

Understanding Deep Learning Equations

January 23, 2025

Chapter 1

Introduction

Chapter 2

Supervised Learning

$$\mathbf{y} = \mathbf{f}[\mathbf{x}]. \quad (2.1)$$

$$\mathbf{y} = \mathbf{f}[\mathbf{x}, \phi]. \quad (2.2)$$

$$\hat{\phi} = \underset{\phi}{\operatorname{argmin}} [L[\phi]]. \quad (2.3)$$

$$\begin{aligned} y &= \mathbf{f}[x, \phi] \\ &= \phi_0 + \phi_1 x. \end{aligned} \quad (2.4)$$

$$\begin{aligned} L[\phi] &= \sum_{i=1}^I (\mathbf{f}[x_i, \phi] - y_i)^2 \\ &= \sum_{i=1}^I (\phi_0 + \phi_1 x_i - y_i)^2. \end{aligned} \quad (2.5)$$

$$\begin{aligned} \hat{\phi} &= \underset{\phi}{\operatorname{argmin}} [L[\phi]] \\ &= \underset{\phi}{\operatorname{argmin}} \left[\sum_{i=1}^I (\mathbf{f}[x_i, \phi] - y_i)^2 \right] \\ &= \underset{\phi}{\operatorname{argmin}} \left[\sum_{i=1}^I (\phi_0 + \phi_1 x_i - y_i)^2 \right]. \end{aligned} \quad (2.6)$$

Chapter 3

Shallow neural networks

$$\begin{aligned} y &= f[x, \phi] \\ &= \phi_0 + \phi_1 a[\theta_{10} + \theta_{11}x] + \phi_2 a[\theta_{20} + \theta_{21}x] + \phi_3 a[\theta_{30} + \theta_{31}x]. \end{aligned} \quad (3.1)$$

$$a[z] = \text{ReLU}[z] = \begin{cases} 0 & z < 0 \\ z & z \geq 0 \end{cases}. \quad (3.2)$$

$$\begin{aligned} h_1 &= a[\theta_{10} + \theta_{11}x] \\ h_2 &= a[\theta_{20} + \theta_{21}x] \\ h_3 &= a[\theta_{30} + \theta_{31}x], \end{aligned} \quad (3.3)$$

$$y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3. \quad (3.4)$$

$$h_d = a[\theta_{d0} + \theta_{d1}x], \quad (3.5)$$

$$y = \phi_0 + \sum_{d=1}^D \phi_d h_d. \quad (3.6)$$

$$\begin{aligned} h_1 &= a[\theta_{10} + \theta_{11}x] \\ h_2 &= a[\theta_{20} + \theta_{21}x] \\ h_3 &= a[\theta_{30} + \theta_{31}x] \\ h_4 &= a[\theta_{40} + \theta_{41}x], \end{aligned} \quad (3.7)$$

$$\begin{aligned}
y_1 &= \phi_{10} + \phi_{11}h_1 + \phi_{12}h_2 + \phi_{13}h_3 + \phi_{14}h_4 \\
y_2 &= \phi_{20} + \phi_{21}h_1 + \phi_{22}h_2 + \phi_{23}h_3 + \phi_{24}h_4.
\end{aligned} \tag{3.8}$$

$$\begin{aligned}
h_1 &= \text{a}[\theta_{10} + \theta_{11}x_1 + \theta_{12}x_2] \\
h_2 &= \text{a}[\theta_{20} + \theta_{21}x_1 + \theta_{22}x_2] \\
h_3 &= \text{a}[\theta_{30} + \theta_{31}x_1 + \theta_{32}x_2],
\end{aligned} \tag{3.9}$$

$$y = \phi_0 + \phi_1h_1 + \phi_2h_2 + \phi_3h_3. \tag{3.10}$$

$$h_d = \text{a} \left[\theta_{d0} + \sum_{i=1}^{D_i} \theta_{di}x_i \right], \tag{3.11}$$

$$y_j = \phi_{j0} + \sum_{d=1}^D \phi_{jd}h_d, \tag{3.12}$$

$$\text{HardSwish}[z] = \begin{cases} 0 & z < -3 \\ z(z+3)/6 & -3 \leq z \leq 3 \\ z & z > 3 \end{cases}. \tag{3.13}$$

$$\text{ReLU}[\alpha \cdot z] = \alpha \cdot \text{ReLU}[z]. \tag{3.14}$$

$$\text{heaviside}[z] = \begin{cases} 0 & z < 0 \\ 1 & z \geq 0 \end{cases} \quad \text{rect}[z] = \begin{cases} 0 & z < 0 \\ 1 & 0 \leq z \leq 1 \\ 0 & z > 1 \end{cases}. \tag{3.15}$$

Chapter 4

Deep neural networks

$$\begin{aligned}h_1 &= \mathbf{a}[\theta_{10} + \theta_{11}x] \\h_2 &= \mathbf{a}[\theta_{20} + \theta_{21}x] \\h_3 &= \mathbf{a}[\theta_{30} + \theta_{31}x],\end{aligned}\tag{4.1}$$

$$y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3.\tag{4.2}$$

$$\begin{aligned}h'_1 &= \mathbf{a}[\theta'_{10} + \theta'_{11}y] \\h'_2 &= \mathbf{a}[\theta'_{20} + \theta'_{21}y] \\h'_3 &= \mathbf{a}[\theta'_{30} + \theta'_{31}y],\end{aligned}\tag{4.3}$$

$$y' = \phi'_0 + \phi'_1 h'_1 + \phi'_2 h'_2 + \phi'_3 h'_3.\tag{4.4}$$

$$\begin{aligned}h'_1 &= \mathbf{a}[\theta'_{10} + \theta'_{11}y] = \mathbf{a}[\theta'_{10} + \theta'_{11}\phi_0 + \theta'_{11}\phi_1 h_1 + \theta'_{11}\phi_2 h_2 + \theta'_{11}\phi_3 h_3] \\h'_2 &= \mathbf{a}[\theta'_{20} + \theta'_{21}y] = \mathbf{a}[\theta'_{20} + \theta'_{21}\phi_0 + \theta'_{21}\phi_1 h_1 + \theta'_{21}\phi_2 h_2 + \theta'_{21}\phi_3 h_3] \\h'_3 &= \mathbf{a}[\theta'_{30} + \theta'_{31}y] = \mathbf{a}[\theta'_{30} + \theta'_{31}\phi_0 + \theta'_{31}\phi_1 h_1 + \theta'_{31}\phi_2 h_2 + \theta'_{31}\phi_3 h_3],\end{aligned}\tag{4.5}$$

$$\begin{aligned}h'_1 &= \mathbf{a}[\psi_{10} + \psi_{11}h_1 + \psi_{12}h_2 + \psi_{13}h_3] \\h'_2 &= \mathbf{a}[\psi_{20} + \psi_{21}h_1 + \psi_{22}h_2 + \psi_{23}h_3] \\h'_3 &= \mathbf{a}[\psi_{30} + \psi_{31}h_1 + \psi_{32}h_2 + \psi_{33}h_3],\end{aligned}\tag{4.6}$$

$$\begin{aligned}
h_1 &= \mathbf{a}[\theta_{10} + \theta_{11}x] \\
h_2 &= \mathbf{a}[\theta_{20} + \theta_{21}x] \\
h_3 &= \mathbf{a}[\theta_{30} + \theta_{31}x],
\end{aligned} \tag{4.7}$$

$$\begin{aligned}
h'_1 &= \mathbf{a}[\psi_{10} + \psi_{11}h_1 + \psi_{12}h_2 + \psi_{13}h_3] \\
h'_2 &= \mathbf{a}[\psi_{20} + \psi_{21}h_1 + \psi_{22}h_2 + \psi_{23}h_3] \\
h'_3 &= \mathbf{a}[\psi_{30} + \psi_{31}h_1 + \psi_{32}h_2 + \psi_{33}h_3],
\end{aligned} \tag{4.8}$$

$$y' = \phi'_0 + \phi'_1 h'_1 + \phi'_2 h'_2 + \phi'_3 h'_3. \tag{4.9}$$

$$\begin{aligned}
y' &= \phi'_0 + \phi'_1 \mathbf{a}[\psi_{10} + \psi_{11}\mathbf{a}[\theta_{10} + \theta_{11}x] + \psi_{12}\mathbf{a}[\theta_{20} + \theta_{21}x] + \psi_{13}\mathbf{a}[\theta_{30} + \theta_{31}x]] \\
&\quad + \phi'_2 \mathbf{a}[\psi_{20} + \psi_{21}\mathbf{a}[\theta_{10} + \theta_{11}x] + \psi_{22}\mathbf{a}[\theta_{20} + \theta_{21}x] + \psi_{23}\mathbf{a}[\theta_{30} + \theta_{31}x]] \\
&\quad + \phi'_3 \mathbf{a}[\psi_{30} + \psi_{31}\mathbf{a}[\theta_{10} + \theta_{11}x] + \psi_{32}\mathbf{a}[\theta_{20} + \theta_{21}x] + \psi_{33}\mathbf{a}[\theta_{30} + \theta_{31}x]],
\end{aligned} \tag{4.10}$$

$$\begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} = \mathbf{a} \left[\begin{bmatrix} \theta_{10} \\ \theta_{20} \\ \theta_{30} \end{bmatrix} + \begin{bmatrix} \theta_{11} \\ \theta_{21} \\ \theta_{31} \end{bmatrix} x \right], \tag{4.11}$$

$$\begin{bmatrix} h'_1 \\ h'_2 \\ h'_3 \end{bmatrix} = \mathbf{a} \left[\begin{bmatrix} \psi_{10} \\ \psi_{20} \\ \psi_{30} \end{bmatrix} + \begin{bmatrix} \psi_{11} & \psi_{12} & \psi_{13} \\ \psi_{21} & \psi_{22} & \psi_{23} \\ \psi_{31} & \psi_{32} & \psi_{33} \end{bmatrix} \begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} \right], \tag{4.12}$$

