2} \ell(d_n, c_\beta(f_\alpha(\mathbf{x}_n))) + \frac{1}{N_1} \sum_{(\mathbf{x}_n, \mathbf{y}_n) \in \mathcal{D}_1} \ell(\mathbf{y}_n, g_\gamma(f_\alpha(\mathbf{x}_n))) \quad (19.19)$$

The objective in Equation (19.19) minimizes the loss on the desired task of classifying y , but maximizes

1 the loss on the auxiliary task of classifying the domain label d . This can be implemented by the
2 **gradient sign reversal** trick, and is related to GANs (Section 26.7.6).
3

4 19.5.4 Unsupervised techniques for label shift

5 In this section, we describe an approach known as **black box shift estimation**, due to [LWS18],
6 which can be used to tackle the label shift problem in an unsupervised way. We assume that the
7 only thing that changes in the target distribution is the label prior, i.e., if the source distribution is
8 denoted by $p(\mathbf{x}, \mathbf{y})$ and target distribution is denoted by $q(\mathbf{x}, \mathbf{y})$, we assume $q(\mathbf{x}, \mathbf{y}) = p(\mathbf{x}|\mathbf{y})q(\mathbf{y})$.
9

10 First note that, for any deterministic function $f : \mathcal{X} \rightarrow \mathcal{Y}$, we have
11

$$\underline{12} \quad p(\mathbf{x}|y) = q(\mathbf{x}|y) \implies p(f(\mathbf{x})|y) = q(f(\mathbf{x})|y) \implies p(\hat{y}|y) = q(\hat{y}|y) \quad (19.20)$$

13 where $\hat{y} = f(\mathbf{x})$ is the predicted label. Let $\mu_i = q(\hat{y} = i)$ be the empirical fraction of times the model
14 predicts class i on the test set, and let $q(y = i)$ be the true but unknown label distribution on the
15 test set, and let $C_{ij} = p(\hat{y} = i|y = j)$ be the class confusion matrix estimated on the training set.
16

17 Then we have

$$\underline{18} \quad \mu_{\hat{y}} = \sum_y q(\hat{y}|y)q(y) = \sum_y p(\hat{y}|y)q(y) = \sum_y p(\hat{y}, y) \frac{q(y)}{p(y)} \quad (19.21)$$

19 We can write this in matrix-vector form as follows:
20

$$\underline{21} \quad \mu_i = \sum_i C_{ij}q_j, \implies \boldsymbol{\mu} = \mathbf{C}\mathbf{q} \quad (19.22)$$

22 Hence we can solve $\mathbf{q} = \mathbf{C}^{-1}\boldsymbol{\mu}$, providing that \mathbf{C} is not singular (this will be the case if \mathbf{C} is strongly
23 diagonal, i.e., the model predicts class y_i correctly more often than any other class y_j). We also
24 require that for every $q(y) > 0$ we have $p(y) > 0$, which means we see every label at training time.

25 Once we know the new label distribution, $q(\mathbf{y})$, we can adjust our discriminative classifier to take
26 the new label prior into account as follows:

$$\underline{27} \quad q(y|\mathbf{x}) = \frac{q(\mathbf{x}|y)q(y)}{q(\mathbf{x})} = \frac{p(\mathbf{x}|y)q(y)}{q(\mathbf{x})} = \frac{p(y|\mathbf{x})p(\mathbf{x})}{p(y)} \frac{q(y)}{q(\mathbf{x})} = p(y|\mathbf{x}) \frac{q(y)}{p(y)} \frac{p(\mathbf{x})}{q(\mathbf{x})} \quad (19.23)$$

28 We can safely ignore the $\frac{p(\mathbf{x})}{q(\mathbf{x})}$ term, which is constant wrt y , and we can plug in our estimates of the
29 label distributions to compute the $\frac{q(y)}{p(y)}$.

30 In summary, there are three requirements for this method: (1) The confusion matrix is invertible;
31 (2) no new labels at test time; (3) the only thing that changes is the label prior. If these three
32 conditions hold, the above approach is a valid estimator. See [LWS18] for the details.
33

34

35 19.5.5 Test-time adaptation

36 In some settings, it is possible to continuously update the model parameters. This allows the model
37 to adapt to changes in the input distribution. This is called **test time adaptation** or **TTA**. The
38 difference from the unsupervised domain adaptation methods of Section 19.5.3 is that, in the online
39 setting, we just have the model which was trained on the source, and not to the source distribution.
40

41

In [Sun+20] they proposed an approach called **TTT** (“test-time training”) for adapting a discriminative model. In this approach, a self-supervised proxy task is used to create pseudo-labels, which can then be used to adapt the model at run time. In more detail, suppose we create a Y-structured network, where we first perform feature extraction, $\mathbf{x} \rightarrow \mathbf{h}$, and then use \mathbf{h} to predict the output \mathbf{y} and some proxy output \mathbf{r} , such as the angle of rotation of the input image. The rotation angle is known if we use data augmentation. Hence we can apply this technique at test time, even if \mathbf{y} is unknown, and update the $\mathbf{x} \rightarrow \mathbf{h} \rightarrow \mathbf{r}$ part of the network, which influences the prediction for \mathbf{y} via the shared bottleneck (feature layer) \mathbf{h} .

Of course, if the proxy output, such as the rotation angle, is not known, we cannot use proxy-supervised learning methods such as TTT. In [Wan+20a], they propose an approach, inspired by semi-supervised learning methods, which they call **TENT**, which stands for “test-time adaptation by entropy minimization”. The idea is to update the classifier parameters to minimize the entropy of the predictive distribution on a batch of test examples. In [Goy+22], they give a justification for this heuristic from the meta-learning perspective. In [ZL21], they present a Bayesian version of TENT, which they call **BACS**, which stands for “Bayesian adaptation under covariate shift”. In [ZLF21], they propose a method called **MEMO** (“marginal entropy minimization with one test point”) that can be used for any architecture. The idea is, once again, to apply data augmentation at test time to the input \mathbf{x} , to create a set of inputs, $\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_B$. Now we update the parameters so as to minimize the predictive entropy produced by the averaged distribution

$$\bar{p}(\mathbf{y}|\mathbf{x}, \mathbf{w}) = \frac{1}{B} \sum_{b=1}^B p(\mathbf{y}|\tilde{\mathbf{x}}_b, \mathbf{w}) \quad (19.24)$$

This ensures that the model gives the same predictions for each perturbation of the input, and that the predictions are confident (low entropy).

An alternative to entropy based methods is to use pseudo-labels (predicted outputs from the source model), and then to self-train on these. This has been shown to work well for unsupervised domain adaptation [KML20]. In the context of TTA, [Che+22] use this technique, combined with contrastive learning, to adapt a classifier to a new distribution.

19.6 Learning from multiple distributions

In Section 19.2, we discussed the setting in which a model is trained on a single source distribution, and then evaluated on a distinct target distribution. In this section, we generalize this to a setting in which the model is trained on data from $J \geq 2$ source distributions, before being tested on data from a target distribution. This includes a variety of different problem settings, depending on the value of J , as we summarize in Figure 19.11.

19.6.1 Multi-task learning

In **multi-task learning** (MTL) [Car97], we have labeled data from J different distributions, $\mathcal{D}^j = \{(\mathbf{x}_n^j, \mathbf{y}_n^j) : n = 1 : N_j\}$, and the goal is to learn a model that predicts well on all J of them simultaneously, where $f(\mathbf{x}, j) : \mathcal{X} \rightarrow \mathcal{Y}_j$ is the output for the j 'th task. For example, we might want to map a color image of size $H \times W \times 3$ to a set of semantic labels per pixel, $\mathcal{Y}^1 = \{1, \dots, C\}^{HW}$, as well as a set of predicted depth values per pixel, $\mathcal{Y}^2 = \mathbb{R}^{HW}$. We can do this using ERM where we

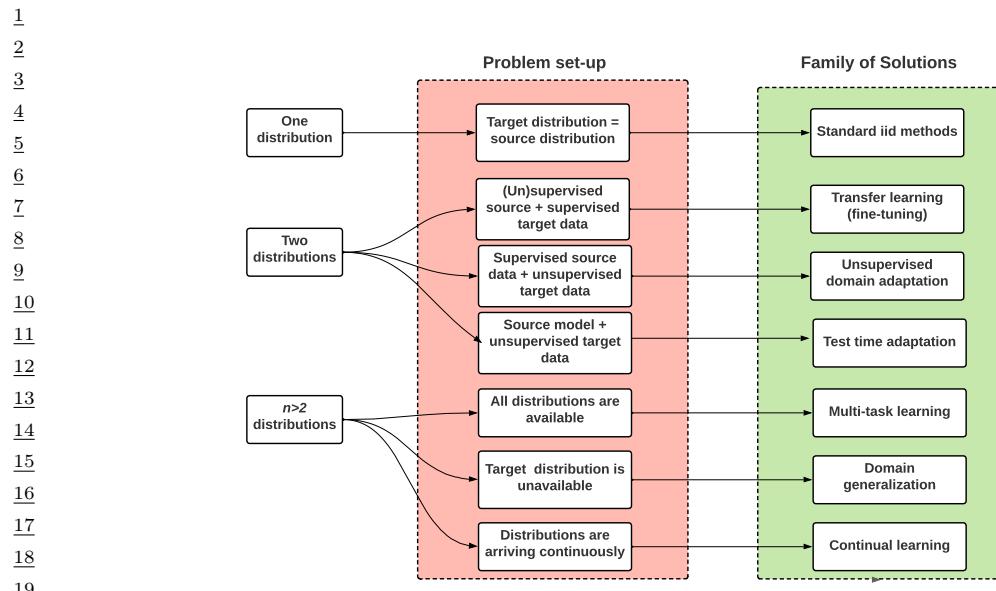


Figure 19.11: Schematic overview of techniques for learning from 1 or more different distributions. Adapted from slide 3 of [Sca21].

have multiple samples for each task:

$$f^* = \operatorname{argmin}_f \sum_{j=1}^J \sum_{n=1}^{N_j} \ell_j(y_n^j, f(x_n^j, j)) \quad (19.25)$$

where ℓ_j is the loss function for task j .

There are many approaches to solving MTL. The simplest is to fit a single model with multiple “output heads”, as illustrated in Figure 19.12. The main challenge is defining a suitable architecture (i.e., deciding which parts of the feature extractor to share across tasks), and how to balance the different losses from each task. See [BLS11] for a theoretical analysis of this problem, and [ZY21] for a more detailed review of neural approaches.

Note that multi-task learning does not always help performance on each task because sometimes there can be “task interference” or “negative transfer” (see e.g., [WZR20]). In such cases, we should use separate networks, rather than using one model with multiple output heads.

19.6.2 Domain generalization

The problem of **domain generalization** assumes we train on J different labeled source distributions or “environments” (also called “domains”), and then test on a new target distribution (denoted by $t = J + 1$). In some cases each environment is just identified with a meaningless integer id. In more realistic settings, each different distribution has associated **meta-data** or **context variables** that characterize the environment in which the data was collected, such as the time, location, imaging device, etc.

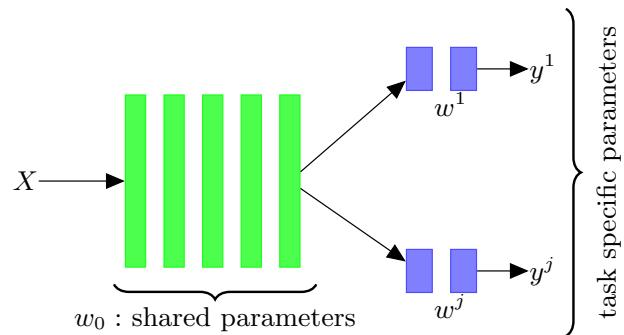


Figure 19.12: Illustration of multi-headed network for multi-task learning.

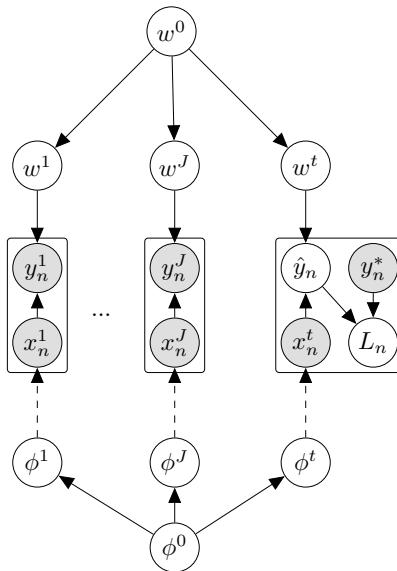


Figure 19.13: Hierarchical Bayesian discriminative model for learning from J different environments (distributions), and then testing on a new target distribution $t = J + 1$. Here \hat{y}_n is the prediction for test example \mathbf{x}_n , y_n^* is the true output, and $\ell_n = \ell(\hat{y}_n, y_n^*)$ is the associated loss. The parameters of the distribution over input features $p_\phi(\mathbf{x})$ are shown with dotted edges, since these distributions do not need to be learned in a discriminative model.

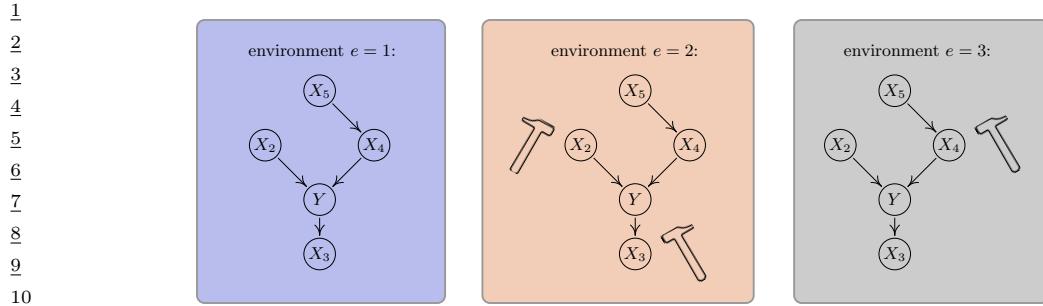


Figure 19.14: Illustration of invariant causal prediction. The hammer symbol represents variables whose distribution is perturbed in the given environment. An invariant predictor must use features $\{X_2, X_4\}$. Considering indirect causes instead of direct ones (e.g. $\{X_2, X_5\}$) or an incomplete set of direct causes (e.g. $\{X_4\}$) may not be sufficient to guarantee invariant prediction. From Figure 1 of [PBM16b]. Used with kind permission of Jonas Peters.

16

17

18

19 Domain generalization (DG) is similar to multi-task learning, but differs in what we want to
20 predict. In particular, in DG, we only care about prediction accuracy on the target distribution, not
21 the J training distribution. Furthermore, we assume we don't have any labeled data from the target
22 distribution. We therefore have to make some assumptions about how $p^t(\mathbf{x}, \mathbf{y})$ relates to $p^j(\mathbf{x}, \mathbf{y})$ for
23 $j = 1 : J$.

24 One way to formalize this is to create a hierarchical Bayesian model, as proposed in [Bax00], and
25 illustrated in Figure 19.13. This encodes the assumption that $p^t(\mathbf{x}, \mathbf{y}) = p(\mathbf{x}|\phi^t)p(\mathbf{y}|\mathbf{x}, \mathbf{w}^t)$ where
26 \mathbf{w}^t is derived from a common "population level" model \mathbf{w}^0 , shared across all distributions, and
27 similarly for ϕ^t . (Note, however, that in a discriminative model, we don't need to model $p(\mathbf{x}|\phi^t)$.)
28 See Section 15.5 for discussion of hierarchical Bayesian GLMs, and Section 17.6 for discussion of
29 hierarchical Bayesian MLPs.

30 Many other techniques have been proposed for DG. Note, however, that [GLP21] found that none
31 of these methods worked consistently better than the baseline approach of performing empirical
32 risk minimization across all the provided datasets. For more information, see e.g., [GLP21; She+21;
33 Wan+21; Chr+21].

34

35 19.6.3 Invariant risk minimization

36

37 One approach to domain generalization that has received a lot of attention is called **invariant**
38 **risk minimization** or **IRM** [Arj+19]. The goal is to learn a predictor that works well across all
39 environments, yet is less prone to depending on the kinds of "spurious features" we discussed in
40 Section 19.2.1.

41 IRM is an extension of an earlier method called **invariant causal prediction** (ICP) [PBM16b].
42 This uses hypothesis testing methods to find the set of predictors (features) that directly cause the
43 outcome in each environment, rather than features that are indirect causes, or are just correlated
44 with the outcome. See Figure 19.14 for an illustration.

45 In [Arj+19], they proposed an extension of ICP to handle the case of high dimensional inputs,
46 where the individual variables do not have any causal meaning (e.g., they correspond to pixels).

47

Their approach requires finding a predictor that works well on average, across all environments, while also being optimal for each individual environment. That is, we want to find

$$f^* = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \sum_{j=1}^J \frac{1}{N_j} \sum_{n=1}^{N_j} \ell(\mathbf{y}_n^j, f(\mathbf{x}_n^j)) \quad (19.26)$$

$$\text{such that } f \in \arg \min_{g \in \mathcal{F}} \frac{1}{N_j} \sum_{n=1}^{N_j} \ell(\mathbf{y}_n^j, g(\mathbf{x}_n^j)) \text{ for all } j \in \mathcal{E} \quad (19.27)$$

where \mathcal{E} is the set of environments, and \mathcal{F} is the set of prediction functions. The intuition behind this is as follows: there may be many functions that achieve low empirical loss on any given environment, since the problem may be underspecified, but if we pick the one that also works well on all environments, it is more likely to rely on causal features rather than spurious features.

Unfortunately, more recent work has shown that the IRM principle often does not work well for covariate shift, both in theory [RRR21] and practice [GLP21], although it can work well in some anti-causal (generative) models [Ahu+21].

19.6.4 Meta-learning

The goal of **meta-learning** is to “learn the learning algorithm” [TP97]. A common way to do this is to provide the meta-learner with a set of datasets from different distributions. This is very similar to domain generalization (Section 19.6.2), except that we partition each training distribution into training and test, so we can “practice” learning to generalize from a training set to a test set. A general review of meta-learning can be found in [Hos+20a]. Here we present a unifying summary based on the hierarchical Bayesian framework proposed in [Gor+19].

19.6.4.1 Meta-learning as probabilistic inference for prediction

We assume there are J tasks (distributions), each of which has a training set $\mathcal{D}_{\text{train}}^j = \{(\mathbf{x}_n^j, \mathbf{y}_n^j) : n = 1 : N^j\}$ and a test set $\mathcal{D}_{\text{test}}^j = \{(\tilde{\mathbf{x}}_m^j, \tilde{\mathbf{y}}_m^j) : m = 1 : M^j\}$. In addition, \mathbf{w}^j are the task specific parameters, and \mathbf{w}^0 are the shared parameters, as shown in Figure 19.15. This is very similar to the domain generalization model in Figure 19.13, except for two differences: first there is the trivial difference due to the use of plate notation; second, in meta learning, we have both training and test partitions for all distributions, whereas in DG, we only have a test set for the target distribution.

We will learn a point estimate for the global parameters \mathbf{w}^0 , since it is shared across all datasets, and thus has little uncertainty. However, we will compute an approximate posterior for \mathbf{w}^j , since each task often has little data. We denote this posterior by $p(\mathbf{w}^j | \mathcal{D}_{\text{train}}^j, \mathbf{w}^0)$. From this, we can compute the posterior predictive distribution for each task:

$$p(\tilde{\mathbf{y}}^j | \tilde{\mathbf{x}}^j, \mathcal{D}_{\text{train}}^j, \mathbf{w}^0) = \int p(\tilde{\mathbf{y}}^j | \tilde{\mathbf{x}}^j, \mathbf{w}^j) p(\mathbf{w}^j | \mathcal{D}_{\text{train}}^j, \mathbf{w}^0) d\mathbf{w}^j \quad (19.28)$$

Since computing the posterior is in general intractable, we will learn an amortized approximation (see Section 10.3.6) to the predictive distribution, denoted by $q_\phi(\tilde{\mathbf{y}}^j | \tilde{\mathbf{x}}^j, \mathcal{D}_{\text{train}}^j, \mathbf{w}^0)$. We choose the parameters of the prior \mathbf{w}^0 and the inference network ϕ to make this *predictive posterior* as accurate

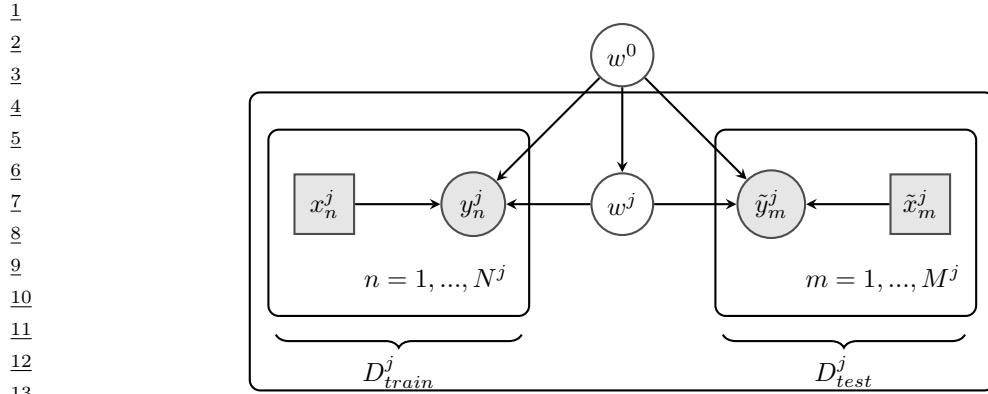


Figure 19.15: Hierarchical Bayesian model for meta-learning. There are J tasks, each of which has a training set $\mathcal{D}^j = \{(\mathbf{x}_n^j, \mathbf{y}_n^j) : n = 1 : N^j\}$ and a test set $\mathcal{D}_{\text{test}}^j = \{(\tilde{\mathbf{x}}_m^j, \tilde{\mathbf{y}}_m^j) : m = 1 : M^j\}$. \mathbf{w}^j are the task specific parameters, and $\boldsymbol{\theta}$ are the shared parameters. Adapted from Figure 1 of [Gor+19].

as possible for any given input dataset:

$$\phi^* = \underset{\phi}{\operatorname{argmin}} \mathbb{E}_{p(\mathcal{D}_{\text{train}}, \tilde{\mathbf{x}})} [D_{\text{KL}}(p(\tilde{\mathbf{y}}|\tilde{\mathbf{x}}, \mathcal{D}_{\text{train}}, \mathbf{w}^0) \| q_\phi(\tilde{\mathbf{y}}|\tilde{\mathbf{x}}, \mathcal{D}_{\text{train}}, \mathbf{w}^0))] \quad (19.29)$$

$$= \underset{\phi}{\operatorname{argmin}} \mathbb{E}_{p(\mathcal{D}_{\text{train}}, \tilde{\mathbf{x}})} [\mathbb{E}_{p(\tilde{\mathbf{y}}|\tilde{\mathbf{x}}, \mathcal{D}_{\text{train}}, \mathbf{w}^0)} [\log q_\phi(\tilde{\mathbf{y}}|\tilde{\mathbf{x}}, \mathcal{D}_{\text{train}}, \mathbf{w}^0)]] \quad (19.30)$$

$$= \underset{\phi}{\operatorname{argmin}} \mathbb{E}_{p(\mathcal{D}_{\text{train}}, \tilde{\mathbf{x}}, \tilde{\mathbf{y}})} \left[\log \int p(\tilde{\mathbf{y}}|\tilde{\mathbf{x}}, \mathbf{w}) q_\phi(\mathbf{w}|\mathcal{D}_{\text{train}}, \mathbf{w}^0) d\mathbf{w} \right] \quad (19.31)$$

where we made the approximation $p(\tilde{\mathbf{y}}|\tilde{\mathbf{x}}, \mathcal{D}_{\text{train}}, \mathbf{w}^0) \approx p(\tilde{\mathbf{y}}|\tilde{\mathbf{x}}, \mathcal{D}_{\text{train}})$. We can then make a Monte Carlo approximation to the outer expectation by sampling J tasks (distributions) from $p(\mathcal{D})$, each of which gets partitioned into a train and test set, $\{(\mathcal{D}_{\text{train}}^j, \mathcal{D}_{\text{test}}^j) \sim p(\mathcal{D}) : j = 1 : J\}$, where $\mathcal{D}_{\text{test}}^j = \{(\tilde{\mathbf{x}}_m^j, \tilde{\mathbf{y}}_m^j)\}$. We can make an MC approximation to the inner expectation (the integral) by drawing S samples from the task-specific parameter posterior $\mathbf{w}_s^j \sim q_\phi(\mathbf{w}^j|\mathcal{D}^j, \mathbf{w}^0)$. The resulting objective has the following form (where we assume each test set has M samples for notational simplicity):

$$\mathcal{L}_{\text{meta}}(\mathbf{w}^0, \phi) = \frac{1}{MJ} \sum_{m=1}^M \sum_{j=1}^J \log \left(\frac{1}{S} \sum_{s=1}^S p(\tilde{\mathbf{y}}_m^j|\tilde{\mathbf{x}}_m^j, \mathbf{w}_s^j) \right) \quad (19.32)$$

Note that this is different from standard (amortized) variational inference, that focuses on approximating the expected accuracy of the *parameter posterior* given all of the data for a task, $\mathcal{D}_{\text{all}}^j = \mathcal{D}_{\text{train}}^j \cup \mathcal{D}_{\text{test}}^j$; rather than focusing on predictive accuracy of a test set given a training set. Indeed, the standard objective has the form

$$\mathcal{L}_{\text{VI}}(\mathbf{w}^0, \phi) = \frac{1}{J} \sum_{j=1}^J \left(\sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{all}}^j} \left[\frac{1}{S} \sum_{s=1}^S \log p(\tilde{\mathbf{y}}^j|\tilde{\mathbf{x}}^j, \mathbf{w}_s^j) \right] - D_{\text{KL}}(q_\phi(\mathbf{w}^j|\mathcal{D}_{\text{all}}^j, \mathbf{w}^0) \| p(\mathbf{w}^j|\mathbf{w}^0)) \right)$$

where $\mathbf{w}_s^j \sim q_\phi(\mathbf{w}^j | \mathcal{D}_{\text{all}}^j)$. We see that the standard formulation takes the average of a log, but the meta-learning formulation takes the log of an average. The latter can give provably better predictive accuracy, as pointed out in [MAD20]. Another difference is that the meta-learning formulation optimizes the forward KL, not reverse KL. Finally, in the meta-learning formulation, we do not have the KL penalty term on the parameter posterior.