$$y' = \phi'_0 + [\phi'_1 \quad \phi'_2 \quad \phi'_3] \begin{bmatrix} h'_1 \\ h'_2 \\ h'_3 \end{bmatrix}, \tag{4.13}$$

$$\begin{aligned}
\mathbf{h} &= \mathbf{a}[\boldsymbol{\theta}_0 + \boldsymbol{\theta}x] \\
\mathbf{h}' &= \mathbf{a}[\boldsymbol{\psi}_0 + \boldsymbol{\Psi}\mathbf{h}] \\
y' &= \phi'_0 + \boldsymbol{\phi}'\mathbf{h}',
\end{aligned} \tag{4.14}$$

$$\begin{aligned}
\mathbf{h}_1 &= \mathbf{a}[\boldsymbol{\beta}_0 + \boldsymbol{\Omega}_0 \mathbf{x}] \\
\mathbf{h}_2 &= \mathbf{a}[\boldsymbol{\beta}_1 + \boldsymbol{\Omega}_1 \mathbf{h}_1] \\
\mathbf{h}_3 &= \mathbf{a}[\boldsymbol{\beta}_2 + \boldsymbol{\Omega}_2 \mathbf{h}_2] \\
&\vdots \\
\mathbf{h}_K &= \mathbf{a}[\boldsymbol{\beta}_{K-1} + \boldsymbol{\Omega}_{K-1} \mathbf{h}_{K-1}] \\
\mathbf{y} &= \boldsymbol{\beta}_K + \boldsymbol{\Omega}_K \mathbf{h}_K.
\end{aligned} \tag{4.15}$$

$$\mathbf{y} = \boldsymbol{\beta}_K + \boldsymbol{\Omega}_K \mathbf{a} \left[\boldsymbol{\beta}_{K-1} + \boldsymbol{\Omega}_{K-1} \mathbf{a} \left[\dots \boldsymbol{\beta}_2 + \boldsymbol{\Omega}_2 \mathbf{a} \left[\boldsymbol{\beta}_1 + \boldsymbol{\Omega}_1 \mathbf{a} \left[\boldsymbol{\beta}_0 + \boldsymbol{\Omega}_0 \mathbf{x} \right] \dots \right] \right] \right]. \tag{4.16}$$

$$N_r = \left(\frac{D}{D_i} + 1 \right)^{D_i(K-1)} \cdot \sum_{j=0}^{D_i} \binom{D}{j}. \tag{4.17}$$

$$\text{ReLU} \left[\boldsymbol{\beta}_1 + \lambda_1 \cdot \boldsymbol{\Omega}_1 \text{ReLU} \left[\boldsymbol{\beta}_0 + \lambda_0 \cdot \boldsymbol{\Omega}_0 \mathbf{x} \right] \right] = \lambda_0 \lambda_1 \cdot \text{ReLU} \left[\frac{1}{\lambda_0 \lambda_1} \boldsymbol{\beta}_1 + \boldsymbol{\Omega}_1 \text{ReLU} \left[\frac{1}{\lambda_0} \boldsymbol{\beta}_0 + \boldsymbol{\Omega}_0 \mathbf{x} \right] \right] \tag{4.18}$$

Chapter 5

Loss functions

$$\begin{aligned}\hat{\phi} &= \operatorname{argmax}_{\phi} \left[\prod_{i=1}^I Pr(\mathbf{y}_i | \mathbf{x}_i) \right] \\ &= \operatorname{argmax}_{\phi} \left[\prod_{i=1}^I Pr(\mathbf{y}_i | \boldsymbol{\theta}_i) \right] \\ &= \operatorname{argmax}_{\phi} \left[\prod_{i=1}^I Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi]) \right].\end{aligned}\tag{5.1}$$

$$Pr(\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_I | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_I) = \prod_{i=1}^I Pr(\mathbf{y}_i | \mathbf{x}_i).\tag{5.2}$$

$$\begin{aligned}\hat{\phi} &= \operatorname{argmax}_{\phi} \left[\prod_{i=1}^I Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi]) \right] \\ &= \operatorname{argmax}_{\phi} \left[\log \left[\prod_{i=1}^I Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi]) \right] \right] \\ &= \operatorname{argmax}_{\phi} \left[\sum_{i=1}^I \log [Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi])] \right].\end{aligned}\tag{5.3}$$

$$\begin{aligned}\hat{\phi} &= \operatorname{argmin}_{\phi} \left[- \sum_{i=1}^I \log [Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi])] \right] \\ &= \operatorname{argmin}_{\phi} [L[\phi]],\end{aligned}\tag{5.4}$$

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} \left[\Pr(\mathbf{y} | \mathbf{f}[\mathbf{x}, \hat{\phi}]) \right]. \quad (5.5)$$

$$\hat{\phi} = \operatorname{argmin}_{\phi} \left[L[\phi] \right] = \operatorname{argmin}_{\phi} \left[- \sum_{i=1}^I \log \left[\Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi]) \right] \right]. \quad (5.6)$$

$$\Pr(y | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{(y - \mu)^2}{2\sigma^2} \right]. \quad (5.7)$$

$$\Pr(y | \mathbf{f}[\mathbf{x}, \phi], \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{(y - \mathbf{f}[\mathbf{x}, \phi])^2}{2\sigma^2} \right]. \quad (5.8)$$

$$\begin{aligned} L[\phi] &= - \sum_{i=1}^I \log \left[\Pr(y_i | \mathbf{f}[\mathbf{x}_i, \phi], \sigma^2) \right] \\ &= - \sum_{i=1}^I \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{(y_i - \mathbf{f}[\mathbf{x}_i, \phi])^2}{2\sigma^2} \right] \right]. \end{aligned} \quad (5.9)$$

$$\begin{aligned} \hat{\phi} &= \operatorname{argmin}_{\phi} \left[- \sum_{i=1}^I \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{(y_i - \mathbf{f}[\mathbf{x}_i, \phi])^2}{2\sigma^2} \right] \right] \right] \\ &= \operatorname{argmin}_{\phi} \left[- \sum_{i=1}^I \left(\log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \right] - \frac{(y_i - \mathbf{f}[\mathbf{x}_i, \phi])^2}{2\sigma^2} \right) \right] \\ &= \operatorname{argmin}_{\phi} \left[- \sum_{i=1}^I - \frac{(y_i - \mathbf{f}[\mathbf{x}_i, \phi])^2}{2\sigma^2} \right] \\ &= \operatorname{argmin}_{\phi} \left[\sum_{i=1}^I (y_i - \mathbf{f}[\mathbf{x}_i, \phi])^2 \right], \end{aligned} \quad (5.10)$$

$$L[\phi] = \sum_{i=1}^I (y_i - \mathbf{f}[\mathbf{x}_i, \phi])^2. \quad (5.11)$$

$$\hat{y} = \operatorname{argmax}_y \left[\Pr(y | \mathbf{f}[\mathbf{x}, \hat{\phi}], \sigma^2) \right]. \quad (5.12)$$

$$\hat{\phi}, \hat{\sigma}^2 = \underset{\phi, \sigma^2}{\operatorname{argmin}} \left[- \sum_{i=1}^I \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[- \frac{(y_i - f[\mathbf{x}_i, \phi])^2}{2\sigma^2} \right] \right] \right]. \quad (5.13)$$

$$\begin{aligned} \mu &= f_1[\mathbf{x}, \phi] \\ \sigma^2 &= f_2[\mathbf{x}, \phi]^2, \end{aligned} \quad (5.14)$$

$$\hat{\phi} = \underset{\phi}{\operatorname{argmin}} \left[- \sum_{i=1}^I \left(\log \left[\frac{1}{\sqrt{2\pi f_2[\mathbf{x}_i, \phi]^2}} \right] - \frac{(y_i - f_1[\mathbf{x}_i, \phi])^2}{2f_2[\mathbf{x}_i, \phi]^2} \right) \right]. \quad (5.15)$$

$$Pr(y|\lambda) = \begin{cases} 1 - \lambda & y = 0 \\ \lambda & y = 1 \end{cases}, \quad (5.16)$$

$$Pr(y|\lambda) = (1 - \lambda)^{1-y} \cdot \lambda^y. \quad (5.17)$$

$$\operatorname{sig}[z] = \frac{1}{1 + \exp[-z]}. \quad (5.18)$$

$$Pr(y|\mathbf{x}) = (1 - \operatorname{sig}[f[\mathbf{x}, \phi]])^{1-y} \cdot \operatorname{sig}[f[\mathbf{x}, \phi]]^y. \quad (5.19)$$

$$L[\phi] = \sum_{i=1}^I -(1 - y_i) \log [1 - \operatorname{sig}[f[\mathbf{x}_i, \phi]]] - y_i \log [\operatorname{sig}[f[\mathbf{x}_i, \phi]]]. \quad (5.20)$$

$$Pr(y = k) = \lambda_k. \quad (5.21)$$

$$\operatorname{softmax}_k[\mathbf{z}] = \frac{\exp[z_k]}{\sum_{k'=1}^K \exp[z_{k'}]}, \quad (5.22)$$

$$Pr(y = k|\mathbf{x}) = \operatorname{softmax}_k[\mathbf{f}[\mathbf{x}, \phi]]. \quad (5.23)$$

$$\begin{aligned}
L[\phi] &= -\sum_{i=1}^I \log \left[\text{softmax}_{y_i} [\mathbf{f}[\mathbf{x}_i, \phi]] \right] \\
&= -\sum_{i=1}^I \left(f_{y_i}[\mathbf{x}_i, \phi] - \log \left[\sum_{k'=1}^K \exp [f_{k'}[\mathbf{x}_i, \phi]] \right] \right), \tag{5.24}
\end{aligned}$$

$$Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}, \phi]) = \prod_d Pr(y_d|\mathbf{f}_d[\mathbf{x}, \phi]), \tag{5.25}$$

$$L[\phi] = -\sum_{i=1}^I \log [Pr(\mathbf{y}_i|\mathbf{f}[\mathbf{x}_i, \phi])] = -\sum_{i=1}^I \sum_d \log [Pr(y_{id}|\mathbf{f}_d[\mathbf{x}_i, \phi])]. \tag{5.26}$$

$$D_{KL}[q||p] = \int_{-\infty}^{\infty} q(z) \log[q(z)] dz - \int_{-\infty}^{\infty} q(z) \log[p(z)] dz. \tag{5.27}$$

$$q(y) = \frac{1}{I} \sum_{i=1}^I \delta[y - y_i], \tag{5.28}$$

$$\begin{aligned}
\hat{\boldsymbol{\theta}} &= \underset{\boldsymbol{\theta}}{\text{argmin}} \left[\int_{-\infty}^{\infty} q(y) \log[q(y)] dy - \int_{-\infty}^{\infty} q(y) \log[Pr(y|\boldsymbol{\theta})] dy \right] \\
&= \underset{\boldsymbol{\theta}}{\text{argmin}} \left[-\int_{-\infty}^{\infty} q(y) \log[Pr(y|\boldsymbol{\theta})] dy \right], \tag{5.29}
\end{aligned}$$

$$\begin{aligned}
\hat{\boldsymbol{\theta}} &= \underset{\boldsymbol{\theta}}{\text{argmin}} \left[-\int_{-\infty}^{\infty} \left(\frac{1}{I} \sum_{i=1}^I \delta[y - y_i] \right) \log[Pr(y|\boldsymbol{\theta})] dy \right] \\
&= \underset{\boldsymbol{\theta}}{\text{argmin}} \left[-\frac{1}{I} \sum_{i=1}^I \log[Pr(y_i|\boldsymbol{\theta})] \right] \\
&= \underset{\boldsymbol{\theta}}{\text{argmin}} \left[-\sum_{i=1}^I \log[Pr(y_i|\boldsymbol{\theta})] \right]. \tag{5.30}
\end{aligned}$$

$$\hat{\phi} = \underset{\phi}{\text{argmin}} \left[-\sum_{i=1}^I \log[Pr(y_i|\mathbf{f}[\mathbf{x}_i, \phi])] \right]. \tag{5.31}$$

$$\text{sig}[z] = \frac{1}{1 + \exp[-z]}. \quad (5.32)$$

$$L = -(1 - y) \log \left[1 - \text{sig}[\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}]] \right] - y \log \left[\text{sig}[\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}]] \right], \quad (5.33)$$

$$Pr(y|\mu, \kappa) = \frac{\exp[\kappa \cos[y - \mu]]}{2\pi \cdot \text{Bessel}_0[\kappa]}, \quad (5.34)$$

$$Pr(y|\lambda, \mu_1, \mu_2, \sigma_1^2, \sigma_2^2) = \frac{\lambda}{\sqrt{2\pi\sigma_1^2}} \exp \left[\frac{-(y - \mu_1)^2}{2\sigma_1^2} \right] + \frac{1 - \lambda}{\sqrt{2\pi\sigma_2^2}} \exp \left[\frac{-(y - \mu_2)^2}{2\sigma_2^2} \right], \quad (5.35)$$

$$Pr(y = k) = \frac{\lambda^k e^{-\lambda}}{k!}. \quad (5.36)$$

Chapter 6

Fitting models

$$\hat{\phi} = \underset{\phi}{\operatorname{argmin}} [L[\phi]]. \quad (6.1)$$

$$\frac{\partial L}{\partial \phi} = \begin{bmatrix} \frac{\partial L}{\partial \phi_0} \\ \frac{\partial L}{\partial \phi_1} \\ \vdots \\ \frac{\partial L}{\partial \phi_N} \end{bmatrix}. \quad (6.2)$$

$$\phi \longleftarrow \phi - \alpha \cdot \frac{\partial L}{\partial \phi}, \quad (6.3)$$

$$\begin{aligned} y &= \mathbf{f}[x, \phi] \\ &= \phi_0 + \phi_1 x. \end{aligned} \quad (6.4)$$

$$\begin{aligned} L[\phi] &= \sum_{i=1}^I \ell_i = \sum_{i=1}^I (\mathbf{f}[x_i, \phi] - y_i)^2 \\ &= \sum_{i=1}^I (\phi_0 + \phi_1 x_i - y_i)^2, \end{aligned} \quad (6.5)$$

$$\frac{\partial L}{\partial \phi} = \frac{\partial}{\partial \phi} \sum_{i=1}^I \ell_i = \sum_{i=1}^I \frac{\partial \ell_i}{\partial \phi}, \quad (6.6)$$

$$\frac{\partial \ell_i}{\partial \boldsymbol{\phi}} = \begin{bmatrix} \frac{\partial \ell_i}{\partial \phi_0} \\ \frac{\partial \ell_i}{\partial \phi_1} \end{bmatrix} = \begin{bmatrix} 2(\phi_0 + \phi_1 x_i - y_i) \\ 2x_i(\phi_0 + \phi_1 x_i - y_i) \end{bmatrix}. \quad (6.7)$$

$$f[x, \boldsymbol{\phi}] = \sin[\phi_0 + 0.06 \cdot \phi_1 x] \cdot \exp\left(-\frac{(\phi_0 + 0.06 \cdot \phi_1 x)^2}{32.0}\right). \quad (6.8)$$

$$L[\boldsymbol{\phi}] = \sum_{i=1}^I (f[x_i, \boldsymbol{\phi}] - y_i)^2. \quad (6.9)$$

$$\phi_{t+1} \leftarrow \phi_t - \alpha \cdot \sum_{i \in \mathcal{B}_t} \frac{\partial \ell_i[\phi_t]}{\partial \phi}, \quad (6.10)$$

$$\begin{aligned} \mathbf{m}_{t+1} &\leftarrow \beta \cdot \mathbf{m}_t + (1 - \beta) \sum_{i \in \mathcal{B}_t} \frac{\partial \ell_i[\phi_t]}{\partial \phi} \\ \phi_{t+1} &\leftarrow \phi_t - \alpha \cdot \mathbf{m}_{t+1}, \end{aligned} \quad (6.11)$$

$$\begin{aligned} \mathbf{m}_{t+1} &\leftarrow \beta \cdot \mathbf{m}_t + (1 - \beta) \sum_{i \in \mathcal{B}_t} \frac{\partial \ell_i[\phi_t - \alpha \beta \cdot \mathbf{m}_t]}{\partial \phi} \\ \phi_{t+1} &\leftarrow \phi_t - \alpha \cdot \mathbf{m}_{t+1}, \end{aligned} \quad (6.12)$$

$$\begin{aligned} \mathbf{m}_{t+1} &\leftarrow \frac{\partial L[\phi_t]}{\partial \phi} \\ \mathbf{v}_{t+1} &\leftarrow \left(\frac{\partial L[\phi_t]}{\partial \phi} \right)^2. \end{aligned} \quad (6.13)$$

$$\phi_{t+1} \leftarrow \phi_t - \alpha \cdot \frac{\mathbf{m}_{t+1}}{\sqrt{\mathbf{v}_{t+1}} + \epsilon}, \quad (6.14)$$

$$\begin{aligned} \mathbf{m}_{t+1} &\leftarrow \beta \cdot \mathbf{m}_t + (1 - \beta) \frac{\partial L[\phi_t]}{\partial \phi} \\ \mathbf{v}_{t+1} &\leftarrow \gamma \cdot \mathbf{v}_t + (1 - \gamma) \left(\frac{\partial L[\phi_t]}{\partial \phi} \right)^2, \end{aligned} \quad (6.15)$$

$$\begin{aligned}\tilde{\mathbf{m}}_{t+1} &\leftarrow \frac{\mathbf{m}_{t+1}}{1 - \beta^{t+1}} \\ \tilde{\mathbf{v}}_{t+1} &\leftarrow \frac{\mathbf{v}_{t+1}}{1 - \gamma^{t+1}}.\end{aligned}\tag{6.16}$$

$$\phi_{t+1} \leftarrow \phi_t - \alpha \cdot \frac{\tilde{\mathbf{m}}_{t+1}}{\sqrt{\tilde{\mathbf{v}}_{t+1}} + \epsilon}.\tag{6.17}$$

$$\begin{aligned}\mathbf{m}_{t+1} &\leftarrow \beta \cdot \mathbf{m}_t + (1 - \beta) \sum_{i \in \mathcal{B}_t} \frac{\partial \ell_i[\phi_t]}{\partial \phi} \\ \mathbf{v}_{t+1} &\leftarrow \gamma \cdot \mathbf{v}_t + (1 - \gamma) \left(\sum_{i \in \mathcal{B}_t} \frac{\partial \ell_i[\phi_t]}{\partial \phi} \right)^2,\end{aligned}\tag{6.18}$$

$$\mathbf{H}[\phi] = \begin{bmatrix} \frac{\partial^2 L}{\partial \phi_0^2} & \frac{\partial^2 L}{\partial \phi_0 \partial \phi_1} & \cdots & \frac{\partial^2 L}{\partial \phi_0 \partial \phi_N} \\ \frac{\partial^2 L}{\partial \phi_1 \partial \phi_0} & \frac{\partial^2 L}{\partial \phi_1^2} & \cdots & \frac{\partial^2 L}{\partial \phi_1 \partial \phi_N} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 L}{\partial \phi_N \partial \phi_0} & \frac{\partial^2 L}{\partial \phi_N \partial \phi_1} & \cdots & \frac{\partial^2 L}{\partial \phi_N^2} \end{bmatrix}.\tag{6.19}$$

$$\mathbf{H}[\phi] = \begin{bmatrix} \frac{\partial^2 L}{\partial \phi_0^2} & \frac{\partial^2 L}{\partial \phi_0 \partial \phi_1} \\ \frac{\partial^2 L}{\partial \phi_1 \partial \phi_0} & \frac{\partial^2 L}{\partial \phi_1^2} \end{bmatrix},\tag{6.20}$$

$$Pr(y = 1|x) = \text{sig}[\phi_0 + \phi_1 x],\tag{6.21}$$

$$\text{sig}[z] = \frac{1}{1 + \exp[-z]}.\tag{6.22}$$

$$f[x, \phi] = \phi_0 + \phi_1 a[\theta_{10} + \theta_{11}x] + \phi_2 a[\theta_{20} + \theta_{21}x] + \phi_3 a[\theta_{30} + \theta_{31}x].\tag{6.23}$$