Below we show how this framework includes several common approaches to meta-learning.

19.6.4.2 Neural processes

In the special case that the task-specific inference network computes a point estimate, $q(\mathbf{w}^j | \mathcal{D}^j, \mathbf{w}^0) = \delta(\mathbf{w}^j - \mathcal{A}_\phi(\mathcal{D}^j, \mathbf{w}^0))$, the posterior predictive distribution becomes

$$q(\tilde{\mathbf{y}}^j | \tilde{\mathbf{x}}^j, \mathcal{D}^j, \mathbf{w}^0) = \int p(\tilde{\mathbf{y}}^j | \tilde{\mathbf{x}}^j, \mathbf{w}^j) q(\mathbf{w}^j | \mathcal{D}^j, \mathbf{w}^0) d\mathbf{w}^j = p(\tilde{\mathbf{y}}^j | \tilde{\mathbf{x}}^j, \mathcal{A}_\phi(\mathcal{D}^j, \mathbf{w}^0), \mathbf{w}^0) \quad (19.34)$$

where $\mathcal{A}_\phi(\mathcal{D}^j, \mathbf{w}^0)$ is a function that takes in a set, and returns some parameters. We can evaluate this predictive distribution empirically, and directly optimize it (wrt ϕ and \mathbf{w}^0) using standard supervised maximum likelihood methods. This approach is called a **neural process** [Gar+18e; Gar+18d]. (See [DGF20] for a good tutorial.)

However, a more common approach to meta-learning is to specify the form of the estimator $\mathcal{A}_\phi(\mathcal{D}^j, \mathbf{w}^0)$ using various heuristics, some of which we discuss below.

19.6.4.3 Gradient-based meta-learning (MAML)

In **gradient-based meta-learning**, we define the task specific inference procedure as follows:

$$\hat{\mathbf{w}}^j = \mathcal{A}(\mathcal{D}^j, \mathbf{w}^0) = \mathbf{w}^0 + \eta \nabla_{\mathbf{w}} \log \sum_{n=1}^{N^j} p(\mathbf{y}_n^j | \mathbf{x}_n^j, \mathbf{w})|_{\mathbf{w}^0} \quad (19.35)$$

That is, we set the task specific parameters to be shared parameters \mathbf{w}^0 , modified by one step along the gradient of the log conditional likelihood. This approach is called **model-agnostic meta-learning** or **MAML** [FAL17]. It is also possible to take multiple gradient steps, by feeding the gradient into an RNN [RL17].

19.6.4.4 Metric-based few-shot learning (prototypical networks)

Now suppose \mathbf{w}^0 correspond to the parameters of a shared neural feature extractor, $h_{\mathbf{w}^0}(\mathbf{x})$, and the task specific parameters are the weights and biases of the last linear layer of a classifier, $\mathbf{w}^j = \{\mathbf{w}_c^j, b_c^j\}_{c=1}^C$. Let us compute the average of the feature vectors for each class in each task's training set:

$$\boldsymbol{\mu}_c^j = \frac{1}{|\mathcal{D}_c^j|} \sum_{\mathbf{x}_n^c \in \mathcal{D}_c^j} h_{\mathbf{w}^0}(\mathbf{x}_n^c) \quad (19.36)$$

1 Now define the task specific inference procedure as follows. We first compute the vector containing
2 the centroid and norm for each class:

3

$$\hat{\mathbf{w}}^j = \mathcal{A}(\mathcal{D}^j, \mathbf{w}^0) = [\boldsymbol{\mu}_c^j, -\frac{1}{2} \|\boldsymbol{\mu}_c^j\|^2]_{c=1}^C \quad (19.37)$$

4

5 The predictive distribution becomes

6

7

$$q(\tilde{y}^j = c | \tilde{\mathbf{x}}^j, \mathcal{D}^j, \mathbf{w}^0) \propto \exp(-d(h_{\mathbf{w}^0}(\tilde{\mathbf{x}}), \boldsymbol{\mu}_c^j)) = \exp\left(h_{\mathbf{w}^0}(\tilde{\mathbf{x}})^T \boldsymbol{\mu}_c^j - \frac{1}{2} \|\boldsymbol{\mu}_c^j\|^2\right) \quad (19.38)$$

8

9 where $d(\mathbf{u}, \mathbf{v})$ is the Euclidean distance. This is equivalent to the technique known as **prototypical**
10 **networks** [SSZ17].

11

12 19.7 Continual learning

13

14 In this section, we discuss **continual learning** (see e.g., [Had+20; Del+21; Qu+21; LCR21; Mai+22;
15 Lin+22]), also called **life-long learning** (see e.g., [Thr98; CL18]), in which the system learns from a
16 sequence of different distributions, p_1, p_2, \dots . In particular, at each time step t , the model receives a
17 batch of labeled data,

18

19

$$\mathcal{D}_t = \{(\mathbf{x}_n, \mathbf{y}_n) : n = 1 : N_t, \mathbf{x}_n \sim p_t(\mathbf{x}), \mathbf{y}_n = f_t(\mathbf{x}_n)\} \quad (19.39)$$

20

21 where $p_t(\mathbf{x})$ is the unknown input distribution, and $f_t : \mathcal{X}_t \rightarrow \mathcal{Y}_t$ is the unknown prediction function.
22 (We focus on noise-free outputs, for notational simplicity.) The learner is then expected to update its
23 belief state about the true function f_t , and to use its beliefs to make predictions on an independent
24 test set,

25

26

$$\mathcal{D}_t^{\text{test}} = \{(\mathbf{x}_n, \mathbf{y}_n) : n = 1 : N_t^{\text{test}}, \mathbf{x}_n \sim p_t^{\text{test}}(\mathbf{x}), \mathbf{y}_n = f_t^{\text{test}}(\mathbf{x}_n)\} \quad (19.40)$$

27

28 Depending on how we assume $p_t(\mathbf{x})$ and f_t evolve over time, and how the test set is defined, we can
29 create a variety of different CL scenarios, as we discuss below.

30

31 19.7.1 Domain drift

32

33 The problem of **domain drift** refers to the setting in which $p_t(\mathbf{x})$ changes over time (i.e., covariate
34 shift), but the functional mapping $f_t : \mathcal{X} \rightarrow \mathcal{Y}$ is constant. For example, the vision system of a
35 self driving car may have to classify cars vs pedestrians under shifting lighting conditions (see e.g.,
36 [Sun+22]).

37 To evaluate such a model, we assume $f_t^{\text{test}} = f_t$ and define $p_t^{\text{test}}(\mathbf{x})$ to be the current input
38 distribution p_t (e.g., if it is currently night time, we want the detector to work well on dark images).
39 Alternatively we can define $p_t^{\text{test}}(\mathbf{x})$ to be the union of all the input distributions seen so far,
40 $p_t^{\text{test}} = \cup_{s=1}^T p_s$ (e.g., we want the detector to work well on dark and light images). This latter
41 assumption is illustrated in Figure 19.16.

42

43

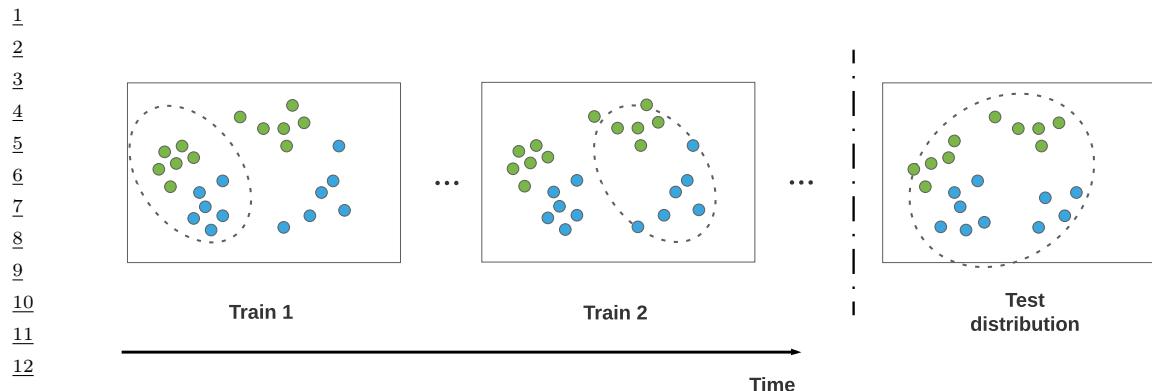


Figure 19.16: An illustration of domain drift.

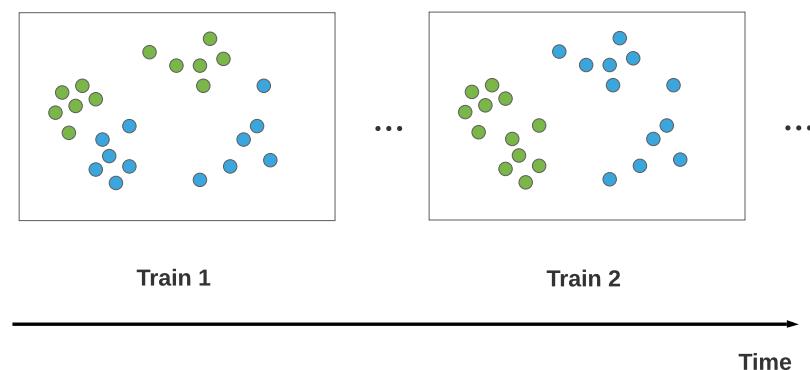


Figure 19.17: An illustration of concept drift.

19.7.2 Concept drift

The problem of **concept drift** refers to the setting where the functional mapping $f_t : \mathcal{X} \rightarrow \mathcal{Y}$ changes over time, but the input distribution $p_t(\mathbf{x})$ is constant [WK96]. For example, we can imagine a setting in which people engage in certain behaviors, and at step t some of these are classified as illegal, and at step $t' > t$, the definition of what is legal changes, and hence the decision boundary changes. This is illustrated in Figure 19.17.

As another example, we might initially be faced with a sort-by-color task, where red objects go on the left and blue objects on the right, and then a sort-by-shape task, where square objects go on the left and circular objects go on the right.⁵ We can think of this as a problem where $p(y|\mathbf{x}, \text{task})$ is stationary, but the task is unobserved, so $p(y|\mathbf{x})$ changes.

In the concept drift scenario, we see that the prediction for the same underlying input point $\mathbf{x} \in \mathcal{X}$ will change depending on when the prediction is performed. This means that the test distribution

⁵. This example is from Mike Mozer.

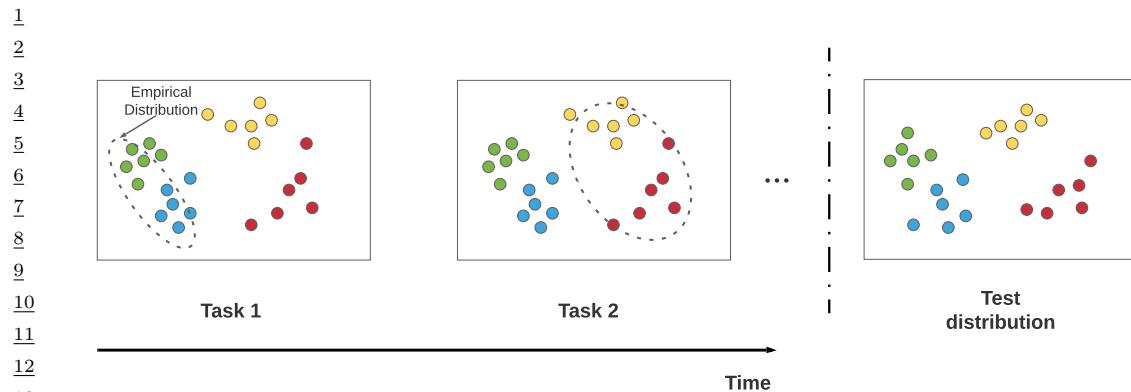


Figure 19.18: An illustration of class incremental learning. Adapted from Figure 1 of [LCR21].

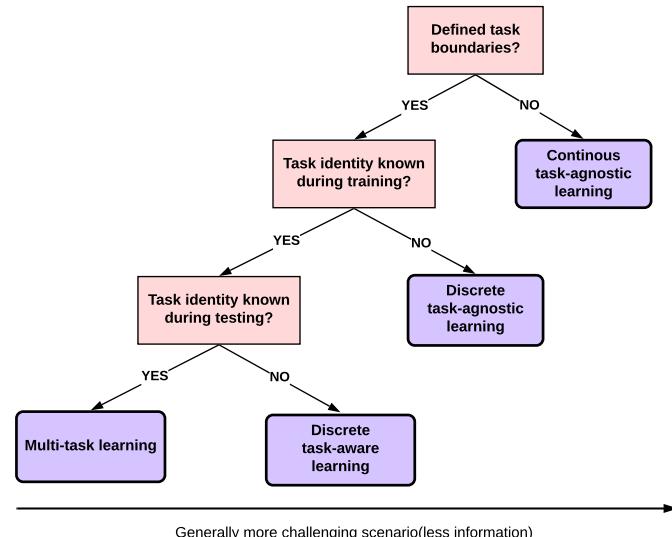


Figure 19.19: Different kinds of incremental learning. Adapted from Figure 1 of [Zen+18].

also needs to change over time for meaningful identification. Alternatively, we can “tag” each input with the corresponding time stamp or task id.

19.7.3 Task incremental learning

A very widely studied form of continual learning focuses on the setting in which new class labels are “revealed” over time. That is, there is assumed to be a true static prediction function $f : \mathcal{X} \rightarrow \mathcal{Y}$, but at step t , the learner only sees samples from $(\mathcal{X}, \mathcal{Y}_t)$, where $\mathcal{Y}_t \subset \mathcal{Y}$. For example, \mathcal{X} may be

the space of images, and \mathcal{Y}_1 might be {cats, dogs}, and \mathcal{Y}_2 might be {cars, bikes, trucks}. Learning to classify with an increasing number of categories is called **class incremental learning** (see e.g., [Mas+20]). This is also called **task incremental learning**, since each distribution is considered as a different **task**. See Figure 19.18 for an illustration.

The problem of class incremental learning has been studied under a variety of different assumptions, as discussed in [Hsu+18; VT18; FG18; Del+21]. The most common scenarios are shown in Figure 19.19. If we assume there are no well defined boundaries between tasks, we have **continuous task-agnostic learning** (see e.g., [SKM21; Zen+21]). If there are well defined boundaries (i.e., discontinuous changes of the training distribution), then we can distinguish two subcases. If the boundaries are not known during training (similar to detecting distribution shift), we have **discrete task-agnostic learning**. Finally, if the boundaries are given to the training algorithm, we have a **task-aware learning** problem.

A common experimental setup in the task-aware setting is to define each task to be a different version of the MNIST dataset, e.g., with all 10 classes present but with the pixels randomly permuted (this is called **permuted MNIST**) or with a subset of 2 classes present at each step (this is called **split MNIST**).⁶ In the task-aware setting, the task label may or may not be known at test time. If it is, the problem is essentially equivalent to multi-task learning (see Section 19.6.1). If it is not, the model must predict the task and corresponding class label within that task (which is a standard supervised problem with a hierarchical label space); this is commonly done by using a multi-headed DNN, with CT outputs, where C is the number of classes, and T is the number of tasks.

In the multi-headed approach, the number of “heads” is usually specified as input to the algorithm, because the softmax imposes a sum-to-one constraint that prevents incremental estimation of the output weights in the open-class setting. An alternative approach is to wait until a new class label is encountered for the first time, and then train the model with an enlarged output head. This requires storing past data from each class, as well as data for the new class (see e.g., [PTD20]). Alternatively, we can use generative classifiers where we do not need to worry about “output heads”. If we use a “deep” nearest neighbor classifier, with a shared feature extractor (embedding function), the main challenge is to efficiently update the stored prototypes for past classes as the feature extractor parameters change (see e.g., [DLT21]). If we fit a separate generative model per class (e.g., a VAE, as in [VLT21]), then online learning becomes easier, but the method may be less sample efficient.

At the time of writing, most of the CL literature focuses on the task-aware setting. However, from a practical point of view, the assumption that task boundaries are provided at training or test time is very unrealistic. For example, consider the problem of training a robot to perform various activities: The data just streams in, and the robot must learn what to do, without anyone telling it that it is now being given an example from a new task or distribution (see e.g., [Fon+21; Woł+21]). Thus future research should focus on the task-agnostic setting, with either discrete or continuous changes.

19.7.4 Catastrophic forgetting

In the class incremental learning literature, it is common to train on a sequence of tasks, but to test (at each step) on all tasks. In this scenario, there are two main possible failure modes. The first possible problem is called “**catastrophic forgetting**” (see e.g., [Rob95b; Fre99; Kir+17]). This

⁶ In the split MNIST setup, for task 1, digits (0,1) get labeled as (0,1), but in task 2, digits (2,3) get labeled as (0,1). So the “meaning” of the output label depends on what task we are solving. Thus the output space is really hierarchical, namely the cross product of task id and class label.

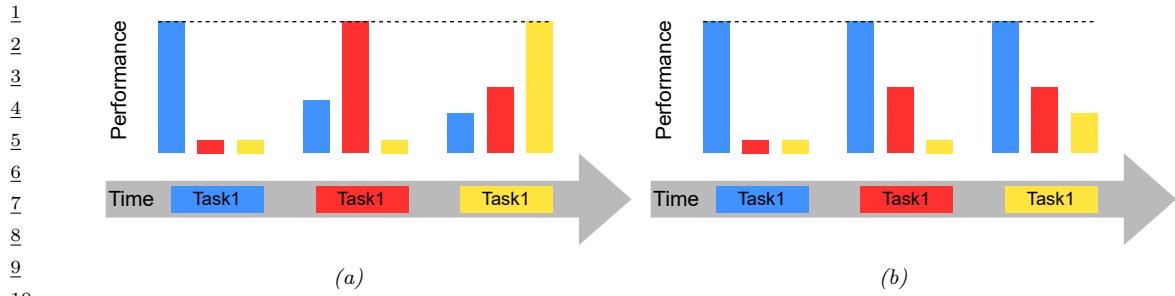


Figure 19.20: Some failure modes in class incremental learning. We train on task 1 (blue) and evaluate on tasks 1–3 (blue, orange, yellow); we then train on task 2 and evaluate on tasks 1–3; etc. (a) Catastrophic forgetting refers to the phenomenon in which performance on a previous task drops when trained on a new task. (b) Too little plasticity (e.g., due to too much regularization) refers to the phenomenon in which only the first task is learned. Adapted from Figure 2 of [Had+20].

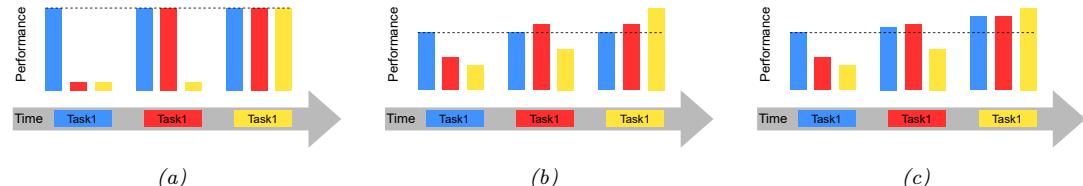


Figure 19.21: What success looks like for class incremental learning. We train on task 1 (blue) and evaluate on tasks 1–3 (blue, orange, yellow); we then train on task 2 and evaluate on tasks 1–3; etc. (a) No forgetting refers to the phenomenon in which performance on previous tasks does not degrade over time. (b) Forward transfer refers to the phenomenon in which training on past tasks improves performance on future tasks beyond what would have been obtained by training from scratch. (c) Backwards transfer refers to the phenomenon in which training on future tasks improves performance on past tasks beyond what would have been obtained by training from scratch. Adapted from Figure 2 of [Had+20].

refers to the phenomenon in which performance on a previous task drops when trained on a new task (see Figure 19.20(a)). Another possible problem is that only the first task is learned, and the model does not adapt to new tasks (see Figure 19.20(b)).

If we avoid these problems, we should expect to see the performance profile in Figure 19.21(a), where performance of incremental training is equal to training on each task separately. However, we might hope to do better by virtue of the fact that we are training on multiple tasks, which are often assumed to be related. In particular, we might hope to see **forward transfer**, in which training on past tasks improves performance on future tasks beyond what would have been obtained by training from scratch (see Figure 19.21(b)). Additionally, we might hope to see **backwards transfer**, in which training on future tasks improves performance on past tasks (see Figure 19.21(c)).

We can quantify the degree of transfer as follows, following [LPR17]. If R_{ij} is the performance on task j after it was trained on task i , R_j^{ind} is the performance on task j when trained just on j , and

there are T tasks, then the amount of forward transfer is

$$\text{FWT} = \frac{1}{T} \sum_{j=1}^T R_{j,j} - R_j^{\text{ind}} \quad (19.41)$$

and the amount of backwards transfer is

$$\text{BWT} = \frac{1}{T} \sum_{j=1}^T R_{T,j} - R_{j,j} \quad (19.42)$$

There are many methods that have been devised to overcome the problem of catastrophic forgetting, but we can group them into three main types. The first is **regularization methods**, which add a loss to preserve information that is relevant to old tasks. (For example, online Bayesian inference is of this type, since the posterior for the parameters is derived from the new data and the past prior; see e.g., the **elastic weight consolidation** method discussed in Section 17.5.1, or the **variational continual learning** method discussed in Supplementary Section 10.2). The second is **memory methods**, which rely on some kind of **experience replay** or **rehearsal** of past data (see e.g., [Hen+21]), or some kind of generative model of past data. The third is **architectural methods**, that add capacity to the network whenever a task boundary is encountered, such as a new class label (see e.g., [Rus+16]).

Of course, these techniques can be combined. For example, we can create a semi-parametric model, in which we store some past data (exemplars) while also learning parameters online in a Bayesian (regularized) way (see e.g., [Kur+20]). The “right” method depends, as usual, on what inductive bias you want to use, and want your computational budget is in terms of time and memory.

19.7.5 Online learning

The problem of **online learning** is similar to continual learning, except the loss metric is different, and we usually assume that learning and evaluation occur at each step. More precisely, we assume the data generating distribution, $p_t^*(\mathbf{x}, \mathbf{y}) = p(\mathbf{x}|\phi_t)p(\mathbf{y}|\mathbf{x}, \mathbf{w}_t)$, evolves over time, as shown in Figure 19.22. At each step t nature generates a data sample, $(\mathbf{x}_t, \mathbf{y}_t) \sim p_t^*$. The agent sees \mathbf{x}_t and is asked to predict \mathbf{y}_t by computing the posterior predictive distribution

$$\hat{p}_{t|t-1} = p(\mathbf{y}|\mathbf{x}_t, \mathcal{D}_{1:t-1}) \quad (19.43)$$

where $\mathcal{D}_{1:t-1} = \{(\mathbf{x}_s, \mathbf{y}_s) : s = 1 : t-1\}$ is all past data. It then incurs a loss of

$$\mathcal{L}_t = \ell(\hat{p}_{t|t-1}, \mathbf{y}_t) \quad (19.44)$$

For classification problems, we often use 0-1 loss, $\mathcal{L}_t = \mathbb{I}(\hat{y}_t \neq y_t)$. In this case the optimal action is to predict the most probable label, $\hat{y}_t = \text{argmax}_c p(y=c|\mathbf{x}_t, \mathcal{D}_{1:t-1})$.

In contrast to the continual learning scenarios studied above, the loss incurred at each step is what matters, rather than loss on a fixed test set. That is, we want to minimize

$$\mathcal{L} = \sum_{t=1}^T \mathcal{L}_t \quad (19.45)$$

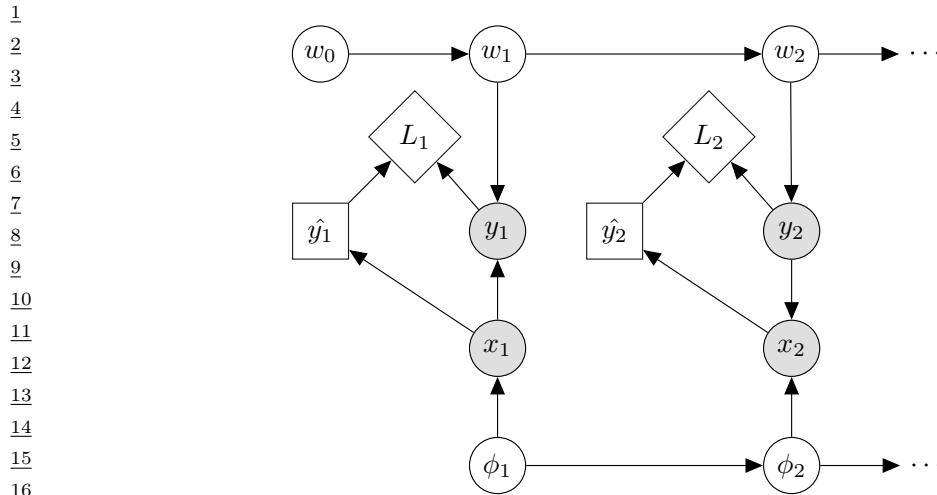


Figure 19.22: Online learning illustrated as an influence diagram (Section 34.2). Here $\hat{y}_t = \operatorname{argmax}_y p(y|\mathbf{x}_t, \mathcal{D}_{1:t-1})$ is the action (MAP predicted output) at time t , and $L_t = \ell(y_t, \hat{y}_t)$ is the corresponding loss (utility) function. The data at time t is assumed to come from a model of the form $p(\mathbf{x}_t, y_t) = p(\mathbf{x}_t|\phi_t)p(y_t|\mathbf{x}_t, \mathbf{w}_t)$. The parameters of the world model can change arbitrarily over time.

Since it is hard to interpret this number, it is common to compare it to the optimal value one could have obtained in hindsight. This yields a quantity called the **regret**:

$$\text{regret} = \sum_{t=1}^T [\ell(\hat{p}_{t|t-1}, \mathbf{y}_t) - \ell(\hat{p}_{t|T}, \mathbf{y}_t)] \quad (19.46)$$

where $\hat{p}_{t|t-1} = p(\mathbf{y}|\mathbf{x}_t, \mathcal{D}_{1:t-1})$ is the online prediction, and $\hat{p}_{t|T} = p(\mathbf{y}|\mathbf{x}_t, \mathcal{D}_{1:T})$ is the optimal estimate at the end of training. It is possible to convert bounds on regret, which are backwards looking, into bounds on risk (i.e., expected future loss), which is forwards looking. See [HT15] for details.

Online learning is very useful for decision and control problems, such as multi-armed bandits (Section 34.4) and reinforcement learning (see Chapter 35), where the agent “lives forever”, and where there is no fixed training phase followed by a test phase. (See e.g., Section 17.5 where we discuss online Bayesian inference for neural networks.)

The previous continual learning scenarios can be derived as special cases of online learning, by defining a suitable sequence of distributions, and by requiring the agent to either train or test at each step on a suitable minibatch of data. (We leave the details of this mapping as an exercise to the reader.)

19.8 Adversarial examples

This section is coauthored with Justin Gilmer.

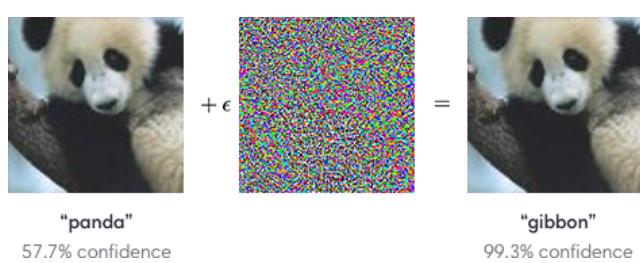


Figure 19.23: Example of an adversarial attack on an image classifier. Left column: original image which is correctly classified. Middle column: small amount of structured noise which is added to the input (magnitude of noise is magnified by 10 \times). Right column: new image, which is confidently misclassified as a “gibbon”, even though it looks just like the original “panda” image. Here $\epsilon = 0.007$. From Figure 1 of [GSS15]. Used with kind permission of Ian Goodfellow.

In Section 19.2, we discussed what happens to a predictive model when the input distribution shifts for some reason. In this section, we consider the case where an adversary deliberately chooses inputs to minimize the performance of a predictive model. That is, suppose an input \mathbf{x} is classified as belonging to class c . We then choose a new input \mathbf{x}_{adv} which minimizes the probability of this label, subject to the constraint that \mathbf{x}_{adv} is “perceptually similar” to the original input \mathbf{x} . This gives rise to the following objective:

$$\mathbf{x}_{\text{adv}} = \underset{\mathbf{x}' \in \Delta(\mathbf{x})}{\operatorname{argmin}} \log p(y = c | \mathbf{x}') \quad (19.47)$$

where $\Delta(\mathbf{x})$ is the set of images that are “similar” to \mathbf{x} (we discuss different notions of similarity below).

Equation (19.47) is an example of an **adversarial attack**. We illustrate this in Figure 19.23. The input image \mathbf{x} is on the left, and is predicted to be a panda with probability 57%. By adding a tiny amount of carefully chosen noise (shown in the middle) to the input, we generate the **adversarial image** \mathbf{x}_{adv} on the right: this “looks like” the input, but is now classified as a gibbon with probability 99%.

The ability to create adversarial images was first noted in [Sze+14]. It is surprisingly easy to create such examples, which seems paradoxical, given the fact that modern classifiers seem to work so well on normal inputs, and the perturbed images “look” the same to humans. We explain this paradox in Section 19.8.5.

The existence of adversarial images also raises security concerns. For example, [Sha+16] showed they could force a face recognition system to misclassify person A as person B , merely by asking person A to wear a pair of sunglasses with a special pattern on them, and [Eyk+18] show that it is possible to attach small “**adversarial stickers**” to traffic signs to classify stop signs as speed limit signs.

Below we briefly discuss how to create adversarial attacks, why they occur, and how we can try to defend against them. We focus on the case of deep neural nets for images, although it is important to note that many other kinds of models (including logistic regression and generative models) can also suffer from adversarial attacks. Furthermore, this is not restricted to the image domain, but occurs with many kinds of high dimensional inputs. For example, [Li+19] contains an audio attack

¹ and [Dal+04; Jia+19] contains a text attack. More details on adversarial examples can be found in
² e.g., [Wiy+19; Yua+19].
³

⁵ 19.8.1 Whitebox (gradient-based) attacks

⁷ To create an adversarial example, we must find a “small” perturbation $\boldsymbol{\delta}$ to add to the input \mathbf{x} to
⁸ create $\mathbf{x}_{\text{adv}} = \mathbf{x} + \boldsymbol{\delta}$ so that $f(\mathbf{x}_{\text{adv}}) = y'$, where $f()$ is the classifier, and y' is the label we want to
⁹ force the system to output. This is known as a **targeted attack**. Alternatively, we may just want
¹⁰ to find a perturbation that causes the current predicted label to change from its current value to any
¹¹ other value, so that $f(\mathbf{x} + \boldsymbol{\delta}) \neq f(\mathbf{x})$, which is known as **untargeted attack**.

¹² In general, we define the objective for the adversary as *maximizing* the following loss:

$$\mathbf{x}_{\text{adv}} = \underset{\mathbf{x}' \in \Delta(\mathbf{x})}{\operatorname{argmax}} \mathcal{L}(\mathbf{x}', y; \boldsymbol{\theta}) \quad (19.48)$$

¹⁶ where y is the true label. For the untargeted case, we can define $\mathcal{L}(\mathbf{x}', y; \boldsymbol{\theta}) = -\log p(y|\mathbf{x}')$, so we
¹⁷ minimize the probability of the true label; and for the targeted case, we can define $\mathcal{L}(\mathbf{x}', y; \boldsymbol{\theta}) =$
¹⁸ $\log p(y'|\mathbf{x}')$, where we maximize the probability of the desired label $y' \neq y$.

¹⁹ To define what we mean by “small” perturbation, we impose the constraint that $\mathbf{x}_{\text{adv}} \in \Delta(\mathbf{x})$,
²⁰ which is the set of “perceptually similar” images to the input \mathbf{x} . Most of the literature has focused
²¹ on a simplistic setting in which the adversary is restricted to making bounded l_p perturbations of a
²² clean input \mathbf{x} , that is

$$\Delta(\mathbf{x}) = \{\mathbf{x}' : \|\mathbf{x}' - \mathbf{x}\|_p < \epsilon\} \quad (19.49)$$

²⁶ Typically people assume $p = 1$ or $p = 0$. We will discuss more realistic threat models in Section 19.8.3.
²⁷ In this section, we assume that the attacker knows the model parameters $\boldsymbol{\theta}$; this is called a
²⁸ **whitebox attack**, and lets us use gradient based optimization methods. We relax this assumption
²⁹ in Section 19.8.2.)

³⁰ To solve the optimization problem in Equation (19.48), we can use any kind of constrained
³¹ optimization method. In [Sze+14] they used bound-constrained BFGS. [GSS15] proposed the more
³² efficient **fast gradient sign (FGS)** method, which performs iterative updates of the form

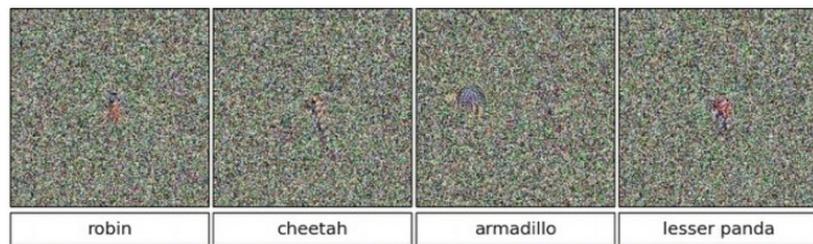
$$\mathbf{x}_{t+1} = \mathbf{x}_t + \boldsymbol{\delta}_t \quad (19.50)$$

$$\boldsymbol{\delta}_t = \epsilon \operatorname{sign}(\nabla_{\mathbf{x}} \log p(y'|\mathbf{x}, \boldsymbol{\theta})|_{\mathbf{x}_t}) \quad (19.51)$$

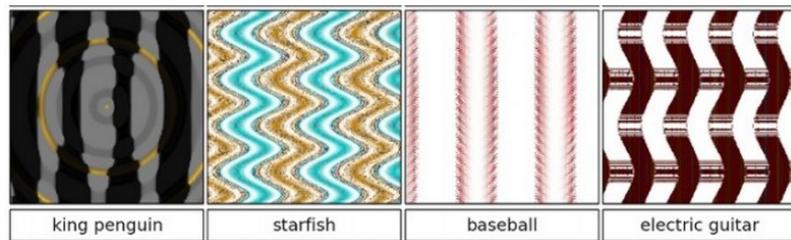
³⁷ where $\epsilon > 0$ is a small learning rate. (Note that this gradient is with respect to the input pixels, not
³⁸ the model parameters.) Figure 19.23 gives an example of this process.

³⁹ More recently, [Mad+18] proposed the more powerful **projected gradient descent (PGD)**
⁴⁰ attack; this can be thought of as an iterated version of FGS. There is no “best” variant of PGD for
⁴¹ solving 19.48. Instead, what matters more is the implementation details, e.g. how many steps are
⁴² used, the step size, and the exact form of the loss. To avoid local minima, we may use random restarts,
⁴³ choosing random points in the constraint space Δ to initialize the optimization. The algorithm
⁴⁴ should be carefully tuned to the specific problem, and the loss should be monitored to check for
⁴⁵ optimization issues. For best practices, see [Car+19].

⁴⁶



10 *Figure 19.24: Images that look like random noise but which cause the CNN to confidently predict a specific*
 11 *class. From Figure 1 of [NYC15]. Used with kind permission of Jeff Clune.*



21 *Figure 19.25: Synthetic images that cause the CNN to confidently predict a specific class. From Figure 1 of*
 22 *[NYC15]. Used with kind permission of Jeff Clune.*

26 19.8.2 Blackbox (gradient-free) attacks

27 In this section, we no longer assume that the adversary knows the parameters θ of the predictive model
 28 f . This is known as a **black box attack**. In such cases, we must use derivative-free optimization
 29 (DFO) methods (see Section 6.9).

30 Evolutionary algorithms (EA) are one class of DFO solvers. These were used in [NYC15] to create
 31 blackbox attacks. Figure 19.24 shows some images that were generated by applying an EA to a
 32 random noise image. These are known as **fooling images**, as opposed to adversarial images, since
 33 they are not visually realistic. Figure 19.25 shows some fooling images that were generated by
 34 applying EA to the parameters of a compositional pattern-producing network (CPPN) [Sta07].⁷ By
 35 suitably perturbing the CPPN parameters, it is possible to generate structured images with high
 36 fitness (classifier score), but which do not look like natural images [Aue12].

37 In [SVK19], they used differential evolution to attack images by modifying a single pixel. This is
 38 equivalent to bounding the ℓ_0 norm of the perturbation, so that $\|\mathbf{x}_{\text{adv}} - \mathbf{x}\|_0 = 1$.

39 In [Pap+17], they learned a differentiable surrogate model of the blackbox, by just querying its
 40 predictions y for different inputs \mathbf{x} . They then used gradient-based methods to generate adversarial
 41 attacks on their surrogate model, and then showed that these attacks transferred to the real model.
 42 In this way, they were able to attack various the image classification APIs of various cloud service

44 7. A CPPN is a set of elementary functions (such as linear, sine, sigmoid, and Gaussian) which can be composed in
 45 order to specify the mapping from each coordinate to the desired color value. CPPN was originally developed as a way
 46 to encode abstract properties such as symmetry and repetition, which are often seen during biological development.

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



14 *Figure 19.26: An adversarially modified image to evade spam detectors. The image is constructed from*
15 *scratch, and does not involve applying a small perturbation to any given image. This is an illustrative example*
16 *of how large the space of possible adversarial inputs Δ can be when the attacker has full control over the input.*
17 *From [Big+11]. Used with kind permission of Battista Biggio.*

18
19
20

21 providers, including Google, Amazon and MetaMind.

22

23 19.8.3 Real world adversarial attacks

24

25 Typically, the space of possible adversarial inputs Δ can be quite large, and will be difficult to exactly
26 define mathematically as it will depend on semantics of the input based on the attacker's goals
27 [BR18]. (The set of variations Δ that we want the model to be invariant to is called the **threat**
28 **model**.)

29 Consider for example of the content constrained threat model discussed in [Gil+18a]. One instance
30 of this threat model involves image spam, where the attacker wishes to upload an image attachment
31 in an email that will not be classified as spam by a detection model. In this case Δ is incredibly
32 large as it consists of all possible images which contain some semantic concept the attacker wishes to
33 upload (in this case an advertisement). To explore Δ , spammers can utilize different fonts, word
34 orientations or add random objects to the background as is the case of the adversarial example in
35 Figure 19.26 (see [Big+11] for more examples). Of course, optimization based methods may still be
36 used here to explore parts of Δ . However, in practice it may be preferable to design an adversarial
37 input by hand as this can be significantly easier to execute with only limited-query black-box access
38 to the underlying classifier.

39

40 19.8.4 Defenses based on robust optimization

41

42 As discussed in Section 19.8.3, securing a system against adversarial inputs in more general threat
43 models seems extraordinarily difficult, due to the vast space of possible adversarial inputs Δ . However,
44 there is a line of research focused on producing models which are invariant to perturbations within
45 a small constraint set $\Delta(\mathbf{x})$, with a focus on l_p -robustness where $\Delta(\mathbf{x}) = \{\mathbf{x}' : \|\mathbf{x} - \mathbf{x}'\|_p < \epsilon\}$.
46 Although solving this toy threat model has little application to security settings, enforcing smoothness

47

priors have in some cases improved robustness to random image corruptions [SHS], lead to models which transfer better [Sal+20], and can bias models towards different features in the data [Yin+19a].

Perhaps the most straightforward method for improving l_p -robustness is to directly optimize for it through **robust optimization** [BTEGN09], also known as **adversarial training** [GSS15]. We define the **adversarial risk** to be

$$\min_{\theta} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim p(\mathbf{x}, \mathbf{y})} \left[\max_{\mathbf{x}' \in \Delta(\mathbf{x})} L(\mathbf{x}', \mathbf{y}; \theta) \right] \quad (19.52)$$

The min max formulation in equation 19.52 poses unique challenges from an optimization perspective—it requires solving both the non-concave inner maximization and the non-convex outer minimization problems. Even worse, the inner max is NP-hard to solve in general [Kat+17]. However, in practice it may be sufficient to compute the gradient of the outer objective $\nabla_{\theta} L(\mathbf{x}_{\text{adv}}, \mathbf{y}; \theta)$ at an approximately maximal point in the inner problem $\mathbf{x}_{\text{adv}} \approx \operatorname{argmax}_{\mathbf{x}} L(\mathbf{x}_{\text{adv}}, \mathbf{y}; \theta)$ [Mad+18]. Currently, best practice is to approximate the inner problem using a few steps of PGD.

Other methods seek to **certify** that a model is robust within a given region $\Delta(x)$. One method for certification uses randomized smoothing [CRK19]—a technique for converting a model robust to random noise into a model which is provably robust to bounded worst-case perturbations in the l_2 -metric. Another class of methods applies specifically for networks with ReLU activations, leveraging the property that the model is locally linear, and that certifying in region defined by linear constraints reduces to solving a series of linear programs, for which standard solvers can be applied [WK18].

19.8.5 Why models have adversarial examples

The existence of adversarial inputs is paradoxical, since modern classifiers seem to do so well on normal inputs. However, the existence of adversarial examples is a natural consequence of the general lack of robustness to distribution shift discussed in Section 19.2. To see this, suppose a model’s accuracy drops on some shifted distribution of inputs $p_{\text{te}}(\mathbf{x})$ that differs from the training distribution $p_{\text{tr}}(\mathbf{x})$; in this case, the model will necessarily be vulnerable to an adversarial attack: if errors exist, there must be a nearest such error. Furthermore, if the input distribution is high dimensional, then we should expect the nearest error to be significantly closer than errors which are sampled randomly from some out-of-distribution $p_{\text{te}}(\mathbf{x})$.

A cartoon illustration of what is going on is shown in Figure 19.27a, where \mathbf{x}_0 is the clean input image, B is an image corrupted by Gaussian noise, and A is an adversarial image. If we assume a linear decision boundary, then the error set E is a half space a certain distance from \mathbf{x}_0 . We can relate the distance to the decision boundary $d(\mathbf{x}_0, E)$ with the error rate in noise at some input \mathbf{x}_0 , denoted by $\mu = \mathbb{P}_{\delta \sim N(0, \sigma I)} [\mathbf{x}_0 + \delta \in E]$. With a linear decision boundary the relationship between these two quantities is determined by

$$d(\mathbf{x}_0, E) = -\sigma \Phi^{-1}(\mu) \quad (19.53)$$

where Φ^{-1} denotes the inverse cdf of the gaussian distribution. When the input dimension is large, this distance will be significantly smaller than the distance to a randomly sampled noisy image $\mathbf{x}_0 + \delta$ for $\delta \sim N(0, \sigma I)$, as the noise term will with high probability have norm $\|\delta\|_2 \approx \sigma \sqrt{d}$. As a concrete example consider the ImageNet dataset, where $d = 224 \times 224 \times 3$ and suppose we set $\sigma = .2$.

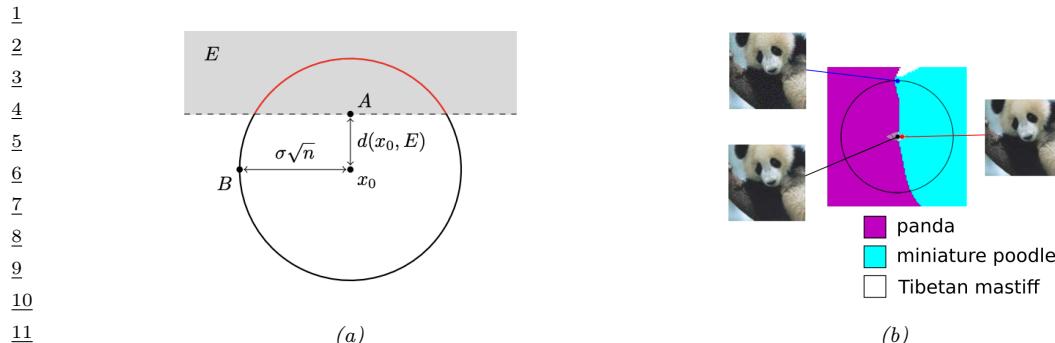


Figure 19.27: (a) When the input dimension n is large and the decision boundary is locally linear, even a small error rate in random noise will imply the existence of small adversarial perturbations. Here, $d(\mathbf{x}_0, E)$ denotes the distance from a clean input \mathbf{x}_0 to an adversarial example (A) while the distance from \mathbf{x}_0 to a random sample $N(0; \sigma^2 I)$ (B) will be approximately $\sigma\sqrt{n}$. As $n \rightarrow \infty$ the ratio of $d(\mathbf{x}_0, A)$ to $d(\mathbf{x}_0, B)$ goes to 0. (b) A 2d slice of the InceptionV3 decision boundary through three points: a clean image (black), an adversarial example (red), and an error in random noise (blue). The adversarial example and the error in noise lie in the same region of the error set which is misclassified as “miniature poodle”, which closely resembles a halfspace as in Figure (a). Used with kind permission of Justin Gilmer.

²² Then if the error rate in noise is just $\mu = .01$, equation 19.53 will imply that $d(\mathbf{x}_0, E) = .5$. Thus the
²³ distance to an adversarial example will be more than 100 times closer than the distance to a typical
²⁴ noisy images, which will be $\sigma\sqrt{d} \approx 77.6$. This phenomenon of small volume error sets being close
²⁵ to most points in a data distribution $p(\mathbf{x})$ is called **concentration of measure**, and is a property
²⁶ common among many high dimensional data distributions [MDM19; Gil+18b].

²⁷ In summary, although the existence of adversarial examples is often discussed as an unexpected
²⁸ phenomenon, there is nothing special about the existence of worst-case errors for ML classifiers—they
²⁹ will always exist as long as errors exist.

PART IV

Generation

20 Generative models: an overview

20.1 Introduction

A **generative model** is a joint probability distribution $p(\mathbf{x})$, for $\mathbf{x} \in \mathcal{X}$. In some cases, the model may be conditioned on inputs or covariates $\mathbf{c} \in \mathcal{C}$, which gives rise to a **conditional generative model** of the form $p(\mathbf{x}|\mathbf{c})$.