Chapter 7

Gradients and initialization

$$\begin{aligned}\mathbf{h}_1 &= \mathbf{a}[\beta_0 + \mathbf{\Omega}_0 \mathbf{x}] \\ \mathbf{h}_2 &= \mathbf{a}[\beta_1 + \mathbf{\Omega}_1 \mathbf{h}_1] \\ \mathbf{h}_3 &= \mathbf{a}[\beta_2 + \mathbf{\Omega}_2 \mathbf{h}_2] \\ \mathbf{f}[\mathbf{x}, \phi] &= \beta_3 + \mathbf{\Omega}_3 \mathbf{h}_3,\end{aligned}\tag{7.1}$$

$$L[\phi] = \sum_{i=1}^I \ell_i.\tag{7.2}$$

$$\phi_{t+1} \leftarrow \phi_t - \alpha \sum_{i \in \mathcal{B}_t} \frac{\partial \ell_i[\phi_t]}{\partial \phi},\tag{7.3}$$

$$\frac{\partial \ell_i}{\partial \beta_k} \quad \text{and} \quad \frac{\partial \ell_i}{\partial \mathbf{\Omega}_k},\tag{7.4}$$

$$\mathbf{f}[x, \phi] = \beta_3 + \omega_3 \cdot \cos\left[\beta_2 + \omega_2 \cdot \exp\left[\beta_1 + \omega_1 \cdot \sin[\beta_0 + \omega_0 \cdot x]\right]\right],\tag{7.5}$$

$$\ell_i = (\mathbf{f}[x_i, \phi] - y_i)^2,\tag{7.6}$$

$$\frac{\partial \ell_i}{\partial \beta_0}, \quad \frac{\partial \ell_i}{\partial \omega_0}, \quad \frac{\partial \ell_i}{\partial \beta_1}, \quad \frac{\partial \ell_i}{\partial \omega_1}, \quad \frac{\partial \ell_i}{\partial \beta_2}, \quad \frac{\partial \ell_i}{\partial \omega_2}, \quad \frac{\partial \ell_i}{\partial \beta_3}, \quad \text{and} \quad \frac{\partial \ell_i}{\partial \omega_3}.\tag{7.7}$$

$$\begin{aligned}
\frac{\partial \ell_i}{\partial \omega_0} = & -2 \left(\beta_3 + \omega_3 \cdot \cos \left[\beta_2 + \omega_2 \cdot \exp \left[\beta_1 + \omega_1 \cdot \sin [\beta_0 + \omega_0 \cdot x_i] \right] - y_i \right) \right. \\
& \cdot \omega_1 \omega_2 \omega_3 \cdot x_i \cdot \cos [\beta_0 + \omega_0 \cdot x_i] \cdot \exp \left[\beta_1 + \omega_1 \cdot \sin [\beta_0 + \omega_0 \cdot x_i] \right] \\
& \left. \cdot \sin \left[\beta_2 + \omega_2 \cdot \exp \left[\beta_1 + \omega_1 \cdot \sin [\beta_0 + \omega_0 \cdot x_i] \right] \right] \right). \tag{7.8}
\end{aligned}$$

$$\begin{aligned}
f_0 &= \beta_0 + \omega_0 \cdot x_i \\
h_1 &= \sin[f_0] \\
f_1 &= \beta_1 + \omega_1 \cdot h_1 \\
h_2 &= \exp[f_1] \\
f_2 &= \beta_2 + \omega_2 \cdot h_2 \\
h_3 &= \cos[f_2] \\
f_3 &= \beta_3 + \omega_3 \cdot h_3 \\
\ell_i &= (f_3 - y_i)^2. \tag{7.9}
\end{aligned}$$

$$\frac{\partial \ell_i}{\partial f_3}, \quad \frac{\partial \ell_i}{\partial h_3}, \quad \frac{\partial \ell_i}{\partial f_2}, \quad \frac{\partial \ell_i}{\partial h_2}, \quad \frac{\partial \ell_i}{\partial f_1}, \quad \frac{\partial \ell_i}{\partial h_1}, \quad \text{and} \quad \frac{\partial \ell_i}{\partial f_0}. \tag{7.10}$$

$$\frac{\partial \ell_i}{\partial f_3} = 2(f_3 - y_i). \tag{7.11}$$

$$\frac{\partial \ell_i}{\partial h_3} = \frac{\partial f_3}{\partial h_3} \frac{\partial \ell_i}{\partial f_3}. \tag{7.12}$$

$$\begin{aligned}
\frac{\partial \ell_i}{\partial f_2} &= \frac{\partial h_3}{\partial f_2} \left(\frac{\partial f_3}{\partial h_3} \frac{\partial \ell_i}{\partial f_3} \right) \\
\frac{\partial \ell_i}{\partial h_2} &= \frac{\partial f_2}{\partial h_2} \left(\frac{\partial h_3}{\partial f_2} \frac{\partial f_3}{\partial h_3} \frac{\partial \ell_i}{\partial f_3} \right) \\
\frac{\partial \ell_i}{\partial f_1} &= \frac{\partial h_2}{\partial f_1} \left(\frac{\partial f_2}{\partial h_2} \frac{\partial h_3}{\partial f_2} \frac{\partial f_3}{\partial h_3} \frac{\partial \ell_i}{\partial f_3} \right) \\
\frac{\partial \ell_i}{\partial h_1} &= \frac{\partial f_1}{\partial h_1} \left(\frac{\partial h_2}{\partial f_1} \frac{\partial f_2}{\partial h_2} \frac{\partial h_3}{\partial f_2} \frac{\partial f_3}{\partial h_3} \frac{\partial \ell_i}{\partial f_3} \right) \\
\frac{\partial \ell_i}{\partial f_0} &= \frac{\partial h_1}{\partial f_0} \left(\frac{\partial f_1}{\partial h_1} \frac{\partial h_2}{\partial f_1} \frac{\partial f_2}{\partial h_2} \frac{\partial h_3}{\partial f_2} \frac{\partial f_3}{\partial h_3} \frac{\partial \ell_i}{\partial f_3} \right). \tag{7.13}
\end{aligned}$$

$$\begin{aligned}\frac{\partial \ell_i}{\partial \beta_k} &= \frac{\partial f_k}{\partial \beta_k} \frac{\partial \ell_i}{\partial f_k} \\ \frac{\partial \ell_i}{\partial \omega_k} &= \frac{\partial f_k}{\partial \omega_k} \frac{\partial \ell_i}{\partial f_k}.\end{aligned}\tag{7.14}$$

$$\frac{\partial f_k}{\partial \beta_k} = 1 \quad \text{and} \quad \frac{\partial f_k}{\partial \omega_k} = h_k.\tag{7.15}$$

$$\frac{\partial f_0}{\partial \beta_0} = 1 \quad \text{and} \quad \frac{\partial f_0}{\partial \omega_0} = x_i.\tag{7.16}$$

$$\begin{aligned}\mathbf{f}_0 &= \beta_0 + \mathbf{\Omega}_0 \mathbf{x}_i \\ \mathbf{h}_1 &= \mathbf{a}[\mathbf{f}_0] \\ \mathbf{f}_1 &= \beta_1 + \mathbf{\Omega}_1 \mathbf{h}_1 \\ \mathbf{h}_2 &= \mathbf{a}[\mathbf{f}_1] \\ \mathbf{f}_2 &= \beta_2 + \mathbf{\Omega}_2 \mathbf{h}_2 \\ \mathbf{h}_3 &= \mathbf{a}[\mathbf{f}_2] \\ \mathbf{f}_3 &= \beta_3 + \mathbf{\Omega}_3 \mathbf{h}_3 \\ \ell_i &= l[\mathbf{f}_3, y_i],\end{aligned}\tag{7.17}$$

$$\frac{\partial \ell_i}{\partial \mathbf{f}_2} = \frac{\partial \mathbf{h}_3}{\partial \mathbf{f}_2} \frac{\partial \mathbf{f}_3}{\partial \mathbf{h}_3} \frac{\partial \ell_i}{\partial \mathbf{f}_3}.\tag{7.18}$$

$$\frac{\partial \ell_i}{\partial \mathbf{f}_1} = \frac{\partial \mathbf{h}_2}{\partial \mathbf{f}_1} \frac{\partial \mathbf{f}_2}{\partial \mathbf{h}_2} \left(\frac{\partial \mathbf{h}_3}{\partial \mathbf{f}_2} \frac{\partial \mathbf{f}_3}{\partial \mathbf{h}_3} \frac{\partial \ell_i}{\partial \mathbf{f}_3} \right)\tag{7.19}$$