There are many kinds of generative model. We give a brief summary in Section 20.2, and go into more detail in subsequent chapters. See also [Tom22] for a recent book on this topic that goes into more depth.

20.2 Types of generative model

There are many kinds of generative model, some of which we list in Table 20.1. At a high level, we can distinguish between **deep generative models** (DGM) — which use deep neural network to learn a complex mapping from a single latent vector \mathbf{z} to the observed data \mathbf{x} — and more “classical” **probabilistic graphical models** (PGM), that map a set of interconnected latent variables $\mathbf{z}_1, \dots, \mathbf{z}_L$ to the observed variables $\mathbf{x}_1, \dots, \mathbf{x}_D$ using simpler, often linear, mappings. Of course, many hybrids are possible. For example, PGMs can use neural networks, and DGMs can use structured state spaces. We discuss PGMs in general terms in Chapter 4, and give examples in Chapter 28, Chapter 29, Chapter 30. In this part of the book, we mostly focus on DGMs.

The main kinds of DGM are: **variational autoencoders (VAE)**, **autoregressive (AR)** models, **normalizing flows**, **diffusion models**, **energy based models (EBM)**, and **generative adversarial networks (GAN)**. We can categorize these models in terms of the following criteria (see Figure 20.1 for a visual summary):

- Density: does the model support pointwise evaluation of the probability density function $p(\mathbf{x})$, and if so, is this fast or slow, exact, approximate or a bound, etc? For **implicit models**, such as GANs, there is no well-defined density $p(\mathbf{x})$. For other models, we can only compute a lower bound on the density (VAEs), or an approximation to the density (EBMs, UPGMs).
- Sampling: does the model support generating new samples, $\mathbf{x} \sim p(\mathbf{x})$, and if so, is this fast or slow, exact or approximate? Directed PGMs, VAEs and GANs all support fast sampling. However, undirected PGMs, EBMs, AR, diffusion and flows are slow for sampling.
- Training: what kind of method is used for parameter estimation? For some models (such as AR, flows and directed PGMs), we can perform exact maximum likelihood estimation (MLE), although

1
2
3
4
5
6
7
8
9

Model	Chapter	Density	Sampling	Training	Latents	Architecture
PGM-D	Section 4.2	Exact, fast	Fast	MLE	Optional	Sparse DAG
PGM-U	Section 4.3	Approx, slow	Slow	MLE-A	Optional	Sparse graph
VAE	Chapter 21	LB, fast	Fast	MLE-LB	\mathbb{R}^L	Encoder-Decoder
AR	Chapter 22	Exact, fast	Slow	MLE	None	Sequential
Flows	Chapter 23	Exact, slow/fast	Slow	MLE	\mathbb{R}^D	Invertible
EBM	Chapter 24	Approx, slow	Slow	MLE-A	Optional	Discriminative
Diffusion	Chapter 25	LB	Slow	MLE-LB	\mathbb{R}^D	Encoder-Decoder
GAN	Chapter 26	NA	Fast	Min-max	\mathbb{R}^L	Generator-Discriminator

10 *Table 20.1: Characteristics of common kinds of generative model. Here D is the dimensionality of the observed
11 \mathbf{x} , and L is the dimensionality of the latent \mathbf{z} , if present. (We usually assume $L \ll D$, although overcomplete
12 representations can have $L \gg D$.) Abbreviations: Approx = approximate, AR = autoregressive, EBM =
13 Energy Based Model, GAN = generative adversarial network, MLE = maximum likelihood estimation, MLE-A
14 = MLE (Approximate), MLE-LB = MLE (Lower Bound), NA = not available, PGM = probabilistic graphical
15 model, PGM-D = directed PGM, PGM-U = undirected PGM, VAE = variational autoencoder.*

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

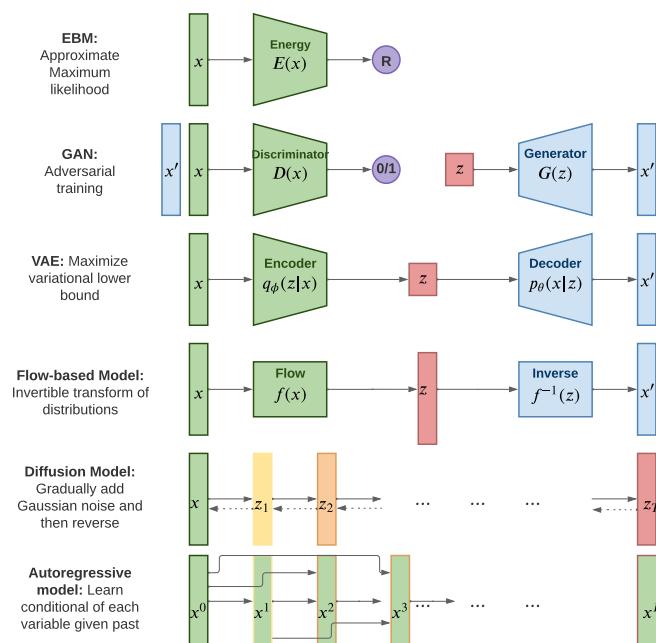
43

44

45

46

47



41 *Figure 20.1: Summary of various kinds of deep generative model. Here \mathbf{x} is the observed data, \mathbf{z} is the latent
42 code, and \mathbf{x}' is a sample from the model. AR models do not have a latent code \mathbf{z} . For diffusion models and
43 flow models, the size of \mathbf{z} is the same as \mathbf{x} . For AR models, x^d is the d 'th dimension of \mathbf{x} . R represents
44 real-valued output, 0/1 represents binary output. Adapted from Figure 1 of [Wen21].*



Figure 20.2: Synthetic faces from a score-based generative model (Section 24.3.5). From Figure 12 of [Song+21]. Used with kind permission of Yang Song.

the objective is usually non-convex, so we can only reach a local optimum. For other models, we cannot tractably compute the likelihood. In the case of VAEs, we maximize a lower bound on the likelihood; in the case of EBMs and UGMs, we maximize an approximation to the likelihood. For GANs we have to use min-max training, which can be unstable, and there is no clear objective function to monitor.

- Latents: does the model use a latent vector \mathbf{z} to generate \mathbf{x} or not, and if so, is it the same size as \mathbf{x} or is it a potentially compressed representation? For example, AR models do not use latents; flows and diffusion use latents, but they are not compressed.¹ Graphical models, including EBMs, may or may not use latents.
- Architecture: what kind of neural network should we use, and are there restrictions? For flows, we are restricted to using invertible neural networks where each layer has a tractable Jacobian. For EBMs, we can use any model we like. The other models have different restrictions.

20.3 Goals of generative modeling

There are several different kinds of tasks that we can use generative models for, as we discuss below.

20.3.1 Generating data

One of the main goals of generative models is to generate (create) new data samples. For example, if we fit a model $p(\mathbf{x})$ to images of faces, we can sample new faces from it, as illustrated in Figure 20.2.²

1. Flow models define a latent vector \mathbf{z} that has the same size as \mathbf{x} , although the internal deterministic computation may use vectors that are larger or smaller than the input (see e.g., the DenseFlow paper [GGS21]).
 2. These images were made with a technique called score-based generative modeling (Section 24.3.5), although similar results can be obtained using many other techniques. See for example <https://this-person-does-not-exist.com/en> which shows results from a GAN model (Chapter 26).



(a) Teddy bears swimming at the Olympics 400m Butterfly event.



(b) A cute corgi lives in a house made out of sushi.



(c) A cute sloth holding a small treasure chest. A bright golden glow is coming from the chest.

Figure 20.3: Some 1024×1024 images generated from text prompts by the Imagen diffusion model (Section 25.2.3). From Figure 1 of [Sah+22]. Used with kind permission of William Chan.



(a)



(b)



(c)

Figure 20.4: Converting a gray-scale image (left) to a colorized one (middle) using a conditional GAN (Chapter 26). Right column shows the original image. (Photo depicts author’s wife, Margaret Murphy, photo taken July 2022.) Generated by <https://github.com/jantic/DeOldify>.

*Similar methods can be used to create samples of text, audio, etc. When this technology is abused to make fake content, they are called **deep fakes**. For a review of this topic, see e.g., [Ngu+19].*

*To control what is generated, it is useful to use a **conditional generative model** of the form $p(\mathbf{x}|\mathbf{c})$. Here are some examples:*

- \mathbf{c} = text prompt, \mathbf{x} = image. This is a **text-to-image** model (see Figure 20.3 and Figure 22.8 for examples).*
- \mathbf{c} = image, \mathbf{x} = text. This is an **image-to-text** model, which is useful for **image captioning**.*
- \mathbf{c} = image, \mathbf{x} = image. This is an **image-to-image** model (see Figure 20.4 and Figure 25.4 for examples).*
- \mathbf{c} = sequence of sounds, \mathbf{x} = sequence of words. This is a **speech-to-text** model, which is useful for **automatic speech recognition (ASR)**.*

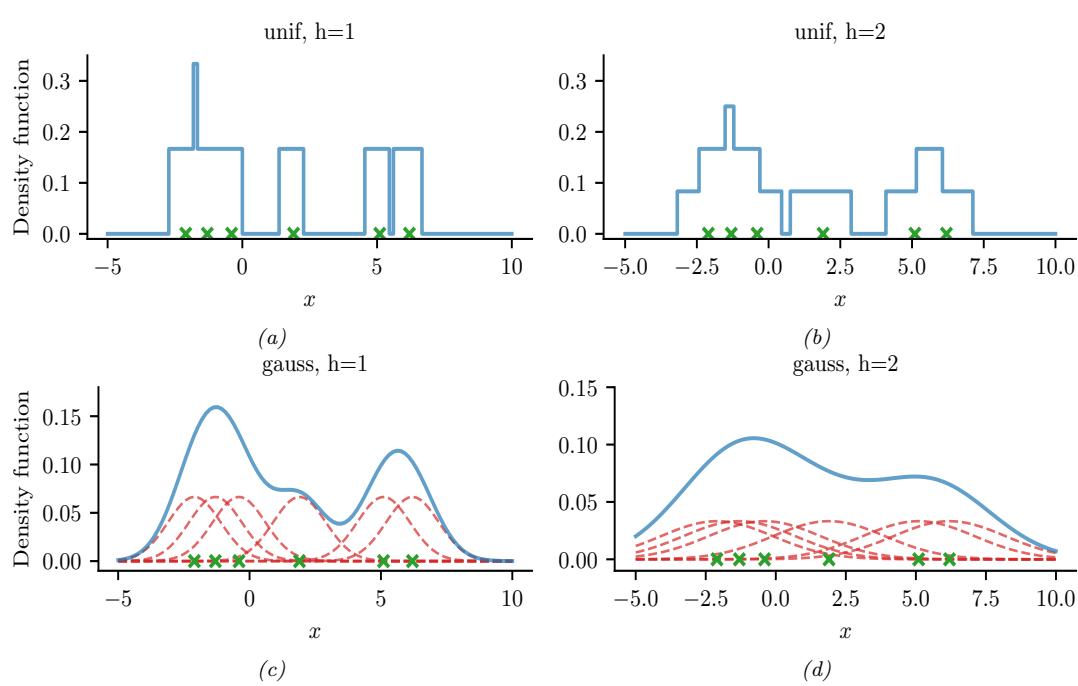


Figure 20.5: A nonparametric (Parzen) density estimator in 1d estimated from 6 data points, denoted by x . Top row: uniform kernel. Bottom row: Gaussian kernel. Left column: bandwidth parameter $h = 1$. Right column: bandwidth parameter $h = 2$. Adapted from http://en.wikipedia.org/wiki/Kernel_density_estimation. Generated by `parzen_window_demo.ipynb`.

- \mathbf{c} = sequence of English words, \mathbf{x} = sequence of French words. This is a **sequence-to-sequence** model, which is useful for **machine translation**.
- \mathbf{c} = initial prompt, \mathbf{x} = continuation of the text. This is another sequence-to-sequence model, which is useful for automatic **text generation** (see Figure 22.4 for an example).

Note that, in the conditional case, we sometimes denote the inputs by \mathbf{x} and the outputs by \mathbf{y} . In this case the model has the familiar form $p(\mathbf{y}|\mathbf{x})$. In the special case that \mathbf{y} denotes a low dimensional quantity, such as a integer class label, $y \in \{1, \dots, C\}$, we get a predictive (discriminative) model. The main difference between a discriminative model and a conditional generative model is this: in a discriminative model, we assume there is one correct output, whereas in a conditional generative model, we assume there may be multiple correct outputs. This makes it harder to evaluate generative models, as we discuss in Section 20.4.

20.3.2 Density estimation

The task of **density estimation** refers to evaluating the probability of an observed data vector, i.e., computing $p(\mathbf{x})$. This can be useful for outlier detection (Section 19.3.2), data compression (Section 5.4), generative classifiers, model comparison, etc.

	Data sample	Variables			Missing values replaced by means		
		A	B	C	A	B	C
1	1	6	6	NA	2	6	7.5
2	2	NA	6	0	9	6	0
3	3	NA	6	NA	9	6	7.5
4	4	10	10	10	10	10	10
5	5	10	10	10	10	10	10
6	6	10	10	10	10	10	10
12	Average	9	8	7.5	9	8	7.5

Figure 20.6: Missing data imputation using the mean of each column.

A simple approach to this problem, which works in low dimensions, is to use **kernel density estimation** or **KDE**, which has the form

$$p(\mathbf{x}|\mathcal{D}) = \frac{1}{N} \sum_{n=1}^N \mathcal{K}_h(\mathbf{x} - \mathbf{x}_n) \quad (20.1)$$

Here $\mathcal{D} = \{\mathbf{v}, \mathbf{x}_1, \dots, \mathbf{x}_N\}$ is the data, and \mathcal{K}_h is a density kernel with **bandwidth** h , which is a function $\mathcal{K} : \mathbb{R} \rightarrow \mathbb{R}_+$ such that $\int \mathcal{K}(x)dx = 1$ and $\int x\mathcal{K}(x)dx = 0$. We give a 1d example of this in Figure 20.5: in the top row, we use a uniform (boxcar) kernel, and in the bottom row we use a Gaussian kernel.

In higher dimensions, KDE suffers from the **curse of dimensionality** (see e.g., [AHK01]), and we need to use parametric density models $p_\theta(\mathbf{x})$ of some kind.

20.3.3 Imputation

The task of **imputation** refers to “filling in” missing values of a data vector or data matrix. For example, suppose \mathbf{X} is an $N \times D$ matrix of data (think of a spreadsheet) in which some entries, call them \mathbf{X}_m , may be missing, while the rest, \mathbf{X}_o , are observed. A simple way to fill in the missing data is to use the mean value of each feature, $\mathbb{E}[x_d]$; this is called **mean value imputation**, and is illustrated in Figure 20.6. However, this ignores dependencies between the variables within each row, and does not return any measure of uncertainty.

We can generalize this by fitting a generative model to the observed data, $p(\mathbf{X}_o)$, and then computing samples from $p(\mathbf{X}_m|\mathbf{X}_o)$. This is called **multiple imputation**. A generative model can be used to fill in more complex data types, such as **in-painting** occluded pixels in an image (see Figure 25.4).

See Section 21.3.4 for a more general discussion of missing data.

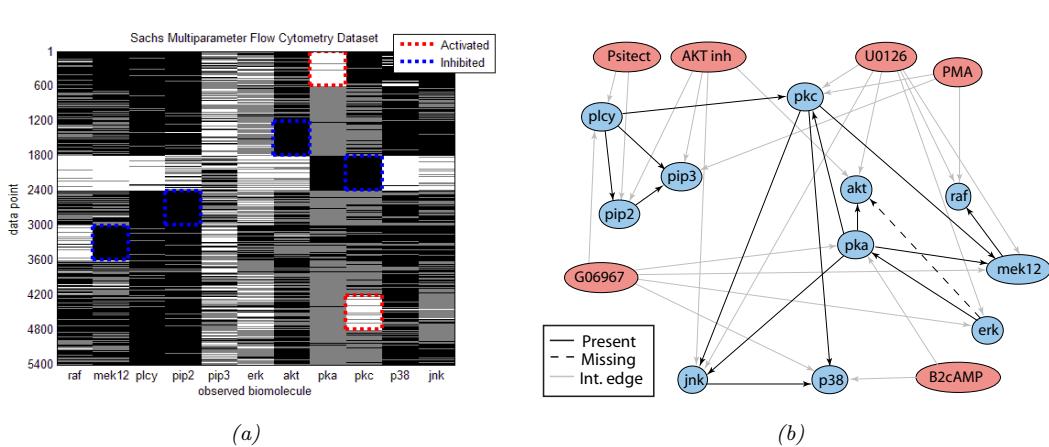


Figure 20.7: (a) A design matrix consisting of 5400 data points (rows) measuring the state (using flow cytometry) of 11 proteins (columns) under different experimental conditions. The data has been discretized into 3 states: low (black), medium (grey) and high (white). Some proteins were explicitly controlled using activating or inhibiting chemicals. (b) A directed graphical model representing dependencies between various proteins (blue circles) and various experimental interventions (pink ovals), which was inferred from this data. We plot all edges for which $p(G_{ij} = 1|\mathcal{D}) > 0.5$. Dotted edges are believed to exist in nature but were not discovered by the algorithm (1 false negative). Solid edges are true positives. The light colored edges represent the effects of intervention. From Figure 6d of [EM07].

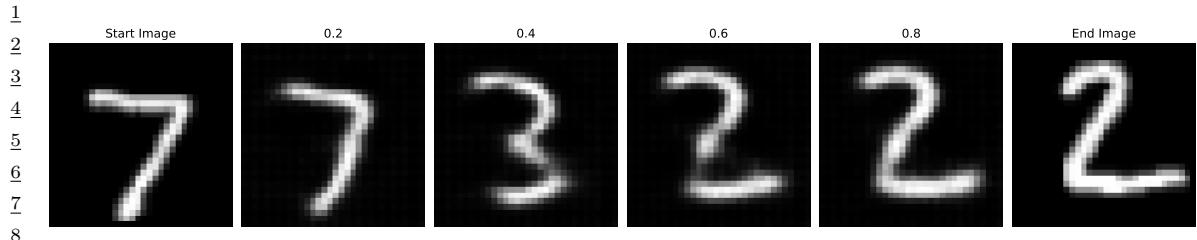
20.3.4 Structure discovery

Some kinds of generative models have latent variables \mathbf{z} , which are assumed to be the “causes” that generated the observed data \mathbf{x} . We can use Bayes rule to invert the model to compute $p(\mathbf{z}|\mathbf{x}) \propto p(\mathbf{z})p(\mathbf{x}|\mathbf{z})$. This can be useful for discovering latent, low-dimensional patterns in the data.

For example, suppose we perturb various proteins in a cell and measure the resulting phosphorylation state using a technique known as flow cytometry, as in [Sac+05]. An example of such a dataset is shown in Figure 20.7(a). Each row represents a data sample $\mathbf{x}_n \sim p(\cdot|\mathbf{a}_n, \mathbf{z})$, where $\mathbf{x} \in \mathbb{R}^{11}$ is a vector of outputs (phosphorylations), $\mathbf{a} \in \{0, 1\}^6$ is a vector of input actions (perturbations) and \mathbf{z} is the unknown cellular signaling network structure. We can infer the graph structure $p(\mathbf{z}|\mathcal{D})$ using graphical model structure learning techniques (see Section 30.3). In particular, we can use the dynamic programming method described in [EM07] to get the result is shown in Figure 20.7(b). Here we plot the median graph, which includes all edges for which $p(z_{ij} = 1|\mathcal{D}) > 0.5$. (For a more recent approach to this problem, see e.g., [Bro+20b].)

20.3.5 Latent space interpolation

One of the most interesting abilities of certain latent variable models is the ability to generate samples that have certain desired properties by interpolating between existing data points in latent space. To explain how this works, let \mathbf{x}_1 and \mathbf{x}_2 be two inputs (e.g. images), and let $\mathbf{z}_1 = e(\mathbf{x}_1)$ and $\mathbf{z}_2 = e(\mathbf{x}_2)$ be their latent encodings. (The method used for computing these will depend on the type of model; we discuss the details in later chapters.) We can regard \mathbf{z}_1 and \mathbf{z}_2 as two “anchors” in



10 *Figure 20.8: Interpolation between two MNIST images in the latent space of a β -VAE (with $\beta = 0.5$).*
 11 *Generated by [mnist_vae_ae_comparison.ipynb](#).*



21 *Figure 20.9: Interpolation between two CelebA images in the latent space of a β -VAE (with $\beta = 0.5$).*
 22 *Generated by [celeba_vae_ae_comparison.ipynb](#).*

23
 24 latent space. We can now generate new images that interpolate between these points by computing
 25 $\mathbf{z} = \lambda \mathbf{z}_1 + (1 - \lambda) \mathbf{z}_2$, where $0 \leq \lambda \leq 1$, and then decoding by computing $\mathbf{x}' = d(\mathbf{z})$, where $d()$ is the
 26 decoder. This is called **latent space interpolation**, and will generate data that combines semantic
 27 features from both \mathbf{x}_1 and \mathbf{x}_2 . (The justification for taking a linear interpolation is that the learned
 28 manifold often has approximately zero curvature, as shown in [SKTF18]. However, sometimes it is
 29 better to use nonlinear interpolation [Whi16; MB21; Fad+20].)

30 We can see an example of this process in Figure 20.8, where we use a β -VAE model (Section 21.3.1)
 31 fit to the MNIST dataset. We see that the model is able to produce plausible interpolations between
 32 the digit 7 and the digit 2. As a more interesting example, we can fit a β -VAE to the the **CelebA**
 33 dataset [Liu+15].³ The results are shown in Figure 20.9, and look reasonable. (We can get much
 34 better quality if we use a larger model trained on more data for a longer amount of time.)

35 It is also possible to perform interpolation in the latent space of text models, as illustrated in
 36 Figure 21.10.

37 20.3.6 Latent space arithmetic

40 In some cases, we can go beyond interpolation, and can perform **latent space arithmetic**, in which
 41 we can increase or decrease the amount of a desired “semantic factor of variation”. This was first
 42 shown in the **word2vec** model [Mik+13], but it also is possible in other latent variable models. For
 43 example, consider our VAE model fit to CelebA dataset, which has faces of celebrities and some
 44

45 3. CelebA contains about 200k images of famous celebrities. The images are also annotated with 40 attributes. We
 46 reduce the resolution of the images to 64x64, as is conventional. See Section 21.2.6 for the details of the model.



Figure 20.10: Arithmetic in the latent space of a β -VAE (with $\beta = 0.5$). The first column is an input image, with embedding \mathbf{z} . Subsequent columns show the decoding of $\mathbf{z} + s\Delta$, where $s \in \{-2, -1, 0, 1, 2\}$ and $\Delta = \bar{\mathbf{z}}^+ - \bar{\mathbf{z}}^-$ is the difference in the average embeddings of images with or without a certain attribute (here, wearing sunglasses). Generated by `celeba_vae_ae_comparison.ipynb`.

corresponding attributes. Let \mathbf{X}_i^+ be a set of images which have attribute i , and \mathbf{X}_i^- be a set of images which do not have this attribute. Let \mathbf{Z}_i^+ and \mathbf{Z}_i^- be the corresponding embeddings, and $\bar{\mathbf{z}}_i^+$ and $\bar{\mathbf{z}}_i^-$ be the average of these embeddings. We define the offset vector as $\Delta_i = \bar{\mathbf{z}}_i^+ - \bar{\mathbf{z}}_i^-$. If we add some positive multiple of Δ_i to a new point \mathbf{z} , we increase the amount of the attribute i ; if we subtract some multiple of Δ_i , we decrease the amount of the attribute i . [Whi16].

We give an example of this in Figure 20.10. We consider the attribute of wearing sunglasses. The j 'th reconstruction is computed using $\hat{\mathbf{x}}_j = d(\mathbf{z} + s_j\Delta)$, where $\mathbf{z} = e(\mathbf{x})$ is the encoding of the original image, and s_j is a scale factor. When $s_j > 0$ we add sunglasses to the face. When $s_j < 0$ we remove sunglasses; but this also has the side effect of making the face look younger and more female, possibly a result of dataset bias.

20.3.7 Generative design

Another interesting use case for (deep) generative models is **generative design**, in which we use the model to generate candidate objects, such as molecules, which have desired properties (see e.g., [RNA22]). One approach is to fit a VAE to unlabeled samples, and then to perform Bayesian optimization (Section 6.8) in its latent space, as discussed in Section 21.3.6.2.

20.3.8 Model-based reinforcement learning

We discuss reinforcement learning (RL) in Chapter 35. The main success stories of RL to date have been in computer games, where simulators exist and data is abundant. However, in other areas, such as robotics, data is expensive to acquire. In this case, it can be useful to learn a generative “**world model**”, so the agent can do planning and learning “in its head”. See Section 35.4 for more details.

20.3.9 Representation learning

Representation learning refers to learning (possibly uninterpretable) latent factors \mathbf{z} that generate the observed data \mathbf{x} . The primary goal is for these features to be used in “**downstream**” supervised tasks. This is discussed in Chapter 32.

1 **20.3.10 Data compression**

3 Models which can assign high probability to frequently occurring data vectors (e.g., images, sentences),
4 and low probability to rare vectors, can be used for **data compression**, since we can assign shorter
5 codes to the more common items. Indeed, the optimal coding length for a vector \mathbf{x} from some
6 stochastic source $p(\mathbf{x})$ is $l(\mathbf{x}) = -\log p(\mathbf{x})$, as proved by Shannon. See Section 5.4 for details.
7

8 **20.4 Evaluating generative models**

10 This section is coauthored with Mihaela Rosca, Shakir Mohamed and Balaji Lakshminarayanan.
11

12 Evaluating generative models requires metrics which capture

- 13 • **sample quality** - are samples generated by the model a part of the data distribution?
- 14 • **sample diversity** - are samples from the model distribution capturing all modes of the data
distribution?
- 15 • **generalization** - is the model generalizing beyond the training data?

19 There is no known metric which meets all these requirements, but various metrics have been proposed
20 to capture different aspects of the learned distribution, some of which we discuss below.

22 **20.4.1 Likelihood-based evaluation**

24 A standard way to measure how close a model q is to a true distribution p is in terms of the KL
25 divergence (Section 5.1):

$$\text{D}_{\text{KL}}(p \parallel q) = \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} = -\mathbb{H}(p) + \mathbb{H}_{ce}(p, q) \quad (20.2)$$

30 where $\mathbb{H}(p)$ is a constant, and $\mathbb{H}_{ce}(p, q)$ is the cross entropy. If we approximate $p(\mathbf{x})$ by the empirical
31 distribution, we can evaluate the cross entropy in terms of the empirical **negative log likelihood**
32 on the dataset:

$$\text{NLL} = -\frac{1}{N} \sum_{n=1}^N \log q(\mathbf{x}_n) \quad (20.3)$$

37 Usually we care about negative log likelihood on a held-out test set.⁴
38

39 **20.4.1.1 Computing log-likelihood**

41 For models of discrete data, such as language models, it is easy to compute the (negative) log
42 likelihood. However, it is common to measure performance using a quantity called **perplexity**, which
43 is defined as 2^H , where $H = \text{NLL}$ is the cross entropy or negative log likelihood.
44

45 4. In some applications, we report **bits per dimension**, which is the NLL using log base 2, divided by the dimensionality
46 of \mathbf{x} . (To compute this metric, recall that $\log_2 L = \frac{\log_e L}{\log_e 2}$.)

For image and audio models, one complication is that the model is usually a continuous distribution $p(\mathbf{x}) \geq 0$ but the data is usually discrete (e.g., $\mathbf{x} \in \{0, \dots, 255\}^D$ if we use one byte per pixel). Consequently the average log likelihood can be arbitrary large, since the pdf can be bigger than 1. To avoid this it is standard practice to use **uniform dequantization** [TOB16], in which we add uniform random noise to the discrete data, and then treat it as continuous-valued data. This gives a lower bound on the average log likelihood of the discrete model on the original data.

To see this, let \mathbf{z} be a continuous latent variable, and \mathbf{x} be a vector of binary observations computed by rounding, so $p(\mathbf{x}|\mathbf{z}) = \delta(\mathbf{x} - \text{round}(\mathbf{z}))$, computed elementwise. We have $p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$. Let $q(\mathbf{z}|\mathbf{x})$ be a probabilistic inverse of \mathbf{x} , that is, it has support only on values where $p(\mathbf{x}|\mathbf{z}) = 1$. In this case, Jensen's inequality gives

$$\log p(\mathbf{x}) \geq \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z}) + \log p(\mathbf{z}) - \log q(\mathbf{z}|\mathbf{x})] \quad (20.4)$$

$$= \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{z}) - \log q(\mathbf{z}|\mathbf{x})] \quad (20.5)$$

Thus if we model the density of $\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})$, which is a dequantized version of \mathbf{x} , we will get a lower bound on $p(\mathbf{x})$.

20.4.1.2 Challenges with using the likelihood

Unfortunately, there are several challenges with using the likelihood to evaluate generative models, some of which we discuss below.

20.4.1.3 Likelihood can be hard to compute

For many models, computing the likelihood can be computationally expensive, since it requires knowing the normalization constant of the probability model. One solution is to use variational inference (Chapter 10), which provides a way to efficiently compute lower (and sometimes upper) bounds on the log likelihood. Another solution is to use annealed importance sampling (Section 11.5.4.1), which provides a way to estimate the log likelihood using Monte Carlo sampling. However, in the case of implicit generative models, such as GANs (Chapter 26), the likelihood is not even defined, so we need to find evaluation metrics that do not rely on likelihood.

20.4.1.4 Likelihood is not related to sample quality

A more subtle concern with likelihood is that it is often uncorrelated with the perceptual quality of the samples, at least for real-valued data, such as images and sound. In particular, a model can have great log-likelihood but create poor samples and vice versa.

To see why a model can have good likelihoods but create bad samples, consider the following argument from [TOB16]. Suppose q_0 is a density model for D -dimensional data \mathbf{x} which performs arbitrarily well as judged by average log-likelihood, and suppose q_1 is a bad model, such as white noise. Now consider samples generated from the mixture model

$$q_2(\mathbf{x}) = 0.01q_0(\mathbf{x}) + 0.99q_1(\mathbf{x}) \quad (20.6)$$

Clearly 99% of the samples will be poor. However, the log-likelihood per pixel will hardly change between q_2 and q_0 if D is large, since

$$\log q_2(\mathbf{x}) = \log[0.01q_0(\mathbf{x}) + 0.99q_1(\mathbf{x})] \geq \log[0.01q_0(\mathbf{x})] = \log q_0(\mathbf{x}) - 100 \quad (20.7)$$

1 For high-dimensional data, $|\log q_0(\mathbf{x})| \sim D \gg 100$, so $\log q_2(\mathbf{x}) \approx \log q_0(\mathbf{x})$, and hence mixing in the
 2 poor sampler does not significantly impact the log likelihood.
 3

4 Now consider a case where the model has good samples but bad likelihoods. To achieve this,
 5 suppose q is a GMM centered on the training images:

$$6 \quad q(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \mathcal{N}(\mathbf{x} | \mathbf{x}_n, \epsilon^2 \mathbf{I}) \quad (20.8)$$

9 If ϵ is small enough that the Gaussian noise is imperceptible, then samples from this model will look
 10 good, since they correspond to the training set of real images. But this model will almost certainly
 11 have poor likelihood on the test set due to overfitting. (In this case we say the model has effectively
 12 just memorized the training set.)
 13

14 20.4.2 Distances and divergences in feature space

15 Due to the challenges associated with comparing distributions in high dimensional spaces, and the
 16 desire to compare distributions in a semantically meaningful way, it is common to use domain-specific
 17 **perceptual distance metrics**, that measure how similar data vectors are to each other or to the
 18 training data. However, most metrics used to evaluate generative models do not directly compare
 19 raw data (e.g. pixels) but use a neural network to obtain features from the raw data and compare
 20 the feature distribution obtained from model samples with the feature distribution obtained from
 21 the dataset. The neural network used to obtain features can be trained solely for the purpose of
 22 evaluation, or can be pretrained; a common choice is to use a pretrained classifier (see e.g., [Sal+16;
 23 Heu+17b; Bin+18; Kyn+19; SSG18a]).

24 The **Inception Score** [Sal+16] measures the average KL divergence between the marginal distribution of class labels obtained from the samples $p_{\theta}(y) = \int p(y|\mathbf{x})p_{\theta}(\mathbf{x})$ and the distribution $p(y|\mathbf{x})$ obtained from a sample $\mathbf{x} \sim p_{\theta}(\mathbf{x})$. This leads to the following score:

$$25 \quad IS = \exp [\mathbf{E}_{p_{\theta}(\mathbf{x})} D_{\text{KL}}(p(y|\mathbf{x}) \| p_{\theta}(y))] \quad (20.9)$$

26 If a model produces high quality samples from all classes in the dataset, then $p_{\theta}(y)$ will often be
 27 close to uniform, while $p(y|\mathbf{x})$ should be a sharp distribution corresponding to the class associated
 28 with \mathbf{x} ; this leads to a high $D_{\text{KL}}(p(y|\mathbf{x}) \| p_{\theta}(y))$ and thus a high Inception Score score.
 29

30 The Inception Score solely relies on class labels, and thus does not measure overfitting or sample
 31 diversity outside the predefined dataset classes. For example, a model which generates one perfect
 32 example per class would get a perfect Inception Score, despite not capturing the variety of examples
 33 inside a class, as shown in Figure 20.11a. To address this drawback, the **Fréchet Inception Distance**
 34 or **FID** score [Heu+17b] measures the Fréchet distance between two Gaussian distributions on sets
 35 of features of a pre-trained classifier. One Gaussian is obtained by passing model samples through a
 36 pretrained classifier, and the other by passing samples the dataset through the same classifier. If we
 37 assume that the mean and covariance obtained from model features are μ_m and Σ_m and those from
 38 the data are μ_d and Σ_d , then the FID is
 39

$$40 \quad FID = \|\mu_m - \mu_d\|_2^2 + \text{trace}(\Sigma_d + \Sigma_m - 2(\Sigma_d \Sigma_m)^{1/2}) \quad (20.10)$$

41 Since it uses features instead of class logits, the Fréchet distance captures more than modes captured
 42 by class labels, as shown in Figure 20.11b. Unlike the Inception score, a lower score is better since
 43 we want the two distributions to be as close as possible.
 44

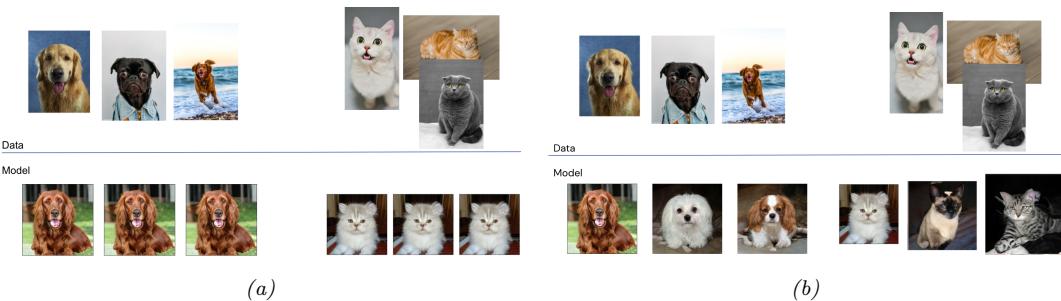


Figure 20.11: (a) Model samples with good (high) inception score are visually realistic. (b) Model samples with good (low) FID score are visually realistic and diverse.

Unfortunately, the Fréchet distance has been shown to have a high bias, with results varying widely based on the number of samples used to compute the score. To mitigate this issue, the **Kernel Inception Distance** has been introduced [Bin+18], which measures the squared MMD (Section 2.7.3) between the features obtained from the data and features obtained from model samples.

20.4.3 Precision and recall metrics

Since the FID only measures the distance between the data and model distributions, it is difficult to use it as a diagnostic tool: a bad (high) FID can indicate that the model is not able to generate high quality data, or that it puts too much mass around the data distribution, or that the model only captures a subset of the data (e.g. in Figure 26.7). Trying to disentangle between these two failure modes has been the motivation to seek individual precision (sample quality) and recall (sample diversity) metrics in the context of generative models [LPO17; Kyn+19]. (The diversity question is especially important in the context of GANs, where mode collapse (Section 26.3.3) can be an issue.)

A common approach is to use nearest neighbors in the feature space of a pretrained classifier to define precision and recall [Kyn+19]. To formalize this, let us define

$$f_k(\phi, \Phi) = \begin{cases} 1 & \text{if } \exists \phi' \in \Phi \text{ s.t. } \|\phi - \phi'\|_2^2 \leq \|\phi' - \text{NN}_k(\phi', \Phi)\|_2^2 \\ 0 & \text{otherwise} \end{cases} \quad (20.11)$$

where Φ is a set of feature vectors and $\text{NN}_k(\phi', \Phi)$ is a function returning the k -th nearest neighbor of ϕ' in Φ . We now define precision and recall as follows:

$$\text{precision}(\Phi_{model}, \Phi_{data}) = \frac{1}{|\Phi_{model}|} \sum_{\phi \in \Phi_{model}} f_k(\phi, \Phi_{data}); \quad (20.12)$$

$$\text{recall}(\Phi_{model}, \Phi_{data}) = \frac{1}{|\Phi_{data}|} \sum_{\phi \in \Phi_{data}} f_k(\phi, \Phi_{model}); \quad (20.13)$$

Precision and recall are always between 0 and 1. Intuitively, the precision metric measures whether samples are as close to data as data is to other data examples, while recall measures whether data is as close to model samples as model samples are to other samples. The parameter k controls

1 how lenient the metrics will be – the higher k , the higher both precision and recall will be. As in
2 classification, precision and recall in generative models can be used to construct a trade-off curve
3 between different models which allows practitioners to make an informed decision regarding which
4 model they want to use.
5

6

7 **20.4.4 Statistical tests**

8 Statistical tests have long been used to determine whether two sets of samples have been generated
9 from the same distribution; these types of statistical tests are called **two sample tests**. Let us
10 define the null hypothesis as the statement that both set of samples are from the same distribution.
11 We then compute a statistic from the data and compare it to a threshold, and based on this we
12 decide whether to reject the null hypothesis. In the context of evaluating implicit generative models
13 such as GANs, statistics based on classifiers [Saj+18] and the MMD [Liu+20b] have been used. For
14 use in scenarios with high dimensional input spaces, which are ubiquitous in the era of deep learning,
15 two sample tests have been adapted to use learned features instead of raw data.

16 Like all other evaluation metrics for generative models, statistical tests have their own advantages
17 and disadvantages: while users can specify Type 1 error – the chance they allow that the null
18 hypothesis is wrongly rejected – statistical tests tend to be computationally expensive and thus
19 cannot be used to monitor progress in training; hence they are best used to compare fully trained
20 models.

21

22 **20.4.5 Challenges with using pretrained classifiers**

23 While popular and convenient, evaluation metrics that rely on pretrained classifiers (such as IS, FID,
24 nearest neighbors in feature space, and statistical tests in feature space) have significant drawbacks.
25 One might not have a pretrained classifier available for the dataset at hand, so classifiers trained on
26 other datasets are used. Given the well known challenges with neural network generalization (see
27 Section 17.4), the features of a classifier trained on images from one dataset might not be reliable
28 enough to provide a fine grained signal of quality for samples obtained from a model trained on a
29 different dataset. If the generative model is trained on the same dataset as the pre-trained classifier
30 but the model is not capturing the data distribution perfectly, we are presenting the pre-trained
31 classifier with out-of-distribution data and relying on its features to obtain score to evaluate our
32 models. Far from being purely theoretical concerns, these issues have been studied extensively and
33 have been shown to affect evaluation in practice [RV19; BS18].
34

35

36 **20.4.6 Using model samples to train classifiers**

37 Instead of using pretrained classifiers to evaluate samples, one can train a classifier on samples from
38 conditional generative models, and then see how good these classifiers are at classifying data. For
39 example, does adding synthetic (sampled) data to the real data help? This is closer to a reliable
40 evaluation of generative model samples, since ultimately, the performance of generative models is
41 dependent on the downstream task they are trained for. If used for semi supervised learning, one
42 should assess how much adding samples to a classifier dataset helps with test accuracy. If used for
43 model based reinforcement learning, one should assess how much the generative model helps with
44 agent performance. For examples of this approach, see e.g., [SSM18; SSA18; RV19; SS20b].
45

46



Figure 20.12: Illustration of nearest neighbors in feature space: in the top left we have the query sample generated using BigGAN, and the rest of the images are its nearest neighbors from the dataset. The nearest neighbors search is done in the feature space of a pretrained classifier. From Figure 13 of [BDS18]. Used with kind permission of Andy Brock.

20.4.7 Assessing overfitting

Many of the metrics discussed so far capture the sample quality and diversity, but do not capture overfitting to the training data. To capture overfitting, often a visual inspection is performed: a set of samples is generated from the model and for each sample its closest K nearest neighbors in the feature space of a pretrained classifier are obtained from the dataset. While this approach requires manually assessing samples, it is a simple way to test whether a model is simply memorizing the data. We show an example in Figure 20.12: since the model sample in the top left is quite different than its neighbors from the dataset (remaining images), we can conclude the sample is not simply memorised from the dataset. Similarly, sample diversity can be measured by approximating the support of the learned distribution by looking for similar samples in a large sample pool — as in the pigeonhole principle — but it is expensive and often requires manual human assessment[AZ17].

For likelihood-based models — such as variational autoencoders Chapter 21, autoregressive models Chapter 22, and normalising flows Chapter 23 — we can assess memorisation by seeing how much the log-likelihood of a model changes when a sample is included in the model’s training set or not [BW21].

20.4.8 Human evaluation

One approach to evaluate generative models is to use human evaluation, by presenting samples from the model alongside samples from the data distribution, and ask human raters to compare the quality

1 of the samples [Zho+19b]. Human evaluation is a suitable metric if the model is used to create art or
2 other data for human display, or if reliable automated metrics are hard to obtain. However, human
3 evaluation can be difficult to standardize, hard to automate and can be expensive or cumbersome to
4 set up.
5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

21 Variational autoencoders

21.1 Introduction

In this chapter, we discuss generative models of the form

$$\mathbf{z} \sim p_{\theta}(\mathbf{z}) \tag{21.1}$$

$$\mathbf{x}|\mathbf{z} \sim \text{Expfam}(\mathbf{x}|d_{\theta}(\mathbf{z})) \tag{21.2}$$

where $p(\mathbf{z})$ is some kind of prior on the latent code \mathbf{z} , $d_{\theta}(\mathbf{z})$ is a deep neural network, known as the **decoder**, and $\text{Expfam}(\mathbf{x}|\boldsymbol{\eta})$ is an exponential family distribution, such as a Gaussian or product of Bernoullis. This is called a **deep latent variable model** or **DLVM**. When the prior is Gaussian (as is often the case), this model is called a **deep latent Gaussian model** or **DLGM**.

Posterior inference (i.e., computing $p_{\theta}(\mathbf{z}|\mathbf{x})$) is computationally intractable, as is computing the marginal likelihood

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p_{\theta}(\mathbf{z}) d\mathbf{z} \tag{21.3}$$

Hence we need to resort to approximate inference. For most of this chapter, we will use **amortized inference**, which we discussed in Section 10.3.6. This trains another model, $q_{\phi}(\mathbf{z}|\mathbf{x})$, called the **recognition network** or **inference network**, simultaneously with the generative model to do approximate posterior inference. This combination is called a **variational autoencoder** or **VAE** [KW14; RMW14b; KW19a], since it can be thought of as a probabilistic version of a deterministic autoencoder, discussed in Section 16.3.3.

In this chapter, we introduce the basic VAE, as well as some extensions. Note that the literature on VAE-like methods is vast¹, so we will only discuss a small subset of the ideas that have been explored.

21.2 VAE basics

In this section, we discuss the basics of variational autoencoders.

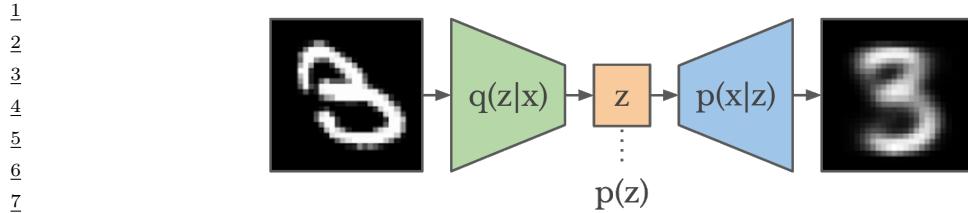


Figure 21.1: Schematic illustration of a VAE. From a figure in [Haf18]. Used with kind permission of Danijar Hafner.

21.2.1 Modeling assumptions

In the simplest setting, a VAE defines a generative model of the form

$$p_{\theta}(z, x) = p_{\theta}(z)p_{\theta}(x|z) \quad (21.4)$$

where $p_{\theta}(z)$ is usually a Gaussian, and $p_{\theta}(x|z)$ is usually a product of exponential family distributions (e.g., Gaussians or Bernoullis), with parameters computed by a neural network decoder, $d_{\theta}(z)$. For example, for binary observations, we can use

$$p_{\theta}(x|z) = \prod_{d=1}^D \text{Ber}(x_d | \sigma(d_{\theta}(z))) \quad (21.5)$$

In addition, a VAE fits a recognition model

$$q_{\phi}(z|x) = q(z|e_{\phi}(x)) \approx p_{\theta}(z|x) \quad (21.6)$$

to perform approximate posterior inference. Here $q_{\phi}(z|x)$ is usually a Gaussian, with parameters computed by a neural network encoder $e_{\phi}(x)$:

$$q_{\phi}(z|x) = \mathcal{N}(z|\mu, \text{diag}(\exp(\ell))) \quad (21.7)$$

$$(\mu, \ell) = e_{\phi}(x) \quad (21.8)$$

where $\ell = \log \sigma$. The model can be thought of as encoding the input x into a stochastic latent bottleneck z and then decoding it to approximately reconstruct the input, as shown in Figure 21.1.

The idea of training an inference network to “invert” a generative network, rather than running an optimization algorithm to infer the latent code, is called amortized inference, and is discussed in Section 10.3.6. This idea was first proposed in the **Helmholtz machine** [Day+95]. However, that paper did not present a single unified objective function for inference and generation, but instead used the wake-sleep (Section 10.6) method for training. By contrast, the VAE optimizes a variational lower bound on the log-likelihood, which means that convergence to a locally optimal MLE of the parameters is guaranteed.

We can use other approaches to fitting the DLGM (see e.g., [Hof17; DF19]). However, learning an inference network to fit the DLGM is often faster and can have some regularization benefits (see e.g., [KP20]).²

⁴⁵ 1. For example, the website <https://github.com/matthewvowels1/Awesome-VAEs> lists over 900 papers.

⁴⁶ 2. Combining a generative model with an inference model in this way results in what has been called a “monference”

21.2.2 Evidence lower bound (ELBO)

When fitting the model, our goal is to maximize the marginal likelihood

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|z)p_{\theta}(z) dz \quad (21.9)$$

Unfortunately, computing this quantity is intractable. However, we can use an inference network to compute an approximate posterior, $q_{\phi}(z|\mathbf{x})$, and hence a lower bound to the marginal likelihood. This idea is discussed in Section 10.1.2, but we repeat the details here, using slightly different notation.

First note that we have the following decomposition:

$$\log p_{\theta}(\mathbf{x}) = \mathbb{E}_{q_{\phi}(z|\mathbf{x})} [\log p_{\theta}(\mathbf{x})] \quad (21.10)$$

$$= \mathbb{E}_{q_{\phi}(z|\mathbf{x})} \left[\log \left(\frac{p_{\theta}(\mathbf{x}, z)}{p_{\theta}(z|\mathbf{x})} \right) \right] \quad (21.11)$$

$$= \mathbb{E}_{q_{\phi}(z|\mathbf{x})} \left[\log \left(\frac{p_{\theta}(\mathbf{x}, z)}{q_{\phi}(z|\mathbf{x})} \frac{q_{\phi}(z|\mathbf{x})}{p_{\theta}(z|\mathbf{x})} \right) \right] \quad (21.12)$$

$$= \underbrace{\mathbb{E}_{q_{\phi}(z|\mathbf{x})} \left[\log \left(\frac{p_{\theta}(\mathbf{x}, z)}{q_{\phi}(z|\mathbf{x})} \right) \right]}_{\mathcal{L}_{\theta, \phi}(\mathbf{x})} + \underbrace{\mathbb{E}_{q_{\phi}(z|\mathbf{x})} \left[\log \left(\frac{q_{\phi}(z|\mathbf{x})}{p_{\theta}(z|\mathbf{x})} \right) \right]}_{D_{\text{KL}}(q_{\phi}(z|\mathbf{x}) \| p_{\theta}(z|\mathbf{x}))} \quad (21.13)$$

The second term in Equation (21.13) is non-negative, and hence

$$\mathcal{L}_{\theta, \phi}(\mathbf{x}) \leq \log p_{\theta}(\mathbf{x}) \quad (21.14)$$

The quantity $\log p_{\theta}(\mathbf{x})$ is the log marginal likelihood, also called the **evidence**. Hence $\mathcal{L}_{\theta, \phi}(\mathbf{x})$ is called the **evidence lower bound** or **ELBO**.