$$\frac{\partial \ell_i}{\partial \mathbf{f}_0} = \frac{\partial \mathbf{h}_1}{\partial \mathbf{f}_0} \frac{\partial \mathbf{f}_1}{\partial \mathbf{h}_1} \left(\frac{\partial \mathbf{h}_2}{\partial \mathbf{f}_1} \frac{\partial \mathbf{f}_2}{\partial \mathbf{h}_2} \frac{\partial \mathbf{h}_3}{\partial \mathbf{f}_2} \frac{\partial \mathbf{f}_3}{\partial \mathbf{h}_3} \frac{\partial \ell_i}{\partial \mathbf{f}_3} \right).\tag{7.20}$$

$$\frac{\partial \mathbf{f}_3}{\partial \mathbf{h}_3} = \frac{\partial}{\partial \mathbf{h}_3} (\beta_3 + \mathbf{\Omega}_3 \mathbf{h}_3) = \mathbf{\Omega}_3^T.\tag{7.21}$$

$$\begin{aligned}\frac{\partial \ell_i}{\partial \beta_k} &= \frac{\partial \mathbf{f}_k}{\partial \beta_k} \frac{\partial \ell_i}{\partial \mathbf{f}_k} \\ &= \frac{\partial}{\partial \beta_k} (\beta_k + \mathbf{\Omega}_k \mathbf{h}_k) \frac{\partial \ell_i}{\partial \mathbf{f}_k} \\ &= \frac{\partial \ell_i}{\partial \mathbf{f}_k},\end{aligned}\tag{7.22}$$

$$\begin{aligned}
\frac{\partial \ell_i}{\partial \boldsymbol{\Omega}_k} &= \frac{\partial \mathbf{f}_k}{\partial \boldsymbol{\Omega}_k} \frac{\partial \ell_i}{\partial \mathbf{f}_k} \\
&= \frac{\partial}{\partial \boldsymbol{\Omega}_k} (\boldsymbol{\beta}_k + \boldsymbol{\Omega}_k \mathbf{h}_k) \frac{\partial \ell_i}{\partial \mathbf{f}_k} \\
&= \frac{\partial \ell_i}{\partial \mathbf{f}_k} \mathbf{h}_k^T.
\end{aligned} \tag{7.23}$$

$$\begin{aligned}
\mathbf{f}_0 &= \boldsymbol{\beta}_0 + \boldsymbol{\Omega}_0 \mathbf{x}_i \\
\mathbf{h}_k &= \mathbf{a}[\mathbf{f}_{k-1}] & k \in \{1, 2, \dots, K\} \\
\mathbf{f}_k &= \boldsymbol{\beta}_k + \boldsymbol{\Omega}_k \mathbf{h}_k. & k \in \{1, 2, \dots, K\}
\end{aligned} \tag{7.24}$$

$$\begin{aligned}
\frac{\partial \ell_i}{\partial \boldsymbol{\beta}_k} &= \frac{\partial \ell_i}{\partial \mathbf{f}_k} & k \in \{K, K-1, \dots, 1\} \\
\frac{\partial \ell_i}{\partial \boldsymbol{\Omega}_k} &= \frac{\partial \ell_i}{\partial \mathbf{f}_k} \mathbf{h}_k^T & k \in \{K, K-1, \dots, 1\} \\
\frac{\partial \ell_i}{\partial \mathbf{f}_{k-1}} &= \mathbb{I}[\mathbf{f}_{k-1} > 0] \odot \left(\boldsymbol{\Omega}_k^T \frac{\partial \ell_i}{\partial \mathbf{f}_k} \right), & k \in \{K, K-1, \dots, 1\}
\end{aligned} \tag{7.25}$$

$$\begin{aligned}
\frac{\partial \ell_i}{\partial \boldsymbol{\beta}_0} &= \frac{\partial \ell_i}{\partial \mathbf{f}_0} \\
\frac{\partial \ell_i}{\partial \boldsymbol{\Omega}_0} &= \frac{\partial \ell_i}{\partial \mathbf{f}_0} \mathbf{x}_i^T.
\end{aligned} \tag{7.26}$$

$$\begin{aligned}
\mathbf{f}_k &= \boldsymbol{\beta}_k + \boldsymbol{\Omega}_k \mathbf{h}_k \\
&= \boldsymbol{\beta}_k + \boldsymbol{\Omega}_k \mathbf{a}[\mathbf{f}_{k-1}],
\end{aligned} \tag{7.27}$$

$$\begin{aligned}
\mathbf{h} &= \mathbf{a}[\mathbf{f}], \\
\mathbf{f}' &= \boldsymbol{\beta} + \boldsymbol{\Omega} \mathbf{h}
\end{aligned} \tag{7.28}$$

$$\begin{aligned}
\mathbb{E}[f'_i] &= \mathbb{E} \left[\beta_i + \sum_{j=1}^{D_h} \Omega_{ij} h_j \right] \\
&= \mathbb{E} [\beta_i] + \sum_{j=1}^{D_h} \mathbb{E} [\Omega_{ij} h_j] \\
&= \mathbb{E} [\beta_i] + \sum_{j=1}^{D_h} \mathbb{E} [\Omega_{ij}] \mathbb{E} [h_j] \\
&= 0 + \sum_{j=1}^{D_h} 0 \cdot \mathbb{E} [h_j] = 0,
\end{aligned} \tag{7.29}$$

$$\begin{aligned}
\sigma_{f'_i}^2 &= \mathbb{E}[f_i'^2] - \mathbb{E}[f'_i]^2 \\
&= \mathbb{E} \left[\left(\beta_i + \sum_{j=1}^{D_h} \Omega_{ij} h_j \right)^2 \right] - 0 \\
&= \mathbb{E} \left[\left(\sum_{j=1}^{D_h} \Omega_{ij} h_j \right)^2 \right] \\
&= \sum_{j=1}^{D_h} \mathbb{E} [\Omega_{ij}^2] \mathbb{E} [h_j^2] \\
&= \sum_{j=1}^{D_h} \sigma_\Omega^2 \mathbb{E} [h_j^2] = \sigma_\Omega^2 \sum_{j=1}^{D_h} \mathbb{E} [h_j^2],
\end{aligned} \tag{7.30}$$

$$\sigma_{f'_i}^2 = \sigma_\Omega^2 \sum_{j=1}^{D_h} \frac{\sigma_f^2}{2} = \frac{1}{2} D_h \sigma_\Omega^2 \sigma_f^2. \tag{7.31}$$

$$\sigma_\Omega^2 = \frac{2}{D_h}, \tag{7.32}$$

$$\sigma_\Omega^2 = \frac{2}{D_{h'}}, \tag{7.33}$$

$$\sigma_\Omega^2 = \frac{4}{D_h + D_{h'}}. \tag{7.34}$$

$$\begin{aligned}
y &= \phi_0 + \phi_1 \mathbf{a} \left[\psi_{01} + \psi_{11} \mathbf{a} [\theta_{01} + \theta_{11} x] + \psi_{21} \mathbf{a} [\theta_{02} + \theta_{12} x] \right] \\
&\quad + \phi_2 \mathbf{a} \left[\psi_{02} + \psi_{12} \mathbf{a} [\theta_{01} + \theta_{11} x] + \psi_{22} \mathbf{a} [\theta_{02} + \theta_{12} x] \right],
\end{aligned} \tag{7.35}$$

$$\ell_i = (y_i - \mathbf{f}[\mathbf{x}_i, \phi])^2. \tag{7.36}$$

$$\ell_i = -(1 - y_i) \log \left[1 - \text{sig}[\mathbf{f}[\mathbf{x}_i, \phi]] \right] - y_i \log \left[\text{sig}[\mathbf{f}[\mathbf{x}_i, \phi]] \right], \tag{7.37}$$

$$\text{sig}[z] = \frac{1}{1 + \exp[-z]}. \tag{7.38}$$

$$\frac{\partial \mathbf{z}}{\partial \mathbf{h}} = \mathbf{\Omega}^T, \tag{7.39}$$

$$\text{Heaviside}[z] = \begin{cases} 0 & z < 0 \\ 1 & z \geq 0 \end{cases}, \tag{7.40}$$

$$\text{rect}[z] = \begin{cases} 0 & z < 0 \\ 1 & 0 \leq z \leq 1 \\ 0 & z > 1 \end{cases}. \tag{7.41}$$

$$\frac{\partial \ell}{\partial \mathbf{\Omega}} = \frac{\partial \ell}{\partial \mathbf{f}} \mathbf{h}^T. \tag{7.42}$$

$$\mathbf{a}[z] = \text{ReLU}[z] = \begin{cases} \alpha \cdot z & z < 0 \\ z & z \geq 0 \end{cases}, \tag{7.43}$$

$$y = \exp \left[\exp[x] + \exp[x]^2 \right] + \sin[\exp[x] + \exp[x]^2]. \tag{7.44}$$

$$\begin{aligned}
f_1 &= \exp[x] \\
f_2 &= f_1^2 \\
f_3 &= f_1 + f_2 \\
f_4 &= \exp[f_3] \\
f_5 &= \sin[f_3] \\
y &= f_4 + f_5.
\end{aligned} \tag{7.45}$$

$$\frac{\partial y}{\partial f_5}, \frac{\partial y}{\partial f_4}, \frac{\partial y}{\partial f_3}, \frac{\partial y}{\partial f_2}, \frac{\partial y}{\partial f_1} \text{ and } \frac{\partial y}{\partial x}, \quad (7.46)$$

$$\frac{\partial f_1}{\partial x}, \frac{\partial f_2}{\partial x}, \frac{\partial f_3}{\partial x}, \frac{\partial f_4}{\partial x}, \frac{\partial f_5}{\partial x}, \text{ and } \frac{\partial y}{\partial x}, \quad (7.47)$$

$$b = \text{ReLU}[a] = \begin{cases} 0 & a < 0 \\ a & a \geq 0 \end{cases}, \quad (7.48)$$

Chapter 8

Measuring performance

$$\mu[x] = \mathbb{E}_y[y[x]] = \int y[x] Pr(y|x) dy, \quad (8.1)$$

$$\begin{aligned} L[x] &= (f[x, \phi] - y[x])^2 \\ &= \left((f[x, \phi] - \mu[x]) + (\mu[x] - y[x]) \right)^2 \\ &= (f[x, \phi] - \mu[x])^2 + 2(f[x, \phi] - \mu[x])(\mu[x] - y[x]) + (\mu[x] - y[x])^2, \end{aligned} \quad (8.2)$$

$$\begin{aligned} \mathbb{E}_y[L[x]] &= \mathbb{E}_y \left[(f[x, \phi] - \mu[x])^2 + 2(f[x, \phi] - \mu[x])(\mu[x] - y[x]) + (\mu[x] - y[x])^2 \right] \\ &= (f[x, \phi] - \mu[x])^2 + 2(f[x, \phi] - \mu[x])(\mu[x] - \mathbb{E}_y[y[x]]) + \mathbb{E}_y[(\mu[x] - y[x])^2] \\ &= (f[x, \phi] - \mu[x])^2 + 2(f[x, \phi] - \mu[x]) \cdot 0 + \mathbb{E}_y[(\mu[x] - y[x])^2] \\ &= (f[x, \phi] - \mu[x])^2 + \sigma^2, \end{aligned} \quad (8.3)$$

$$f_\mu[x] = \mathbb{E}_{\mathcal{D}}[f[x, \phi[\mathcal{D}]]]. \quad (8.4)$$

$$\begin{aligned} (f[x, \phi[\mathcal{D}]] - \mu[x])^2 &= \left((f[x, \phi[\mathcal{D}]] - f_\mu[x]) + (f_\mu[x] - \mu[x]) \right)^2 \\ &= (f[x, \phi[\mathcal{D}]] - f_\mu[x])^2 + 2(f[x, \phi[\mathcal{D}]] - f_\mu[x])(f_\mu[x] - \mu[x]) + (f_\mu[x] - \mu[x])^2. \end{aligned} \quad (8.5)$$

$$\mathbb{E}_{\mathcal{D}}[(f[x, \phi[\mathcal{D}]] - \mu[x])^2] = \mathbb{E}_{\mathcal{D}}[(f[x, \phi[\mathcal{D}]] - f_\mu[x])^2] + (f_\mu[x] - \mu[x])^2, \quad (8.6)$$

$$\mathbb{E}_{\mathcal{D}}\left[\mathbb{E}_y[L[x]]\right] = \mathbb{E}_{\mathcal{D}}\left[\underbrace{\left(\mathbf{f}[x, \phi[\mathcal{D}]] - \mathbf{f}_{\mu}[x]\right)^2}_{\text{variance}} + \underbrace{\left(\mathbf{f}_{\mu}[x] - \mu[x]\right)^2}_{\text{bias}} + \underbrace{\sigma^2}_{\text{noise}}\right]. \quad (8.7)$$

$$\text{Vol}[r] = \frac{r^D \pi^{D/2}}{\Gamma[D/2 + 1]}, \quad (8.8)$$

Chapter 9

Regularization

$$\begin{aligned}\hat{\phi} &= \operatorname{argmin}_{\phi} [L[\phi]] \\ &= \operatorname{argmin}_{\phi} \left[\sum_{i=1}^I \ell_i[\mathbf{x}_i, \mathbf{y}_i] \right],\end{aligned}\tag{9.1}$$

$$\hat{\phi} = \operatorname{argmin}_{\phi} \left[\sum_{i=1}^I \ell_i[\mathbf{x}_i, \mathbf{y}_i] + \lambda \cdot g[\phi] \right],\tag{9.2}$$

$$\hat{\phi} = \operatorname{argmax}_{\phi} \left[\prod_{i=1}^I \operatorname{Pr}(\mathbf{y}_i | \mathbf{x}_i, \phi) \right].\tag{9.3}$$

$$\hat{\phi} = \operatorname{argmax}_{\phi} \left[\prod_{i=1}^I \operatorname{Pr}(\mathbf{y}_i | \mathbf{x}_i, \phi) \operatorname{Pr}(\phi) \right].\tag{9.4}$$

$$\hat{\phi} = \operatorname{argmin}_{\phi} \left[\sum_{i=1}^I \ell_i[\mathbf{x}_i, \mathbf{y}_i] + \lambda \sum_j \phi_j^2 \right],\tag{9.5}$$

$$\frac{d\phi}{dt} = -\frac{\partial L}{\partial \phi}.\tag{9.6}$$

$$\phi_{t+1} = \phi_t - \alpha \frac{\partial L[\phi_t]}{\partial \phi},\tag{9.7}$$

$$\tilde{L}_{GD}[\phi] = L[\phi] + \frac{\alpha}{4} \left\| \frac{\partial L}{\partial \phi} \right\|^2. \quad (9.8)$$

$$\begin{aligned} \tilde{L}_{SGD}[\phi] &= \tilde{L}_{GD}[\phi] + \frac{\alpha}{4B} \sum_{b=1}^B \left\| \frac{\partial L_b}{\partial \phi} - \frac{\partial L}{\partial \phi} \right\|^2 \\ &= L[\phi] + \frac{\alpha}{4} \left\| \frac{\partial L}{\partial \phi} \right\|^2 + \frac{\alpha}{4B} \sum_{b=1}^B \left\| \frac{\partial L_b}{\partial \phi} - \frac{\partial L}{\partial \phi} \right\|^2. \end{aligned} \quad (9.9)$$

$$L = \frac{1}{I} \sum_{i=1}^I \ell_i[\mathbf{x}_i, y_i] \quad \text{and} \quad L_b = \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}_b} \ell_i[\mathbf{x}_i, y_i]. \quad (9.10)$$

$$Pr(\phi | \{\mathbf{x}_i, \mathbf{y}_i\}) = \frac{\prod_{i=1}^I Pr(\mathbf{y}_i | \mathbf{x}_i, \phi) Pr(\phi)}{\int \prod_{i=1}^I Pr(\mathbf{y}_i | \mathbf{x}_i, \phi) Pr(\phi) d\phi}, \quad (9.11)$$

$$Pr(\mathbf{y} | \mathbf{x}, \{\mathbf{x}_i, \mathbf{y}_i\}) = \int Pr(\mathbf{y} | \mathbf{x}, \phi) Pr(\phi | \{\mathbf{x}_i, \mathbf{y}_i\}) d\phi. \quad (9.12)$$

$$\phi \longleftarrow (1 - \lambda')\phi - \alpha \frac{\partial L}{\partial \phi}, \quad (9.13)$$

$$\phi_1 = \phi_0 + \alpha \cdot \mathbf{g}[\phi_0], \quad (9.14)$$

$$\frac{d\phi}{dt} = \mathbf{g}[\phi]. \quad (9.15)$$

$$\frac{d\phi}{dt} \approx \mathbf{g}[\phi] + \alpha \mathbf{g}_1[\phi] + \dots, \quad (9.16)$$

$$\begin{aligned} \phi[\alpha] &\approx \phi + \alpha \frac{d\phi}{dt} + \frac{\alpha^2}{2} \frac{d^2\phi}{dt^2} \Big|_{\phi=\phi_0} \\ &\approx \phi + \alpha (\mathbf{g}[\phi] + \alpha \mathbf{g}_1[\phi]) + \frac{\alpha^2}{2} \left(\frac{\partial \mathbf{g}[\phi]}{\partial \phi} \frac{d\phi}{dt} + \alpha \frac{\partial \mathbf{g}_1[\phi]}{\partial \phi} \frac{d\phi}{dt} \right) \Big|_{\phi=\phi_0} \\ &= \phi + \alpha (\mathbf{g}[\phi] + \alpha \mathbf{g}_1[\phi]) + \frac{\alpha^2}{2} \left(\frac{\partial \mathbf{g}[\phi]}{\partial \phi} \mathbf{g}[\phi] + \alpha \frac{\partial \mathbf{g}_1[\phi]}{\partial \phi} \mathbf{g}[\phi] \right) \Big|_{\phi=\phi_0} \\ &\approx \phi + \alpha \mathbf{g}[\phi] + \alpha^2 \left(\mathbf{g}_1[\phi] + \frac{1}{2} \frac{\partial \mathbf{g}[\phi]}{\partial \phi} \mathbf{g}[\phi] \right) \Big|_{\phi=\phi_0}, \end{aligned} \quad (9.17)$$

$$\mathbf{g}_1[\phi] = -\frac{1}{2} \frac{\partial \mathbf{g}[\phi]}{\partial \phi} \mathbf{g}[\phi]. \quad (9.18)$$

$$\begin{aligned} \frac{d\phi}{dt} &\approx \mathbf{g}[\phi] + \alpha \mathbf{g}_1[\phi] \\ &= -\frac{\partial L}{\partial \phi} - \frac{\alpha}{2} \left(\frac{\partial^2 L}{\partial \phi^2} \right) \frac{\partial L}{\partial \phi}. \end{aligned} \quad (9.19)$$

$$L_{GD}[\phi] = L[\phi] + \frac{\alpha}{4} \left\| \frac{\partial L}{\partial \phi} \right\|^2, \quad (9.20)$$

$$Pr(\phi) = \prod_{j=1}^J \text{Norm}_{\phi_j}[0, \sigma_\phi^2], \quad (9.21)$$

$$\phi \longleftarrow (1 - \lambda)\phi - \alpha \frac{\partial L}{\partial \phi}, \quad (9.22)$$

$$\tilde{L}[\phi] = L[\phi] + \frac{\lambda}{2\alpha} \sum_k \phi_k^2, \quad (9.23)$$