We can rewrite the ELBO in 3 equivalent ways as follows:

$$\mathcal{L}_{\theta, \phi}(\mathbf{x}) = \mathbb{E}_{q_{\phi}(z|\mathbf{x})} [\log p_{\theta}(\mathbf{x}, z) - \log q_{\phi}(z|\mathbf{x})] \quad (21.15)$$

$$= \mathbb{E}_{q_{\phi}(z|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|z) + \log p_{\theta}(z)] + \mathbb{H}(q_{\phi}(z|\mathbf{x})) \quad (21.16)$$

$$= \mathbb{E}_{q_{\phi}(z|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|z)] - D_{\text{KL}}(q_{\phi}(z|\mathbf{x}) \| p_{\theta}(z)) \quad (21.17)$$

We can interpret this last objective as the expected log likelihood plus a regularization term, that ensures the (per-sample) posterior is “well behaved” (does not deviate too far from the prior in terms of KL divergence).

The tightness of this lower bound is controlled by the **variational gap**, which is given by $D_{\text{KL}}(q_{\phi}(z|\mathbf{x}) \| p_{\theta}(z|\mathbf{x}))$. A better approximate posterior results in a tighter bound. When the KL goes to zero, the posterior is exact, so any improvements to the ELBO directly translate to improvements in the likelihood of the data, as in the EM algorithm (see Section 6.6.3).

An alternative to maximizing the ELBO is to minimize the **negative ELBO**, also called the **variational free energy**, given by

$$\mathcal{L}(\theta, \phi; \mathbf{x}) = \mathbb{E}_{q_{\phi}(z|\mathbf{x})} [-\log p_{\theta}(\mathbf{x}|z)] + D_{\text{KL}}(q_{\phi}(z|\mathbf{x}) \| p_{\theta}(z)) \quad (21.18)$$

i.e., model-inference hybrid. See the blog by Jacob Andreas, <http://blog.jacobandreas.net/conference.html>, for further discussion.

1 **21.2.3 Evaluating the ELBO**

3 We can approximate the ELBO by sampling from the posterior, $\mathbf{z}_s \sim q_{\phi}(\mathbf{z}|\mathbf{x})$, and then evaluating
4 the expectation term using a Monte Carlo estimate,
5

6
$$\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z})] \approx \frac{1}{S} \sum_{s=1}^S \log p_{\theta}(\mathbf{x}, \mathbf{z}_s) \quad (21.19)$$

7

9 We can often compute the (differential) entropy term $\mathbb{H}(q_{\phi}(\mathbf{z}|\mathbf{x}))$ analytically, depending on the
10 choice of variational distribution. For example, if we use a Gaussian posterior, $q(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$,
11 we can use Equation (5.95) to compute the entropy:
12

13
$$\mathbb{H}(q_{\phi}(\mathbf{z}|\mathbf{x})) = \frac{1}{2} \ln |\boldsymbol{\Sigma}| + \text{const} \quad (21.20)$$

14

16 Similarly we can often compute the KL term $D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p_{\theta}(\mathbf{z}))$ analytically. For example, if
17 we assume a diagonal Gaussian prior, $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$, and diagonal gaussian posterior, $q(\mathbf{z}|\mathbf{x}) =$
18 $\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}))$, we can use Equation (5.81) to compute the KL in closed form:
19

20
$$D_{\text{KL}}(q \parallel p) = -\frac{1}{2} \sum_{k=1}^K [\log \sigma_k^2 - \sigma_k^2 - \mu_k^2 + 1] \quad (21.21)$$

21

23 where K is the number of latent dimensions.
24

24 In cases where we cannot evaluate the entropy or KL in closed form, we can use just use Monte
25 Carlo to approximate all terms:
26

27
$$\mathbb{L}_{\theta, \phi}(\mathbf{x}) \approx \frac{1}{S} \sum_{s=1}^S [\log p_{\theta}(\mathbf{x}|\mathbf{z}_s) + \log p_{\theta}(\mathbf{z}_s) - \log q_{\phi}(\mathbf{z}_s|\mathbf{x})] \quad (21.22)$$

28

30 We often use a single MC sample, so $S = 1$.
31

32 **21.2.4 Optimizing the ELBO**

34 The ELBO for a single datapoint \mathbf{x} is given in Equation (21.17). The ELBO for the whole dataset,
35 scaled by $N = |\mathcal{D}|$, the number of examples, is given by
36

37
$$\mathbb{L}_{\theta, \phi}(\mathcal{D}) = \frac{1}{N} \sum_{\mathbf{x}_n \in \mathcal{D}} \mathbb{L}_{\theta, \phi}(\mathbf{x}_n) = \frac{1}{N} \sum_{\mathbf{x}_n \in \mathcal{D}} [\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}_n)} [\log p_{\theta}(\mathbf{x}_n|\mathbf{z}) + \log p_{\theta}(\mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x}_n)]] \quad (21.23)$$

38

40 Our goal is to maximize this wrt θ and ϕ : the former fits the model to the data, the latter reduces
41 the KL gap between the approximate and true posterior.
42

44 We can create an unbiased minibatch approximation of this objective by sampling examples \mathbf{x} ,
45 and then computing the objective for a given \mathbf{x} . So now we focus on a fixed \mathbf{x} , for brevity, and drop
46 the sum over \mathbf{x}_n .
47

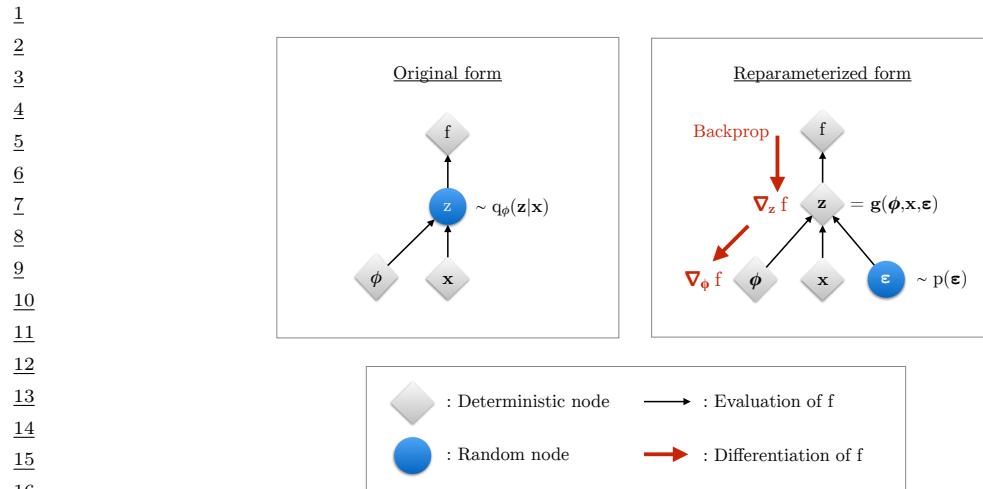


Figure 21.2: Illustration of the reparameterization trick. The objective f depends on the variational parameters ϕ , the observed data x , and the latent random variable $z \sim q_\phi(z|x)$. On the left, we show the standard form of the computation graph. On the right, we show a reparameterized form, in which we move the stochasticity into the noise source ϵ , and compute z deterministically, $z = g(\phi, x, \epsilon)$. The rest of the graph is deterministic, so we can backpropagate the gradient of the scalar f wrt ϕ through z and into ϕ . From Figure 2.3 of [KW19a]. Used with kind permission of Durk Kingma.

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

$$\nabla_\theta \mathbb{L}_{\theta, \phi}(x) = \nabla_\theta \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x, z) - \log q_\phi(z|x)] \quad (21.24)$$

$$= \mathbb{E}_{q_\phi(z|x)} [\nabla_\theta \{\log p_\theta(x, z) - \log q_\phi(z|x)\}] \quad (21.25)$$

$$\approx \nabla_\theta \log p_\theta(x, z^s) \quad (21.26)$$

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

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

$$\nabla_\phi \mathbb{L}_{\theta, \phi}(x) = \nabla_\phi \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x, z) - \log q_\phi(z|x)] \quad (21.27)$$

$$\neq \mathbb{E}_{q_\phi(z|x)} [\nabla_\phi \{\log p_\theta(x, z) - \log q_\phi(z|x)\}] \quad (21.28)$$

However, we can often use the reparameterization trick, which we discuss in Section 21.2.5. (If not, we can use black box VI, which we discuss in Section 10.3.2.)

21.2.5 Using the reparameterization trick to compute ELBO gradients

In this section, we discuss the **reparameterization trick** for taking gradients wrt distributions over continuous latent variables $z \sim q_\phi(z|x)$. We explain this in detail in Section 6.5.4, but we summarize the basic idea here.

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

$$\mathbf{z} = r(\boldsymbol{\epsilon}, \phi, \mathbf{x}) \quad (21.29)$$

For example,

$$\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma})) \iff \mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\epsilon} \odot \boldsymbol{\sigma}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (21.30)$$

Using this, we have

$$\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[f(\mathbf{z})] = \mathbb{E}_{p(\boldsymbol{\epsilon})}[f(\mathbf{z})] \quad \text{s.t. } \mathbf{z} = r(\boldsymbol{\epsilon}, \phi, \mathbf{x}) \quad (21.31)$$

where we define

$$f(\mathbf{z}) = \log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x}) \quad (21.32)$$

Hence

$$\nabla_{\phi} \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[f(\mathbf{z})] = \nabla_{\phi} \mathbb{E}_{p(\boldsymbol{\epsilon})}[f(\mathbf{z})] = \mathbb{E}_{p(\boldsymbol{\epsilon})}[\nabla_{\phi} f(\mathbf{z})] \quad (21.33)$$

which we can approximate with a single Monte Carlo sample. This lets us propagate gradients back through the f function and then into the DNN transformation function r that is used to compute $\mathbf{z} = r(\boldsymbol{\epsilon}, \phi, \mathbf{x})$. See Figure 21.2 for an illustration.

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

$$\log q_\phi(\mathbf{z}|\mathbf{x}) = \log p(\boldsymbol{\epsilon}) - \log \left| \det \left(\frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}} \right) \right| \quad (21.34)$$

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

$$\frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}} = \begin{pmatrix} \frac{\partial z_1}{\partial \epsilon_1} & \dots & \frac{\partial z_1}{\partial \epsilon_k} \\ \vdots & \ddots & \vdots \\ \frac{\partial z_k}{\partial \epsilon_1} & \dots & \frac{\partial z_k}{\partial \epsilon_k} \end{pmatrix} \quad (21.35)$$

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

21.2.5.1 Fully factorized Gaussian

Suppose we have a fully factorized Gaussian posterior:

$$\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (21.36)$$

$$\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon} \quad (21.37)$$

$$(\boldsymbol{\mu}, \log \boldsymbol{\sigma}) = e_\phi(\mathbf{x}) \quad (21.38)$$

Then the Jacobian is $\frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}} = \text{diag}(\boldsymbol{\sigma})$, so

$$\log q_\phi(\mathbf{z}|\mathbf{x}) = \sum_{k=1}^K \log \mathcal{N}(\epsilon_k|0, 1) - \log \sigma_k = \sum_{k=1}^K -\frac{1}{2} \log(2\pi) - \frac{1}{2} \epsilon_k^2 - \log \sigma_k \quad (21.39)$$

21.2.5.2 Full covariance Gaussian

Now consider a full covariance Gaussian posterior:

(21.40)

$$z = \mu + L\epsilon \quad (21.41)$$

where \mathbf{L} is a lower triangular matrix with non-zero entries on the diagonal, which satisfies $\boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}^T$. The Jacobian of this affine transformation is $\frac{\partial \mathbf{z}}{\partial \epsilon} = \mathbf{L}$. Since \mathbf{L} is a triangular matrix, its determinant is the product of its main diagonal, so

$$\log \left| \det \frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}} \right| = \sum_{k=1}^K \log |L_{kk}| \quad (21.42)$$

We need to make the parameters of the transformation r be a function of the inputs \mathbf{x} . One way to do this is to define

$$(\mu, \log \sigma, L') = e_\phi(x) \quad (21.43)$$

$$\mathbf{L} = \mathbf{M} \odot \mathbf{L}' + \text{diag}(\boldsymbol{\sigma}) \quad (21.44)$$

where \mathbf{M} is a masking matrix with 0s on and above the diagonal, and 1s below the diagonal. With this construction, the diagonal entries of \mathbf{L} are given by σ , so

$$\log \left| \det \frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}} \right| = \sum^K \log |L_{kk}| = \sum^K \log \sigma_k \quad (21.45)$$

See Algorithm 36 for the corresponding pseudo code for computing the reparameterized ELBO.

Algorithm 36: Computing a single sample unbiased estimate of the reparameterized ELBO for a VAE with full covariance Gaussian posterior, and factorized Bernoulli likelihood. Based on Algorithm 2 of [KW19a].

- ```

1 $(\mu, \log \sigma, L') = e_\phi(x)$
2 $M = np.triu(np.ones(K), -1)$
3 $L = M \odot L' + \text{diag}(\sigma)$
4 $\epsilon \sim \mathcal{N}(0, I)$
5 $z = L\epsilon + \mu$
6 $p = d_\theta(z)$
7 $\mathcal{L}_{\text{logqz}} = -\sum_{k=1}^K \left[\frac{1}{2}\epsilon_k^2 + \frac{1}{2}\log(2\pi) + \log \sigma_k \right] // \text{ from } q_\phi(z|x)$
8 $\mathcal{L}_{\text{logpz}} = -\sum_{k=1}^K \left[\frac{1}{2}z_k^2 + \frac{1}{2}\log(2\pi) \right] // \text{ from } p_\theta(z)$
9 $\mathcal{L}_{\text{logpx}} = -\sum_{d=1}^D [x_d \log p_d + (1 - x_d) \log(1 - p_d)] // \text{ from } p_\theta(x|z)$
10 $\mathcal{L} = \mathcal{L}_{\text{logqz}} + \mathcal{L}_{\text{logpz}} + \mathcal{L}_{\text{logpx}}$

```

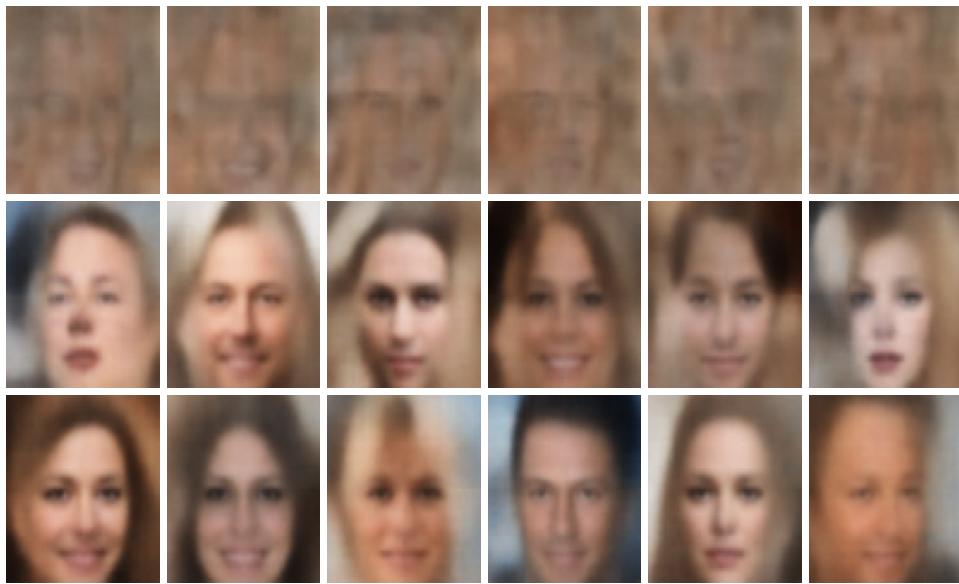


Figure 21.3: Illustration of unconditional image generation using (V)AEs trained on CelebA. Row 1: Deterministic autoencoder. Row 2:  $\beta$ -VAE with  $\beta = 0.5$ . Row 3: VAE (with  $\beta = 1$ ). Generated by [celeba\\_vae\\_ae\\_comparison.ipynb](#).

### 21.2.5.3 Inverse autoregressive flows

In Section 10.4.3, we discuss how to use inverse autoregressive flows to learn more expressive posteriors  $q_\phi(\mathbf{z}|\mathbf{x})$ , leveraging the tractability of the Jacobian of this nonlinear transformation.

### 21.2.6 Comparison of VAEs and autoencoders

VAEs are very similar to deterministic autoencoders (AE). There are 2 main differences: in the AE, the objective is the log likelihood of the reconstruction without any KL term; and in addition, the encoding is deterministic, so the encoder network just needs to compute  $\mathbb{E}[\mathbf{z}|\mathbf{x}]$  and not  $\mathbb{V}[\mathbf{z}|\mathbf{x}]$ . In view of these similarities, one can use the same codebase to implement both methods. However, it is natural to wonder what the benefits and potential drawbacks of the VAE are compared to the deterministic AE.

We shall answer this question by fitting both models to the CelebA dataset. Both models have the same convolutional structure with the following number of hidden channels per convolutional layer in the encoder: (32, 64, 128, 256, 512). The spatial size of each layer is as follows: (32, 16, 8, 4, 2). The final  $2 \times 2 \times 512$  convolutional layer then gets reshaped and passed through a linear layer to generate the mean and (marginal) variance of the stochastic latent vector, which has size 256. The structure of the decoder is the mirror image of the encoder. Each model is trained for 5 epochs with a batch size of 256, which takes about 20 minutes on a GPU.

The main advantage of a VAE over a deterministic autoencoder is that it defines a proper generative model, that can create sensible-looking novel images by decoding prior samples  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . By



*Figure 21.4: Illustration of image reconstruction using (V)AEs trained and applied to CelebA. Row 1: Original images. Row 2: Deterministic autoencoder. Row 3:  $\beta$ -VAE with  $\beta = 0.5$ . Row 4: VAE (with  $\beta = 1$ ). Generated by [celeba\\_vae\\_ae\\_comparison.ipynb](#).*

contrast, an autoencoder only knows how to decode latent codes derived from the training set, so does poorly when fed random inputs. This is illustrated in Figure 21.3.

We can also use both models to reconstruct a given input image. In Figure 21.4, we see that both AE and VAE can reconstruct the input images reasonably well, although the VAE reconstructions are somewhat blurry, for reasons we discuss in Section 21.3.1. We can reduce the amount of blurriness by scaling down the KL penalty term by a factor of  $\beta$ ; this is known as the  $\beta$ -VAE, and is discussed in more detail in Section 21.3.1.

#### 21.2.7 VAEs optimize in an augmented space

In this section, we derive several alternative expressions for the ELBO which shed light on how VAEs work.

First, let us define the joint generative distribution

$$p_{\theta}(x, z) = p_{\theta}(z)p_{\theta}(x|z) \quad (21.46)$$

1 from which we can derive the generative data marginal  
2

3  
4  $p_{\theta}(\mathbf{x}) = \int_{\mathbf{z}} p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z}$  (21.47)  
5

6 and the generative posterior  
7

8  $p_{\theta}(\mathbf{z}|\mathbf{x}) = p_{\theta}(\mathbf{x}, \mathbf{z})/p_{\theta}(\mathbf{x})$  (21.48)  
9

10 Let us also define the joint *inference* distribution

11  
12  $q_{\mathcal{D}, \phi}(\mathbf{z}, \mathbf{x}) = p_{\mathcal{D}}(\mathbf{x}) q_{\phi}(\mathbf{z}|\mathbf{x})$  (21.49)

13 where  
14

15  
16  $p_{\mathcal{D}}(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \delta(\mathbf{x}_n - \mathbf{x})$  (21.50)  
17

18 is the empirical distribution. From this we can derive the inference latent marginal, also called the  
19 **aggregated posterior**:

21  
22  $q_{\mathcal{D}, \phi}(\mathbf{z}) = \int_{\mathbf{x}} q_{\mathcal{D}, \phi}(\mathbf{x}, \mathbf{z}) d\mathbf{x}$  (21.51)  
23

24 and the inference likelihood

25  
26  $q_{\mathcal{D}, \phi}(\mathbf{x}|\mathbf{z}) = q_{\mathcal{D}, \phi}(\mathbf{x}, \mathbf{z})/q_{\mathcal{D}, \phi}(\mathbf{z})$  (21.52)  
27

28 See Figure 21.5 for a visual illustration.

29 Having defined our terms, we can now derive various alternative versions of the ELBO, following  
30 [ZSE19]. First note that the ELBO averaged over all the data is given by

31  
32  $\mathbb{L}_{\theta, \phi} = \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]] - \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p_{\theta}(\mathbf{z}))]$  (21.53)

33  $= \mathbb{E}_{q_{\mathcal{D}, \phi}(\mathbf{x}, \mathbf{z})} [\log p_{\theta}(\mathbf{x}|\mathbf{z}) + \log p_{\theta}(\mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})]$  (21.54)

34  $= \mathbb{E}_{q_{\mathcal{D}, \phi}(\mathbf{x}, \mathbf{z})} \left[ \log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\mathcal{D}, \phi}(\mathbf{x}, \mathbf{z})} + \log p_{\mathcal{D}}(\mathbf{x}) \right]$  (21.55)

35  $= -D_{\text{KL}}(q_{\mathcal{D}, \phi}(\mathbf{x}, \mathbf{z}) \parallel p_{\theta}(\mathbf{x}, \mathbf{z})) + \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [\log p_{\mathcal{D}}(\mathbf{x})]$  (21.56)

36 If we define  $\stackrel{c}{=}$  to mean equal up to additive constants, we can rewrite the above as  
37

38  
39  $\mathbb{L}_{\theta, \phi} \stackrel{c}{=} -D_{\text{KL}}(q_{\phi}(\mathbf{x}, \mathbf{z}) \parallel p_{\theta}(\mathbf{x}, \mathbf{z}))$  (21.57)

40  
41  $\stackrel{c}{=} -D_{\text{KL}}(p_{\mathcal{D}}(\mathbf{x}) \parallel p_{\theta}(\mathbf{x})) - \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p_{\theta}(\mathbf{z}|\mathbf{x}))]$  (21.58)

42 Thus maximizing the ELBO requires minimizing the two KL terms. The first KL term is minimized  
43 by MLE, and the second KL term is minimized by fitting the true posterior. Thus if the posterior  
44 family is limited, there may be a conflict between these objectives.