Chapter 10

Convolutional networks

$$\mathbf{f}[\mathbf{t}[\mathbf{x}]] = \mathbf{f}[\mathbf{x}]. \quad (10.1)$$

$$\mathbf{f}[\mathbf{t}[\mathbf{x}]] = \mathbf{t}[\mathbf{f}[\mathbf{x}]]. \quad (10.2)$$

$$z_i = \omega_1 x_{i-1} + \omega_2 x_i + \omega_3 x_{i+1}, \quad (10.3)$$

$$\begin{aligned} h_i &= \mathbf{a}[\beta + \omega_1 x_{i-1} + \omega_2 x_i + \omega_3 x_{i+1}] \\ &= \mathbf{a}\left[\beta + \sum_{j=1}^3 \omega_j x_{i+j-2}\right], \end{aligned} \quad (10.4)$$

$$h_i = \mathbf{a}\left[\beta_i + \sum_{j=1}^D \omega_{ij} x_j\right]. \quad (10.5)$$

$$h_{ij} = \mathbf{a}\left[\beta + \sum_{m=1}^3 \sum_{n=1}^3 \omega_{mn} x_{i+m-2, j+n-2}\right], \quad (10.6)$$

Chapter 11

Residual networks

$$\begin{aligned}\mathbf{h}_1 &= \mathbf{f}_1[\mathbf{x}, \phi_1] \\ \mathbf{h}_2 &= \mathbf{f}_2[\mathbf{h}_1, \phi_2] \\ \mathbf{h}_3 &= \mathbf{f}_3[\mathbf{h}_2, \phi_3] \\ \mathbf{y} &= \mathbf{f}_4[\mathbf{h}_3, \phi_4],\end{aligned}\tag{11.1}$$

$$\mathbf{y} = \mathbf{f}_4 \left[\mathbf{f}_3 \left[\mathbf{f}_2 \left[\mathbf{f}_1[\mathbf{x}, \phi_1], \phi_2 \right], \phi_3 \right], \phi_4 \right].\tag{11.2}$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{f}_1} = \frac{\partial \mathbf{f}_2}{\partial \mathbf{f}_1} \frac{\partial \mathbf{f}_3}{\partial \mathbf{f}_2} \frac{\partial \mathbf{f}_4}{\partial \mathbf{f}_3}.\tag{11.3}$$

$$\begin{aligned}\mathbf{h}_1 &= \mathbf{x} + \mathbf{f}_1[\mathbf{x}, \phi_1] \\ \mathbf{h}_2 &= \mathbf{h}_1 + \mathbf{f}_2[\mathbf{h}_1, \phi_2] \\ \mathbf{h}_3 &= \mathbf{h}_2 + \mathbf{f}_3[\mathbf{h}_2, \phi_3] \\ \mathbf{y} &= \mathbf{h}_3 + \mathbf{f}_4[\mathbf{h}_3, \phi_4],\end{aligned}\tag{11.4}$$

$$\begin{aligned}\mathbf{y} &= \mathbf{x} + \mathbf{f}_1[\mathbf{x}] \\ &+ \mathbf{f}_2[\mathbf{x} + \mathbf{f}_1[\mathbf{x}]] \\ &+ \mathbf{f}_3 \left[\mathbf{x} + \mathbf{f}_1[\mathbf{x}] + \mathbf{f}_2[\mathbf{x} + \mathbf{f}_1[\mathbf{x}]] \right] \\ &+ \mathbf{f}_4 \left[\mathbf{x} + \mathbf{f}_1[\mathbf{x}] + \mathbf{f}_2[\mathbf{x} + \mathbf{f}_1[\mathbf{x}]] + \mathbf{f}_3 \left[\mathbf{x} + \mathbf{f}_1[\mathbf{x}] + \mathbf{f}_2[\mathbf{x} + \mathbf{f}_1[\mathbf{x}]] \right] \right],\end{aligned}\tag{11.5}$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{f}_1} = \mathbf{I} + \frac{\partial \mathbf{f}_2}{\partial \mathbf{f}_1} + \left(\frac{\partial \mathbf{f}_3}{\partial \mathbf{f}_1} + \frac{\partial \mathbf{f}_2}{\partial \mathbf{f}_1} \frac{\partial \mathbf{f}_3}{\partial \mathbf{f}_2} \right) + \left(\frac{\partial \mathbf{f}_4}{\partial \mathbf{f}_1} + \frac{\partial \mathbf{f}_2}{\partial \mathbf{f}_1} \frac{\partial \mathbf{f}_4}{\partial \mathbf{f}_2} + \frac{\partial \mathbf{f}_3}{\partial \mathbf{f}_1} \frac{\partial \mathbf{f}_4}{\partial \mathbf{f}_3} + \frac{\partial \mathbf{f}_2}{\partial \mathbf{f}_1} \frac{\partial \mathbf{f}_3}{\partial \mathbf{f}_2} \frac{\partial \mathbf{f}_4}{\partial \mathbf{f}_3} \right),\tag{11.6}$$

$$\begin{aligned}
m_h &= \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} h_i \\
s_h &= \sqrt{\frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} (h_i - m_h)^2},
\end{aligned} \tag{11.7}$$

$$h_i \leftarrow \frac{h_i - m_h}{s_h + \epsilon} \quad \forall i \in \mathcal{B}, \tag{11.8}$$

$$h_i \leftarrow \gamma h_i + \delta \quad \forall i \in \mathcal{B}. \tag{11.9}$$

$$\begin{aligned}
f_1 &= \mathbb{E}[z_i] & f_5 &= \sqrt{f_4 + \epsilon} \\
f_{2i} &= z_i - f_1 & f_6 &= 1/f_5 \\
f_{3i} &= f_{2i}^2 & f_{7i} &= f_{2i} \times f_6 \\
f_4 &= \mathbb{E}[f_{3i}] & z'_i &= f_{7i} \times \gamma + \delta,
\end{aligned} \tag{11.10}$$

Chapter 12

Transformers

$$\mathbf{f}[\mathbf{x}] = \mathbf{ReLU}[\boldsymbol{\beta} + \boldsymbol{\Omega}\mathbf{x}], \quad (12.1)$$

$$\mathbf{v}_m = \boldsymbol{\beta}_v + \boldsymbol{\Omega}_v \mathbf{x}_m, \quad (12.2)$$

$$\mathbf{sa}_n[\mathbf{x}_1, \dots, \mathbf{x}_N] = \sum_{m=1}^N a[\mathbf{x}_m, \mathbf{x}_n] \mathbf{v}_m. \quad (12.3)$$

$$\begin{aligned} \mathbf{q}_n &= \boldsymbol{\beta}_q + \boldsymbol{\Omega}_q \mathbf{x}_n \\ \mathbf{k}_m &= \boldsymbol{\beta}_k + \boldsymbol{\Omega}_k \mathbf{x}_m, \end{aligned} \quad (12.4)$$

$$\begin{aligned} a[\mathbf{x}_m, \mathbf{x}_n] &= \text{softmax}_m [\mathbf{k}_\bullet^T \mathbf{q}_n] \\ &= \frac{\exp [\mathbf{k}_m^T \mathbf{q}_n]}{\sum_{m'=1}^N \exp [\mathbf{k}_{m'}^T \mathbf{q}_n]}, \end{aligned} \quad (12.5)$$

$$\begin{aligned} \mathbf{V}[\mathbf{X}] &= \boldsymbol{\beta}_v \mathbf{1}^T + \boldsymbol{\Omega}_v \mathbf{X} \\ \mathbf{Q}[\mathbf{X}] &= \boldsymbol{\beta}_q \mathbf{1}^T + \boldsymbol{\Omega}_q \mathbf{X} \\ \mathbf{K}[\mathbf{X}] &= \boldsymbol{\beta}_k \mathbf{1}^T + \boldsymbol{\Omega}_k \mathbf{X}, \end{aligned} \quad (12.6)$$

$$\mathbf{Sa}[\mathbf{X}] = \mathbf{V}[\mathbf{X}] \cdot \mathbf{Softmax} \left[\mathbf{K}[\mathbf{X}]^T \mathbf{Q}[\mathbf{X}] \right], \quad (12.7)$$

$$\mathbf{Sa}[\mathbf{X}] = \mathbf{V} \cdot \mathbf{Softmax} \left[\mathbf{K}^T \mathbf{Q} \right]. \quad (12.8)$$

$$\mathbf{Sa}[\mathbf{X}] = \mathbf{V} \cdot \mathbf{Softmax} \left[\frac{\mathbf{K}^T \mathbf{Q}}{\sqrt{D_q}} \right]. \quad (12.9)$$

$$\begin{aligned} \mathbf{V}_h &= \beta_{vh} \mathbf{1}^T + \boldsymbol{\Omega}_{vh} \mathbf{X} \\ \mathbf{Q}_h &= \beta_{qh} \mathbf{1}^T + \boldsymbol{\Omega}_{qh} \mathbf{X} \\ \mathbf{K}_h &= \beta_{kh} \mathbf{1}^T + \boldsymbol{\Omega}_{kh} \mathbf{X}. \end{aligned} \quad (12.10)$$

$$\mathbf{Sa}_h[\mathbf{X}] = \mathbf{V}_h \cdot \mathbf{Softmax} \left[\frac{\mathbf{K}_h^T \mathbf{Q}_h}{\sqrt{D_q}} \right], \quad (12.11)$$

$$\mathbf{MhSa}[\mathbf{X}] = \boldsymbol{\Omega}_c \left[\mathbf{Sa}_1[\mathbf{X}]^T, \mathbf{Sa}_2[\mathbf{X}]^T, \dots, \mathbf{Sa}_H[\mathbf{X}]^T \right]^T. \quad (12.12)$$