45

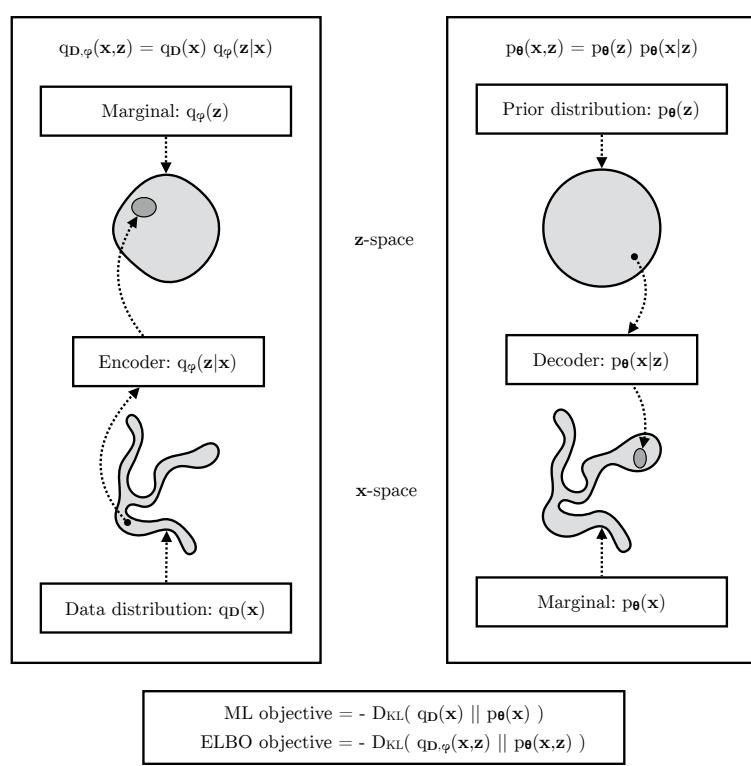


Figure 21.5: The maximum likelihood (ML) objective can be viewed as the minimization of  $D_{\text{KL}}(p_D(x) \parallel p_\theta(x))$ . (Note: in the figure,  $p_D(x)$  is denoted by  $q_D(x)$ .) The ELBO objective is minimization of  $D_{\text{KL}}(q_{D,\phi}(x,z) \parallel p_\theta(x,z))$ , which upper bounds  $D_{\text{KL}}(q_D(x) \parallel p_\theta(x))$ . From Figure 2.4 of [KW19a]. Used with kind permission of Durk Kingma.

Finally, we note that the ELBO can also be written as

$$L_{\theta,\phi} \stackrel{c}{=} -D_{\text{KL}}(q_{D,\phi}(z) \parallel p_\theta(z)) - \mathbb{E}_{q_{D,\phi}(z)} [D_{\text{KL}}(q_\phi(x|z) \parallel p_\theta(x|z))] \quad (21.59)$$

We see from Equation (21.59) that VAEs are trying to minimize the difference between the inference marginal and generative prior,  $D_{\text{KL}}(q_\phi(z) \parallel p_\theta(z))$ , while simultaneously minimizing reconstruction error,  $D_{\text{KL}}(q_\phi(x|z) \parallel p_\theta(x|z))$ . Since  $x$  is typically of much higher dimensionality than  $z$ , the latter term usually dominates. Consequently, if there is a conflict between these two objectives (e.g., due to limited modeling power), the VAE will favor reconstruction accuracy over posterior inference. Thus the learned posterior may not be a very good approximation to the true posterior (see [ZSE19] for further discussion).

### 21.3 VAE generalizations

In this section, we discuss some variants of the basic VAE model.

1 **21.3.1  $\beta$ -VAE**

3 It is often the case that VAEs generate somewhat blurry images, as illustrated in Figure 21.4,  
4 Figure 21.3 and Figure 20.9. This is not the case for models that optimize the exact likelihood, such  
5 as pixelCNNs (Section 22.3.2) and flow models (Chapter 23). To see why VAEs are different, consider  
6 the common case where the decoder is a Gaussian with fixed variance, so

7

$$\log p_{\theta}(\mathbf{x}|\mathbf{z}) = -\frac{1}{2\sigma^2} \|\mathbf{x} - d_{\theta}(\mathbf{z})\|_2^2 + \text{const} \quad (21.60)$$

8

9 Let  $e_{\phi}(\mathbf{x}) = \mathbb{E}[q_{\phi}(\mathbf{z}|\mathbf{x})]$  be the encoding of  $\mathbf{x}$ , and  $\mathcal{X}(\mathbf{z}) = \{\mathbf{x} : e_{\phi}(\mathbf{x}) = \mathbf{z}\}$  be the set of inputs that  
10 get mapped to  $\mathbf{z}$ . For a fixed inference network, the optimal setting of the generator parameters,  
11 when using squared reconstruction loss, is to ensure  $d_{\theta}(\mathbf{z}) = \mathbb{E}[\mathbf{x} : \mathbf{x} \in \mathcal{X}(\mathbf{z})]$ . Thus the decoder  
12 should predict the average of all inputs  $\mathbf{x}$  that map to that  $\mathbf{z}$ , resulting in blurry images.

13 We can solve this problem by increasing the expressive power of the posterior approximation  
14 (avoiding the merging of distinct inputs into the same latent code), or of the generator (by adding  
15 back information that is missing from the latent code), or both. However, an even simpler solution is  
16 to reduce the penalty on the KL term, making the model closer to a deterministic autoencoder:

17

$$\mathcal{L}_{\beta}(\theta, \phi | \mathbf{x}) = \underbrace{-\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\mathcal{L}_E} + \beta \underbrace{D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \| p_{\theta}(\mathbf{z}))}_{\mathcal{L}_R} \quad (21.61)$$

18

19 where  $\mathcal{L}_E$  is the reconstruction error (negative log likelihood), and  $\mathcal{L}_R$  is the KL regularizer. This is  
20 called the  $\beta$ -VAE objective [Hig+17a]. If we set  $\beta = 1$ , we recover the objective used in standard  
21 VAEs; if we set  $\beta = 0$ , we recover the objective used in standard autoencoders.

22 By varying  $\beta$  from 0 to infinity, we can reach different points on the **rate distortion curve**, as  
23 discussed in Section 5.4.2. These points make different tradeoffs between reconstruction error (distortion)  
24 and how much information is stored in the latents about the input (rate of the corresponding  
25 code). By using  $\beta < 1$ , we store more bits about each input, and hence can reconstruct images in a  
26 less blurry way. If we use  $\beta > 1$ , we get a more compressed representation.

27 **21.3.1.1 Disentangled representations**

28 One advantage of using  $\beta > 1$  is that it encourages the learning of a latent representation that is  
29 “disentangled”. Intuitively this means that each latent dimension represents a different **factor of**  
30 **variation** in the input. This is often formalized in terms of the total correlation (Section 5.3.5.1),  
31 which is defined as follows:

32

$$\text{TC}(\mathbf{z}) = \sum_k \mathbb{H}(z_k) - \mathbb{H}(\mathbf{z}) = D_{\text{KL}} \left( p(\mathbf{z}) \| \prod_k p_k(z_k) \right) \quad (21.62)$$

33

34 This is zero iff the components of  $\mathbf{z}$  are all mutually independent, and hence disentangled. In [AS18],  
35 they prove that using  $\beta > 1$  will decrease the TC.

36 Unfortunately, in [Loc+18] they prove that nonlinear latent variable models are unidentifiable, and  
37 therefore for any disentangled representation, there is an equivalent fully entangled representation  
38 with exactly the same likelihood. Thus it is not possible to recover the correct latent representation  
39 without choosing the appropriate inductive bias, via the encoder, decoder, prior, dataset, or learning  
40 algorithm, i.e., merely adjusting  $\beta$  is not sufficient. See Section 32.4.1 for more discussion.

41

---

### 21.3.1.2 Connection with information bottleneck

In this section, we show that the  $\beta$ -VAE is an unsupervised version of the information bottleneck (IB) objective from Section 5.6. If the input is  $\mathbf{x}$ , the hidden bottleneck is  $\mathbf{z}$ , and the target outputs are  $\tilde{\mathbf{x}}$ , then the unsupervised IB objective becomes

$$\mathcal{L}_{\text{UIB}} = \beta \mathbb{I}(\mathbf{z}; \mathbf{x}) - \mathbb{I}(\mathbf{z}; \tilde{\mathbf{x}}) \quad (21.63)$$

$$= \beta \mathbb{E}_{p(\mathbf{x}, \mathbf{z})} \left[ \log \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{x})p(\mathbf{z})} \right] - \mathbb{E}_{p(\mathbf{z}, \tilde{\mathbf{x}})} \left[ \log \frac{p(\mathbf{z}, \tilde{\mathbf{x}})}{p(\mathbf{z})p(\tilde{\mathbf{x}})} \right] \quad (21.64)$$

where

$$p(\mathbf{x}, \mathbf{z}) = p_{\mathcal{D}}(\mathbf{x})p(\mathbf{z}|\mathbf{x}) \quad (21.65)$$

$$p(\mathbf{z}, \tilde{\mathbf{x}}) = \int p_{\mathcal{D}}(\mathbf{x})p(\mathbf{z}|\mathbf{x})p(\tilde{\mathbf{x}}|\mathbf{z})d\mathbf{x} \quad (21.66)$$

Intuitively, the objective in Equation (21.63) means we should pick a representation  $\mathbf{z}$  that can predict  $\tilde{\mathbf{x}}$  reliably, while not memorizing too much information about the input  $\mathbf{x}$ . The tradeoff parameter is controlled by  $\beta$ .

From Equation (5.180), we have the following variational upper bound on this unsupervised objective:

$$\mathcal{L}_{\text{UVIB}} = -\mathbb{E}_{q_{\mathcal{D}, \phi}(\mathbf{z}, \mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] + \beta \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p_{\theta}(\mathbf{z}))] \quad (21.67)$$

which matches Equation (21.61) when averaged over  $\mathbf{x}$ .

### 21.3.2 InfoVAE

In Section 21.2.7, we discussed some drawbacks of the standard ELBO objective for training VAEs, namely the tendency to ignore the latent code when the decoder is powerful (Section 21.4), and the tendency to learn a poor posterior approximation due to the mismatch between the KL terms in data space and latent space (Section 21.2.7). We can fix these problems to some degree by using a generalized objective of the following form:

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}|\mathbf{x}) = -\lambda D_{\text{KL}}(q_{\phi}(\mathbf{z}) \parallel p_{\theta}(\mathbf{z})) - \mathbb{E}_{q_{\phi}(\mathbf{z})} [D_{\text{KL}}(q_{\phi}(\mathbf{x}|\mathbf{z}) \parallel p_{\theta}(\mathbf{x}|\mathbf{z}))] + \alpha \mathbb{I}_q(\mathbf{x}; \mathbf{z}) \quad (21.68)$$

where  $\alpha \geq 0$  controls how much we weight the mutual information  $\mathbb{I}_q(\mathbf{x}; \mathbf{z})$  between  $\mathbf{x}$  and  $\mathbf{z}$ , and  $\lambda \geq 0$  controls the tradeoff between  $\mathbf{z}$ -space KL and  $\mathbf{x}$ -space KL. This is called the **InfoVAE** objective [ZSE19]. If we set  $\alpha = 0$  and  $\lambda = 1$ , we recover the standard ELBO, as shown in Equation (21.59).

Unfortunately, the objective in Equation (21.68) cannot be computed as written, because of the intractable MI term:

$$\mathbb{I}_q(\mathbf{x}; \mathbf{z}) = \mathbb{E}_{q_{\phi}(\mathbf{x}, \mathbf{z})} \left[ \log \frac{q_{\phi}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{x})q_{\phi}(\mathbf{z})} \right] = -\mathbb{E}_{q_{\phi}(\mathbf{x}, \mathbf{z})} \left[ \log \frac{q_{\phi}(\mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \quad (21.69)$$

1 However, using the fact that  $q_\phi(\mathbf{x}|\mathbf{z}) = p_{\mathcal{D}}(\mathbf{x})q_\phi(\mathbf{z}|\mathbf{x})/q_\phi(\mathbf{z})$ , we can rewrite the objective as follows:  
 2

$$\frac{1}{4} \quad \mathbb{L} = \mathbb{E}_{q_\phi(\mathbf{x}, \mathbf{z})} \left[ -\lambda \log \frac{q_\phi(\mathbf{z})}{p_\theta(\mathbf{z})} - \log \frac{q_\phi(\mathbf{x}|\mathbf{z})}{p_\theta(\mathbf{x}|\mathbf{z})} - \alpha \log \frac{q_\phi(\mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \quad (21.70)$$

$$\frac{6}{7} \quad = \mathbb{E}_{q_\phi(\mathbf{x}, \mathbf{z})} \left[ \log p_\theta(\mathbf{x}|\mathbf{z}) - \log \frac{q_\phi(\mathbf{z})^{\lambda+\alpha-1} p_{\mathcal{D}}(\mathbf{x})}{p_\theta(\mathbf{z})^\lambda q_\phi(\mathbf{z}|\mathbf{x})^{\alpha-1}} \right] \quad (21.71)$$

$$\frac{8}{9} \quad = \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})]] - (1-\alpha) \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}) \| p_\theta(\mathbf{z}))] \\ \frac{10}{11} \quad - (\alpha + \lambda - 1) D_{\text{KL}}(q_\phi(\mathbf{z}) \| p_\theta(\mathbf{z})) - \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [\log p_{\mathcal{D}}(\mathbf{x})] \quad (21.72)$$

11 where the last term is a constant we can ignore. The first two terms can be optimized using the  
 12 reparameterization trick. Unfortunately, the last term requires computing  $q_\phi(\mathbf{z}) = \int_{\mathbf{x}} q_\phi(\mathbf{x}, \mathbf{z}) d\mathbf{x}$ ,  
 13 which is intractable. Fortunately, we can easily sample from this distribution, by sampling  $\mathbf{x} \sim p_{\mathcal{D}}(\mathbf{x})$   
 14 and  $\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})$ . Thus  $q_\phi(\mathbf{z})$  is an **implicit probability model**, similar to a GAN (see Chapter 26).  
 15

16 As long as we use a strict divergence, meaning  $D(q, p) = 0$  iff  $q = p$ , then one can show that this  
 17 does not affect the optimality of the procedure. In particular, proposition 2 of [ZSE19] tells us the  
 18 following:

19 **Theorem 1.** Let  $\mathcal{X}$  and  $\mathcal{Z}$  be continuous spaces, and  $\alpha < 1$  (to bound the MI) and  $\lambda > 0$ . For any  
 20 fixed value of  $\mathbb{I}_q(\mathbf{x}; \mathbf{z})$ , the approximate InfoVAE loss, with any strict divergence  $D(q_\phi(\mathbf{z}), p_\theta(\mathbf{z}))$ , is  
 21 globally optimized if  $p_\theta(\mathbf{x}) = p_{\mathcal{D}}(\mathbf{x})$  and  $q_\phi(\mathbf{z}|\mathbf{x}) = p_\theta(\mathbf{z}|\mathbf{x})$ .  
 22

### 23 21.3.2.1 Connection with MMD VAE

24 If we set  $\alpha = 1$ , the InfoVAE objective simplifies to

$$\frac{27}{28} \quad \mathbb{L} \stackrel{c}{=} \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})]] - \lambda D_{\text{KL}}(q_\phi(\mathbf{z}) \| p_\theta(\mathbf{z})) \quad (21.73)$$

29 The **MMD VAE**<sup>3</sup> replaces the KL divergence in the above term with the (squared) maximum mean  
 30 discrepancy or **MMD** divergence defined in Section 2.7.3. (This is valid based on the above theorem.)  
 31 The advantage of this approach over standard InfoVAE is that the resulting objective is tractable. In  
 32 particular, if we set  $\lambda = 1$  and swap the sign we get

$$\frac{33}{34} \quad \mathcal{L} = \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [-\log p_\theta(\mathbf{x}|\mathbf{z})]] + \text{MMD}(q_\phi(\mathbf{z}), p_\theta(\mathbf{z})) \quad (21.74)$$

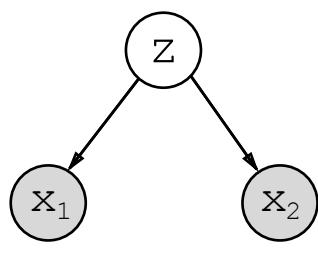
35 As we discuss in Section 2.7.3, we can compute the MMD as follows:  
 36

$$\frac{37}{38} \quad \text{MMD}(p, q) = \mathbb{E}_{p(\mathbf{z}), p(\mathbf{z}')} [\mathcal{K}(\mathbf{z}, \mathbf{z}')] + \mathbb{E}_{q(\mathbf{z}), q(\mathbf{z}')} [\mathcal{K}(\mathbf{z}, \mathbf{z}')] - 2 \mathbb{E}_{p(\mathbf{z}), q(\mathbf{z}')} [\mathcal{K}(\mathbf{z}, \mathbf{z}')] \quad (21.75)$$

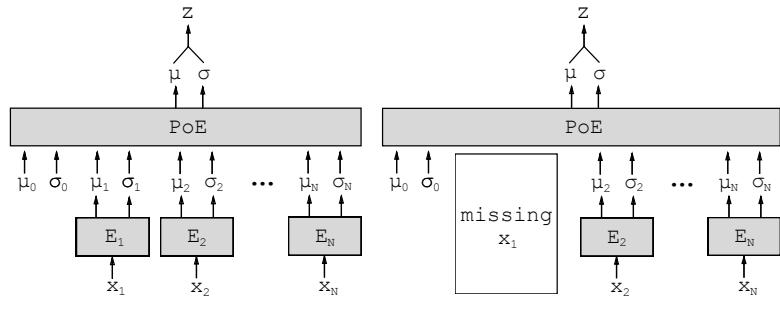
39 where  $\mathcal{K}()$  is some kernel function, such as the RBF kernel,  $\mathcal{K}(\mathbf{z}, \mathbf{z}') = \exp(-\frac{1}{2\sigma^2} \|\mathbf{z} - \mathbf{z}'\|_2^2)$ . Intuitively  
 40 the MMD measures the similarity (in latent space) between samples from the prior and samples from  
 41 the aggregated posterior.

42 In practice, we can implement the MMD objective by using the posterior predicted mean  $\mathbf{z}_n =$   
 43  $e_\phi(\mathbf{x}_n)$  for all  $B$  samples in the current minibatch, and comparing this to  $B$  random samples from  
 44 the  $\mathcal{N}(\mathbf{0}, \mathbf{I})$  prior.  
 45

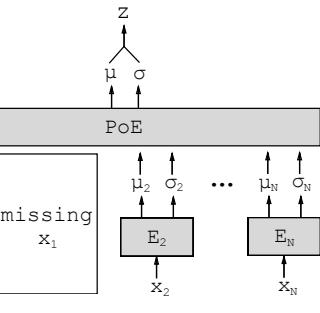
46 3. Proposed in <https://ermongroup.github.io/blog/a-tutorial-on-mmd-variational-autoencoders/>.



(a)



(b)



(c)

Figure 21.6: Illustration of multi-modal VAE. (a) The generative model with  $N = 2$  modalities. (b) The product of experts (PoE) inference network is derived from  $N$  individual Gaussian experts  $E_i$ .  $\mu_0$  and  $\sigma_0$  are parameters of the prior. (c) If a modality is missing, we omit its contribution to the posterior. From Figure 1 of [WG18]. Used with kind permission of Mike Wu.

If we use a Gaussian decoder with fixed variance, the negative log likelihood is just a squared error term:

$$-\log p_{\theta}(\mathbf{x}|\mathbf{z}) = \|\mathbf{x} - d_{\theta}(\mathbf{z})\|_2^2 \quad (21.76)$$

Thus the entire model is deterministic, and just predicts the means in latent space and visible space.

### 21.3.2.2 Connection with $\beta$ -VAEs

If we set  $\alpha = 0$  and  $\lambda = 1$ , we get back the original ELBO. If  $\lambda > 0$  is freely chosen, but we use  $\alpha = 1 - \lambda$ , we get the  $\beta$ -VAE.

### 21.3.2.3 Connection with adversarial autoencoders

If we set  $\alpha = 1$  and  $\lambda = 1$ , and  $D$  is chosen to be the Jensen Shannon divergence (which can be minimized by training a binary discriminator, as explained in Section 26.2.2), then we get a model known as an **adversarial autoencoder** [Mak+15a].

## 21.3.3 Multi-modal VAEs

It is possible to extend VAEs to create joint distributions over different kinds of variables, such as images and text. This is sometimes called a **multimodal VAE** or **MVAE**. Let us assume there are  $M$  modalities. We assume they are conditionally independent given the latent code, and hence the generative model has the form

$$p_{\theta}(\mathbf{x}_1, \dots, \mathbf{x}_M, \mathbf{z}) = p(\mathbf{z}) \prod_{m=1}^M p_{\theta}(\mathbf{x}_m | \mathbf{z}) \quad (21.77)$$

where we treat  $p(\mathbf{z})$  as a fixed prior. See Figure 21.6(a) for an illustration.

1      The standard ELBO is given by  
2

3  
4       $\mathcal{L}_{\theta, \phi}(\mathbf{X}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{X})} \left[ \sum_m \log p_{\theta}(\mathbf{x}_m|\mathbf{z}) \right] - D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{X}) \parallel p(\mathbf{z}))$       (21.78)  
5  
6

7 where  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_M)$  is the observed data. However, the different likelihood terms  $p(\mathbf{x}_m|\mathbf{z})$  may  
8 have different dynamic ranges (e.g., Gaussian pdf for pixels, and categorical pmf for text), so we  
9 introduce weight terms  $\lambda_m \geq 0$  for each likelihood. In addition, let  $\beta \geq 0$  control the amount of KL  
10 regularization. This gives us a weighted version of the ELBO, as follows:

11  
12       $\mathcal{L}_{\theta, \phi}(\mathbf{X}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{X})} \left[ \sum_m \lambda_m \log p_{\theta}(\mathbf{x}_m|\mathbf{z}) \right] - \beta D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{X}) \parallel p(\mathbf{z}))$       (21.79)  
13  
14

15 Often we don't have a lot of paired (aligned) data from all  $M$  modalities. For example, we may  
16 have a lot of images (modality 1), and a lot of text (modality 2), but very few (image, text) pairs.  
17 So it is useful to generalize the loss so it fits the marginal distributions of subsets of the features. Let  
18  $O_m = 1$  if modality  $m$  is observed (i.e.,  $\mathbf{x}_m$  is known), and let  $O_m = 0$  if it is missing or unobserved.  
19 Let  $\mathbf{X} = \{\mathbf{x}_m : O_m = 1\}$  be the visible features. We now use the following objective:  
20

21  
22       $\mathcal{L}_{\theta, \phi}(\mathbf{X}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{X})} \left[ \sum_{m:O_m=1} \lambda_m \log p_{\theta}(\mathbf{x}_m|\mathbf{z}) \right] - \beta D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{X}) \parallel p(\mathbf{z}))$       (21.80)  
23

24      The key problem is how to compute the posterior  $q_{\phi}(\mathbf{z}|\mathbf{X})$  given different subsets of features. In  
25 general this can be hard, since the inference network is a discriminative model that assumes all  
26 inputs are available. For example, if it is trained on (image, text) pairs,  $q_{\phi}(\mathbf{z}|\mathbf{x}_1, \mathbf{x}_2)$ , how can we  
27 compute the posterior just given an image,  $q_{\phi}(\mathbf{z}|\mathbf{x}_1)$ , or just given text,  $q_{\phi}(\mathbf{z}|\mathbf{x}_2)$ ? (This issue arises  
28 in general with VAE when we have missing inputs; we discuss the general case in Section 21.3.4.)  
29

30      Fortunately, based on our conditional independence assumption between the modalities, we can  
31 compute the optimal form for  $q_{\phi}(\mathbf{z}|\mathbf{X})$  given set of inputs by computing the exact posterior under  
32 the model, which is given by

33  
34       $p(\mathbf{z}|\mathbf{X}) = \frac{p(\mathbf{z})p(\mathbf{x}_1, \dots, \mathbf{x}_M|\mathbf{z})}{p(\mathbf{x}_1, \dots, \mathbf{x}_M)} = \frac{p(\mathbf{z})}{p(\mathbf{x}_1, \dots, \mathbf{x}_M)} \prod_{m=1}^M p(\mathbf{x}_m|\mathbf{z})$       (21.81)  
35  
36

37  
38       $= \frac{p(\mathbf{z})}{p(\mathbf{x}_1, \dots, \mathbf{x}_M)} \prod_{m=1}^M \frac{p(\mathbf{z}|\mathbf{x}_m)p(\mathbf{x}_m)}{p(\mathbf{z})}$       (21.82)  
39

40  
41       $\propto p(\mathbf{z}) \prod_{m=1}^M \frac{p(\mathbf{z}|\mathbf{x}_m)}{p(\mathbf{z})} \approx p(\mathbf{z}) \prod_{m=1}^M \tilde{q}(\mathbf{z}|\mathbf{x}_m)$       (21.83)  
42

43 This can be viewed as a product of experts (Section 24.1.1), where each  $\tilde{q}(\mathbf{z}|\mathbf{x}_m)$  is an "expert" for  
44 the  $m$ 'th modality, and  $p(\mathbf{z})$  is the prior. We can compute the above posterior for any subset of  
45 modalities for which we have data by modifying the product over  $m$ . If we use Gaussian distributions  
46 for the prior  $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_0, \boldsymbol{\Lambda}_0^{-1})$  and marginal posterior ratio  $\tilde{q}(\mathbf{z}|\mathbf{x}_m) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_m, \boldsymbol{\Lambda}_m^{-1})$ , then we  
47

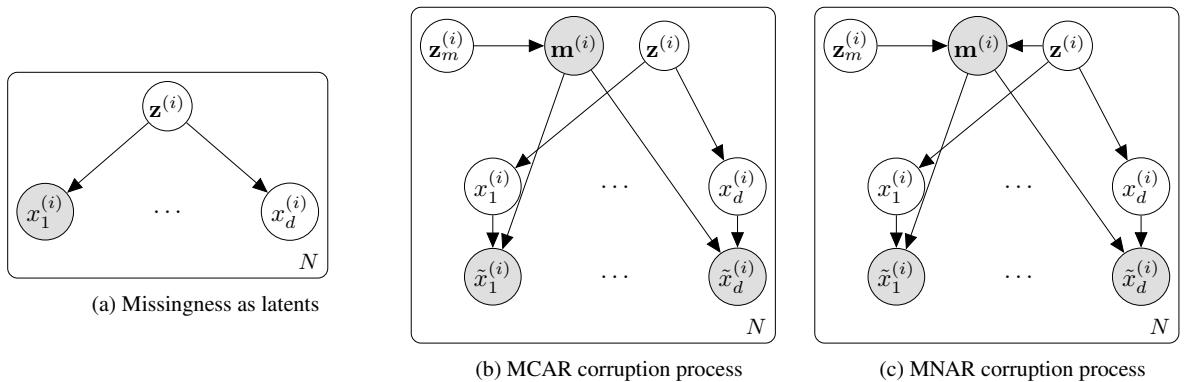


Figure 21.7: Illustration of different VAE variants for handling missing data. From Figure 1 of [CNW20]. Used with kind permission of Mark Collier.

can compute the product of Gaussians using the result from Equation (2.115):

$$\prod_{m=0}^M \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_m, \boldsymbol{\Lambda}_m^{-1}) \propto \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma}), \quad \boldsymbol{\Sigma} = (\sum_m \boldsymbol{\Lambda}_m)^{-1}, \quad \boldsymbol{\mu} = \boldsymbol{\Sigma}(\sum_m \boldsymbol{\Lambda}_m \boldsymbol{\mu}_m) \quad (21.84)$$

Thus the overall posterior precision is the sum of individual expert posterior precisions, and the overall posterior mean is the precision weighted average of the individual expert posterior means. See Figure 21.6(b) for an illustration. For a linear Gaussian (factor analysis) model, we can ensure  $q(\mathbf{z}|\mathbf{x}_m) = p(\mathbf{z}|\mathbf{x}_m)$ , in which case the above solution is the exact posterior [WN18], but in general it will be an approximation.

We need to train the individual expert recognition models  $q(\mathbf{z}|\mathbf{x}_m)$  as well as the joint model  $q(\mathbf{z}|\mathbf{X})$ , so the model knows what to do with fully observed as well as partially observed inputs at test time. In [Ved+18], they propose a somewhat complex “triple ELBO” objective. In [WG18], they propose the simpler approach of optimizing the ELBO for the fully observed feature vector, all the marginals, and a set of  $\mathcal{J}$  randomly chosen joint modalities:

$$\mathbb{L}_{\boldsymbol{\theta}, \boldsymbol{\phi}}(\mathbf{X}) = \mathbb{L}_{\boldsymbol{\theta}, \boldsymbol{\phi}}(\mathbf{x}_1, \dots, \mathbf{x}_M) + \sum_{m=1}^M \mathbb{L}_{\boldsymbol{\theta}, \boldsymbol{\phi}}(\mathbf{x}_m) + \sum_{j \in \mathcal{J}} \mathbb{L}_{\boldsymbol{\theta}, \boldsymbol{\phi}}(\mathbf{X}_j) \quad (21.85)$$

This generalizes nicely to the semi-supervised setting, in which we only have a few aligned (“labeled”) examples from the joint, but have many unaligned (“unlabeled”) examples from the individual marginals. See Figure 21.6(c) for an illustration.

Note that the above scheme can only handle the case of a fixed number of missingness patterns; we generalize to allow for arbitrary missingness in Section 21.3.4.

### 21.3.4 VAEs with missing data

Sometimes we may have **missing data**, in which parts of the data vector  $\mathbf{x} \in \mathbb{R}^D$  may be unknown. In Section 21.3.3 we saw a special case of this when we discussed multimodal VAEs. In this section

1 we allow for arbitrary patterns of missingness.  
2

3 To model the missing data, let  $\mathbf{m} \in \{0, 1\}^D$  be a binary vector where  $m_j = 1$  if  $x_j$  is missing, and  
4  $m_j = 0$  otherwise. Let  $\mathbf{X} = \{\mathbf{x}^{(n)}\}$  and  $\mathbf{M} = \{\mathbf{m}^{(n)}\}$  be  $N \times D$  matrices. Furthermore, let  $\mathbf{X}_o$  be  
5 the observed parts of  $\mathbf{X}$  and  $\mathbf{X}_h$  be the hidden parts. If we assume  $p(\mathbf{M}|\mathbf{X}_o, \mathbf{X}_h) = p(\mathbf{M})$ , we say the  
6 data is **missing completely at random** or **MCAR**, since the missingness does not depend on the  
7 hidden or observed features. If we assume  $p(\mathbf{M}|\mathbf{X}_o, \mathbf{X}_h) = p(\mathbf{M}|\mathbf{X}_o)$ , we say the data is **missing at**  
8 **random** or **MAR**, since the missingness does not depend on the hidden features, but may depend  
9 on the visible features. If neither of these assumptions hold, we say the data is **not missing at**  
10 **random** or **NMAR**.

11 In the MCAR and MAR cases, we can ignore the missingness mechanism, since it tells us nothing  
12 about the hidden features. However, in the NMAR case, we need to model the **missing data**  
13 **mechanism**, since the lack of information may be informative. For example, the fact that someone  
14 did not fill out an answer to a sensitive question on a survey (e.g., “Do you have COVID?”) could be  
15 informative about the underlying value. See e.g., [LR87; Mar08] for more information on missing  
16 data models.

17 In the context of VAEs, we can model the MCAR scenario by treating the missing values as latent  
18 variables. This is illustrated in Figure 21.7(a). Since missing leaf nodes in a directed graphical model  
19 do not affect their parents, we can simply ignore them when computing the posterior  $p(\mathbf{z}^{(i)}|\mathbf{x}_o^{(i)})$ ,  
20 where  $\mathbf{x}_o^{(i)}$  are the observed parts of example  $i$ . However, when using an amortized inference network,  
21 it can be difficult to handle missing inputs, since the model is usually trained to compute  $p(\mathbf{z}^{(i)}|\mathbf{x}_{1:d}^{(i)})$ .  
22 One solution to this is to use the product of experts approach discussed in the context of multi-modal  
23 VAEs in Section 21.3.3. However, this is designed for the case where whole blocks (corresponding to  
24 different modalities) are missing, and will not work well if there are arbitrary missing patterns (e.g.,  
25 pixels that get dropped out due to occlusion or scratches on the lens). In addition, this method will  
26 not work for the NMAR case.

27 An alternative approach, proposed in [CNW20], is to explicitly include the missingness indicators  
28 into the model, as shown in Figure 21.7(b). We assume the model always generates each  $\mathbf{x}_j$  for  
29  $j = 1 : d$ , but we only get to see the “corrupted” versions  $\tilde{\mathbf{x}}_j$ . If  $m_j = 0$  then  $\tilde{\mathbf{x}}_j = \mathbf{x}_j$ , but if  $m_j = 1$ ,  
30 then  $\tilde{\mathbf{x}}_j$  is a special value, such as 0, unrelated to  $\mathbf{x}_j$ . We can model any correlation between the  
31 missingness elements (components of  $\mathbf{m}$ ) by using another latent variable  $\mathbf{z}_m$ . This model can easily  
32 be extended to the NMAR case by letting  $\mathbf{m}$  depend on the latent factors for the observed data,  $\mathbf{z}$ ,  
33 as well as the usual missingness latent factors  $\mathbf{z}_m$ , as shown in Figure 21.7(c).

34 We modify the VAE to be conditional on the missingness pattern, so the VAE decoder has the  
35 form  $p(\mathbf{x}_o|\mathbf{z}, \mathbf{m})$ , and the encoder has the form  $q(\mathbf{z}|\mathbf{x}_o, \mathbf{m})$ . However, we assume the prior is  $p(\mathbf{z})$   
36 as usual, independent of  $\mathbf{m}$ . We can compute a lower bound on the log marginal likelihood of the  
37

38

39

40

41

42

43

44

45

46

47

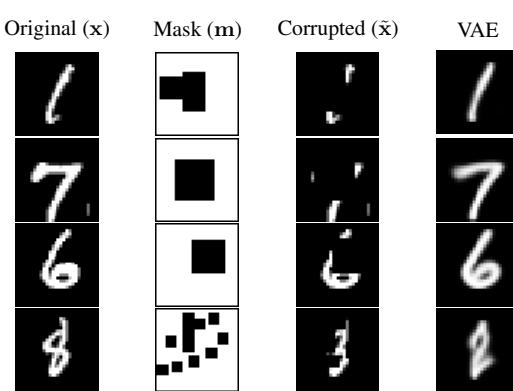


Figure 21.8: Imputing missing pixels given a masked out image using a VAE using a MCAR assumption. From Figure 2 of [CNW20]. Used with kind permission of Mark Collier.

observed data, given the missingness, as follows:

$$\log p(\mathbf{x}_o | \mathbf{m}) = \log \int \int p(\mathbf{x}_o, \mathbf{x}_m | \mathbf{z}, \mathbf{m}) p(\mathbf{z}) d\mathbf{x}_m d\mathbf{z} \quad (21.86)$$

$$= \log \int p(\mathbf{x}_o | \mathbf{z}, \mathbf{m}) p(\mathbf{z}) d\mathbf{z} \quad (21.87)$$

$$= \log \int p(\mathbf{x}_o | \mathbf{z}, \mathbf{m}) p(\mathbf{z}) \frac{q(\mathbf{z} | \tilde{\mathbf{x}}, \mathbf{m})}{q(\mathbf{z} | \tilde{\mathbf{x}})} d\mathbf{z} \quad (21.88)$$

$$= \log \mathbb{E}_{q(\mathbf{z} | \tilde{\mathbf{x}}, \mathbf{m})} \left[ p(\mathbf{x}_o | \mathbf{z}, \mathbf{m}) \frac{p(\mathbf{z})}{q(\mathbf{z} | \tilde{\mathbf{x}}, \mathbf{m})} \right] \quad (21.89)$$

$$\geq \mathbb{E}_{q(\mathbf{z} | \tilde{\mathbf{x}}, \mathbf{m})} [\log p(\mathbf{x}_o | \mathbf{z}, \mathbf{m})] - D_{\text{KL}}(q(\mathbf{z} | \tilde{\mathbf{x}}, \mathbf{m}) \| p(\mathbf{z})) \quad (21.90)$$

We can fit this model in the usual way. See Figure 21.8 for an example.

### 21.3.5 Semi-supervised VAEs

In this section, we discuss how to extend VAEs to the **semi-supervised learning** setting in which we have both labeled data,  $\mathcal{D}_L = \{(\mathbf{x}_n, y_n)\}$ , and unlabeled data,  $\mathcal{D}_U = \{(\mathbf{x}_n)\}$ . We focus on the **M2** model, proposed in [Kin+14a].

The generative model has the following form:

$$p_{\theta}(\mathbf{x}, y) = p_{\theta}(y)p_{\theta}(\mathbf{x}|y) = p_{\theta}(y) \int p_{\theta}(\mathbf{x}|y, \mathbf{z}) p_{\theta}(\mathbf{z}) d\mathbf{z} \quad (21.91)$$

where  $\mathbf{z}$  is a latent variable,  $p_{\theta}(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I})$  is the latent prior,  $p_{\theta}(y) = \text{Cat}(y | \boldsymbol{\pi})$  the label prior, and  $p_{\theta}(\mathbf{x}|y, \mathbf{z}) = p(\mathbf{x}|f_{\theta}(y, \mathbf{z}))$  is the likelihood, such as a Gaussian, with parameters computed by  $f$  (a deep neural network). The main innovation of this approach is to assume that data is generated according to both a latent class variable  $y$  as well as the continuous latent variable  $\mathbf{z}$ . The class variable  $y$  is observed for labeled data and unobserved for unlabeled data.

1 To compute the likelihood for the *labeled* data,  $p_{\theta}(\mathbf{x}, y)$ , we need to marginalize over  $\mathbf{z}$ , which we  
2 can do by using an inference network of the form  
3

4  $q_{\phi}(\mathbf{z}|y, \mathbf{x}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_{\phi}(y, \mathbf{x}), \text{diag}(\boldsymbol{\sigma}_{\phi}(y, \mathbf{x}))$  (21.92)  
5

6 We then use the following variational lower bound  
7

8  $\log p_{\theta}(\mathbf{x}, y) \geq \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}, y)} [\log p_{\theta}(\mathbf{x}|y, \mathbf{z}) + \log p_{\theta}(y) + \log p_{\theta}(\mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x}, y)] = -\mathcal{L}(\mathbf{x}, y)$  (21.93)  
9

10 as is standard for VAEs (see Section 21.2). The only difference is that we observe two kinds of data:  
11  $\mathbf{x}$  and  $y$ .

12 To compute the likelihood for the *unlabeled* data,  $p_{\theta}(\mathbf{x})$ , we need to marginalize over  $\mathbf{z}$  *and*  $y$ ,  
13 which we can do by using an inference network of the form

14  $q_{\phi}(\mathbf{z}, y|\mathbf{x}) = q_{\phi}(\mathbf{z}|\mathbf{x})q_{\phi}(y|\mathbf{x})$  (21.94)  
15

16  $q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_{\phi}(\mathbf{x}), \text{diag}(\boldsymbol{\sigma}_{\phi}(\mathbf{x}))$  (21.95)  
17

18  $q_{\phi}(y|\mathbf{x}) = \text{Cat}(y|\boldsymbol{\pi}_{\phi}(\mathbf{x}))$  (21.96)

19 Note that  $q_{\phi}(y|\mathbf{x})$  acts like a discriminative classifier, that imputes the missing labels. We then use  
20 the following variational lower bound:

21  $\log p_{\theta}(\mathbf{x}) \geq \mathbb{E}_{q_{\phi}(\mathbf{z}, y|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|y, \mathbf{z}) + \log p_{\theta}(y) + \log p_{\theta}(\mathbf{z}) - \log q_{\phi}(\mathbf{z}, y|\mathbf{x})]$  (21.97)  
22

23  $= -\sum_y q_{\phi}(y|\mathbf{x})\mathcal{L}(\mathbf{x}, y) + \mathbb{H}(q_{\phi}(y|\mathbf{x})) = -\mathcal{U}(\mathbf{x})$  (21.98)  
24

25 Note that the discriminative classifier  $q_{\phi}(y|\mathbf{x})$  is only used to compute the log-likelihood of the  
26 unlabeled data, which is undesirable. We can therefore add an extra classification loss on the  
27 supervised data, to get the following overall objective function:

28  $\mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_L} [\mathcal{L}(\mathbf{x}, y)] + \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_U} [\mathcal{U}(\mathbf{x})] + \alpha \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_L} [-\log q_{\phi}(y|\mathbf{x})]$  (21.99)  
29

30 where  $\mathcal{D}_L$  is the labeled data,  $\mathcal{D}_U$  is the unlabeled data, and  $\alpha$  is a hyperparameter that controls the  
31 relative weight of generative and discriminative learning.  
32

### 33 21.3.6 VAEs with sequential encoders/decoders

34 In this section, we discuss VAEs for sequential data, such as text and biosequences, in which the  
35 data  $\mathbf{x}$  is a variable-length sequence, but we have a fixed-sized latent variable  $\mathbf{z} \in \mathbb{R}^K$ . (We consider  
36 the more general case in which  $\mathbf{z}$  is a variable-length sequence of latents — known as **sequential**  
37 **VAE** or **dynamic VAE** — in Section 29.13.) All we have to do is modify the decoder  $p(\mathbf{x}|\mathbf{z})$  and  
38 encoder  $q(\mathbf{z}|\mathbf{x})$  to work with sequences.

39

#### 40 21.3.6.1 Models

41 If we use an RNN for the encoder and decoder of a VAE, we get a model which is called a **VAE-RNN**,  
42 as proposed in [Bow+16a]. In more detail, the generative model is  $p(\mathbf{z}, \mathbf{x}_{1:T}) = p(\mathbf{z})\text{RNN}(\mathbf{x}_{1:T}|\mathbf{z})$ ,  
43 where  $\mathbf{z}$  can be injected as the initial state of the RNN, or as an input to every time step. The  
44

45

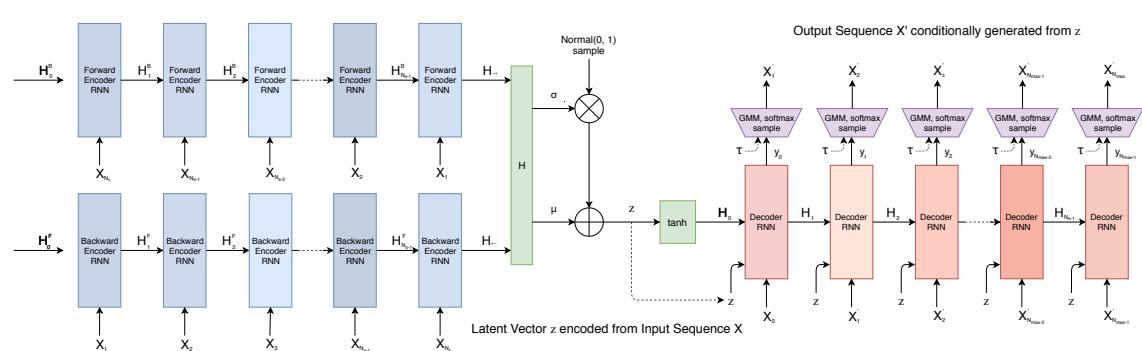


Figure 21.9: Illustration of a VAE with a bidirectional RNN encoder and a unidirectional RNN decoder. The output generator can use a GMM and/or softmax distribution. From Figure 2 of [HE18]. Used with kind permission of David Ha.

inference model is  $q(\mathbf{z}|\mathbf{x}_{1:T}) = \mathcal{N}(\mathbf{z}|\mu(\mathbf{h}), \Sigma(\mathbf{h}))$ , where  $\mathbf{h} = [\mathbf{h}_T^\rightarrow, \mathbf{h}_1^\leftarrow]$  is the output of a bidirectional RNN applied to  $\mathbf{x}_{1:T}$ . See Figure 21.9 for an illustration.

More recently, people have tried to combine transformers with VAEs. For example, in the **Optimus** model of [Li+20], they use a BERT model for the encoder. In more detail, the encoder  $q(\mathbf{z}|\mathbf{x})$  is derived from the embedding vector associated with a dummy token corresponding to the “class label” which is appended to the input sequence  $\mathbf{x}$ . The decoder is a standard autoregressive model (similar to GPT), with one additional input, namely the latent vector  $\mathbf{z}$ . They consider two ways of injecting the latent vector. The simplest approach is to add  $\mathbf{z}$  to the embedding layer of every token in the decoding step, by defining  $\mathbf{h}'_i = \mathbf{h}_i + \mathbf{W}\mathbf{z}$ , where  $\mathbf{h}_i \in \mathbb{R}^H$  is the original embedding for the  $i$ 'th token, and  $\mathbf{W} \in \mathbb{R}^{H \times K}$  is a decoding matrix, where  $K$  is the size of the latent vector. However, they get better results in their experiments by letting all the layers of the decoder attend to the latent code  $\mathbf{z}$ . An easy way to do this is to define the memory vector  $\mathbf{h}_m = \mathbf{W}\mathbf{z}$ , where  $\mathbf{W} \in \mathbb{R}^{LH \times K}$ , where  $L$  is the number of layers in the decoder, and then to append  $\mathbf{h}_m \in \mathbb{R}^{L \times H}$  to all the other embeddings at each layer.

An alternative approach, known as **transformer VAE**, was proposed in [Gre20]. This model uses a **funnel transformer** [Dai+20b] as the encoder, and the T5 [Raf+20a] conditional transformer for the decoder. In addition, it uses an MMD VAE (Section 21.3.2.1) to avoid posterior collapse.

### 21.3.6.2 Applications

In this section, we discuss some applications of VAEs to sequence data.

#### Text

In [Bow+16b], they apply the VAE-RNN model to natural language sentences. (See also [MB16; SSB17] for related work.) Although this does not improve performance in terms of the standard perplexity measures (predicting the next word given the previous words), it does provide a way to infer a semantic representation of the sentence. This can then be used for latent space interpolation, as discussed in Section 20.3.5. The results of doing this with the VAE-RNN are illustrated in

|   |                                   |                                             |
|---|-----------------------------------|---------------------------------------------|
| 1 |                                   | i went to the store to buy some groceries . |
| 2 |                                   | i store to buy some groceries .             |
| 3 | he was silent for a long moment . | i were to buy any groceries .               |
| 4 | he was silent for a moment .      | horses are to buy any groceries .           |
| 5 | it was quiet for a moment .       | horses are to buy any animal .              |
| 6 | it was dark and cold .            | horses the favorite any animal .            |
| 7 | there was a pause .               | horses the favorite favorite animal .       |
| 8 | it was my turn .                  | horses are my favorite animal .             |

(a)

(b)

10 Figure 21.10: (a) Samples from the latent space of a VAE text model, as we interpolate between two sentences  
11 (on first and last line). Note that the intermediate sentences are grammatical, and semantically related to  
12 their neighbors. From Table 8 of [Bow+16b]. (b) Same as (a), but now using a deterministic autoencoder  
13 (with the same RNN encoder and decoder). From Table 1 of [Bow+16b]. Used with kind permission of Sam  
 Bowman.

16 Figure 21.10a. (Similar results are shown in [Li+20], using a VAE-transformer.) By contrast, if  
17 we use a standard deterministic autoencoder, with the same RNN encoder and decoder networks,  
18 we learn a much less meaningful space, as illustrated in Figure 21.10b. The reason is that the  
19 deterministic autoencoder has “holes” in its latent space, which get decoded to nonsensical outputs.  
20

However, because RNNs (and transformers) are powerful decoders, we need to address the problem of posterior collapse, which we discuss in Section 21.4. One common way to avoid this problem is to use KL annealing, but a more effective method is to use the InfoVAE method of Section 21.3.2, which includes adversarial autoencoders (used in [She+20] with an RNN decoder) and MMD autoencoders (used in [Gre20] with a transformer decoder).

26  
25 Sketches

<sup>28</sup> In [HE18], they apply the VAE-RNN model to generate sketches (line drawings) of various animals  
<sup>29</sup> and hand-written characters. They call their model **sketch-rnn**. The training data records the  
<sup>30</sup> sequence of  $(x, y)$  pen positions, as well as whether the pen was touching the paper or not. The  
<sup>31</sup> emission model used a GMM for the real-valued location offsets, and a categorical softmax distribution  
<sup>32</sup> for the discrete state.

<sup>33</sup> Figure 21.11 shows some samples from various class-conditional models. We vary the temperature  
<sup>34</sup> parameter  $\tau$  of the emission model to control the stochasticity of the generator. (More precisely, we  
<sup>35</sup> multiply the GMM variances by  $\tau$ , and divide the discrete probabilities by  $\tau$  before renormalizing.)  
<sup>36</sup> When the temperature is low, the model tries to reconstruct the input as closely as possible. However,  
<sup>37</sup> when the input is untypical of the training set (e.g., a cat with three eyes, or a toothbrush), the  
<sup>38</sup> reconstruction is “regularized” towards a canonical cat with two eyes, while still keeping some features  
<sup>39</sup> of the input.

41

<sup>43</sup> In [GB+18], they use VAE-RNNs to model molecular graph structure, represented as a string using  
<sup>44</sup> the SMILES representation.<sup>4</sup> It is also possible to learn a mapping from the latent space to some

<sup>46</sup> See [https://en.wikipedia.org/wiki/Simplified\\_molecular-input\\_line-entry\\_system](https://en.wikipedia.org/wiki/Simplified_molecular-input_line-entry_system)