$$\begin{aligned} \mathbf{X} &\leftarrow \mathbf{X} + \mathbf{MhSa}[\mathbf{X}] \\ \mathbf{X} &\leftarrow \mathbf{LayerNorm}[\mathbf{X}] \\ \mathbf{x}_n &\leftarrow \mathbf{x}_n + \mathbf{mlp}[\mathbf{x}_n] & \forall n \in \{1, \dots, N\} \\ \mathbf{X} &\leftarrow \mathbf{LayerNorm}[\mathbf{X}], \end{aligned} \quad (12.13)$$

$$\begin{aligned} Pr(\text{It takes great courage to let yourself appear weak}) &= \\ Pr(\text{It}) \times Pr(\text{takes}|\text{It}) \times Pr(\text{great}|\text{It takes}) \times Pr(\text{courage}|\text{It takes great}) \times \\ Pr(\text{to}|\text{It takes great courage}) \times Pr(\text{let}|\text{It takes great courage to}) \times \\ Pr(\text{yourself}|\text{It takes great courage to let}) \times \\ Pr(\text{appear}|\text{It takes great courage to let yourself}) \times \\ Pr(\text{weak}|\text{It takes great courage to let yourself appear}). \end{aligned} \quad (12.14)$$

$$Pr(t_1, t_2, \dots, t_N) = Pr(t_1) \prod_{n=2}^N Pr(t_n | t_1, \dots, t_{n-1}). \quad (12.15)$$

$$\mathbf{Sa}[\mathbf{X}] = \mathbf{V} \cdot \mathbf{Softmax} \left[\frac{\mathbf{K}^T \mathbf{Q}}{\sqrt{D_q}} \right], \quad (12.16)$$

$$\begin{aligned}
\mathbf{V} &= \beta_v \mathbf{1}^T + \boldsymbol{\Omega}_v \mathbf{X} \\
\mathbf{Q} &= \beta_q \mathbf{1}^T + \boldsymbol{\Omega}_q (\mathbf{X} + \boldsymbol{\Pi}) \\
\mathbf{K} &= \beta_k \mathbf{1}^T + \boldsymbol{\Omega}_k (\mathbf{X} + \boldsymbol{\Pi}).
\end{aligned} \tag{12.17}$$

$$\mathbf{Sa}[\mathbf{XP}] = \mathbf{Sa}[\mathbf{X}]\mathbf{P}. \tag{12.18}$$

$$y_i = \text{softmax}_i[\mathbf{z}] = \frac{\exp[z_i]}{\sum_{j=1}^5 \exp[z_j]}, \tag{12.19}$$

$$a[\mathbf{x}_m, \mathbf{x}_n] = \text{softmax}_m [\mathbf{k}_{\bullet}^T \mathbf{q}_n] = \frac{\exp [\mathbf{k}_m^T \mathbf{q}_n]}{\sum_{m'=1}^N \exp [\mathbf{k}_{m'}^T \mathbf{q}_n]}. \tag{12.20}$$

Chapter 13

Graph neural networks

$$\begin{aligned}\mathbf{X}' &= \mathbf{X}\mathbf{P} \\ \mathbf{A}' &= \mathbf{P}^T \mathbf{A} \mathbf{P},\end{aligned}\tag{13.1}$$

$$Pr(y = 1 | \mathbf{X}, \mathbf{A}) = \text{sig} [\beta_K + \boldsymbol{\omega}_K \mathbf{H}_K \mathbf{1} / N],\tag{13.2}$$

$$Pr(y^{(n)} = 1 | \mathbf{X}, \mathbf{A}) = \text{sig} \left[\beta_K + \boldsymbol{\omega}_K \mathbf{h}_K^{(n)} \right].\tag{13.3}$$

$$Pr(y^{(mn)} = 1 | \mathbf{X}, \mathbf{A}) = \text{sig} \left[\mathbf{h}^{(m)T} \mathbf{h}^{(n)} \right].\tag{13.4}$$

$$\begin{aligned}\mathbf{H}_1 &= \mathbf{F}[\mathbf{X}, \mathbf{A}, \phi_0] \\ \mathbf{H}_2 &= \mathbf{F}[\mathbf{H}_1, \mathbf{A}, \phi_1] \\ \mathbf{H}_3 &= \mathbf{F}[\mathbf{H}_2, \mathbf{A}, \phi_2] \\ \vdots &= \vdots \\ \mathbf{H}_K &= \mathbf{F}[\mathbf{H}_{K-1}, \mathbf{A}, \phi_{K-1}],\end{aligned}\tag{13.5}$$

$$\mathbf{H}_{k+1} \mathbf{P} = \mathbf{F}[\mathbf{H}_k \mathbf{P}, \mathbf{P}^T \mathbf{A} \mathbf{P}, \phi_k].\tag{13.6}$$

$$y = \text{sig} [\beta_K + \boldsymbol{\omega}_K \mathbf{H}_K \mathbf{1} / N] = \text{sig} [\beta_K + \boldsymbol{\omega}_K \mathbf{H}_K \mathbf{P} \mathbf{1} / N],\tag{13.7}$$

$$\mathbf{agg}[n, k] = \sum_{m \in \text{ne}[n]} \mathbf{h}_k^{(m)}, \quad (13.8)$$

$$\mathbf{h}_{k+1}^{(n)} = \mathbf{a} \left[\beta_k + \Omega_k \cdot \mathbf{h}_k^{(n)} + \Omega_k \cdot \mathbf{agg}[n, k] \right]. \quad (13.9)$$

$$\begin{aligned} \mathbf{H}_{k+1} &= \mathbf{a} \left[\beta_k \mathbf{1}^T + \Omega_k \mathbf{H}_k + \Omega_k \mathbf{H}_k \mathbf{A} \right] \\ &= \mathbf{a} \left[\beta_k \mathbf{1}^T + \Omega_k \mathbf{H}_k (\mathbf{A} + \mathbf{I}) \right], \end{aligned} \quad (13.10)$$

$$\begin{aligned} \mathbf{H}_1 &= \mathbf{a} \left[\beta_0 \mathbf{1}^T + \Omega_0 \mathbf{X} (\mathbf{A} + \mathbf{I}) \right] \\ \mathbf{H}_2 &= \mathbf{a} \left[\beta_1 \mathbf{1}^T + \Omega_1 \mathbf{H}_1 (\mathbf{A} + \mathbf{I}) \right] \\ \vdots &= \vdots \\ \mathbf{H}_K &= \mathbf{a} \left[\beta_{K-1} \mathbf{1}^T + \Omega_{K-1} \mathbf{H}_{K-1} (\mathbf{A} + \mathbf{I}) \right] \\ \mathbf{f}[\mathbf{X}, \mathbf{A}, \Phi] &= \text{sig} \left[\beta_K + \omega_K \mathbf{H}_K \mathbf{1} / N \right], \end{aligned} \quad (13.11)$$

$$\mathbf{f}[\mathbf{X}, \mathbf{A}, \Phi] = \text{sig} \left[\beta_K \mathbf{1}^T + \omega_K \mathbf{H}_K \right], \quad (13.12)$$

$$\mathbf{H}_{k+1} = \mathbf{a} \left[\beta_k \mathbf{1}^T + \Omega_k \mathbf{H}_k (\mathbf{A} + \mathbf{I}) \right]. \quad (13.13)$$

$$\mathbf{H}_{k+1} = \mathbf{a} \left[\beta_k \mathbf{1}^T + \Omega_k \mathbf{H}_k (\mathbf{A} + (1 + \epsilon_k) \mathbf{I}) \right]. \quad (13.14)$$

$$\begin{aligned} \mathbf{H}_{k+1} &= \mathbf{a} \left[\beta_k \mathbf{1}^T + \Omega_k \mathbf{H}_k \mathbf{A} + \Psi_k \mathbf{H}_k \right] \\ &= \mathbf{a} \left[\beta_k \mathbf{1}^T + \begin{bmatrix} \Omega_k & \Psi_k \end{bmatrix} \begin{bmatrix} \mathbf{H}_k \mathbf{A} \\ \mathbf{H}_k \end{bmatrix} \right] \\ &= \mathbf{a} \left[\beta_k \mathbf{1}^T + \Omega'_k \begin{bmatrix} \mathbf{H}_k \mathbf{A} \\ \mathbf{H}_k \end{bmatrix} \right], \end{aligned} \quad (13.15)$$

$$\mathbf{H}_{k+1} = \left[\mathbf{a} \left[\beta_k \mathbf{1}^T + \Omega_k \mathbf{H}_k \mathbf{A} \right] \right]_{\mathbf{H}_k}. \quad (13.16)$$

$$\mathbf{agg}[n] = \frac{1}{|\mathbf{ne}[n]|} \sum_{m \in \mathbf{ne}[n]} \mathbf{h}_m, \quad (13.17)$$

$$\mathbf{H}_{k+1} = \mathbf{a} \left[\beta_k \mathbf{1}^T + \Omega_k \mathbf{H}_k (\mathbf{A} \mathbf{D}^{-1} + \mathbf{I}) \right]. \quad (13.18)$$

$$\mathbf{agg}[n] = \sum_{m \in \mathbf{ne}[n]} \frac{\mathbf{h}_m}{\sqrt{|\mathbf{ne}[n]| |\mathbf{ne}[m]|}}, \quad (13.19)$$

$$\mathbf{H}_{k+1} = \mathbf{a} \left[\beta_k \mathbf{1}^T + \Omega_k \mathbf{H}_k (\mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2} + \mathbf{I}) \right]. \quad (13.20)$$

$$\mathbf{agg}[n] = \max_{m \in \mathbf{ne}[n]} [\mathbf{h}_m], \quad (13.21)$$

$$\mathbf{H}'_k = \beta_k \mathbf{1}^T + \Omega_k \mathbf{H}_k. \quad (13.22)$$

$$s_{mn} = \mathbf{a} \left[\phi_k^T \begin{bmatrix} \mathbf{h}'_m \\ \mathbf{h}'_n \end{bmatrix} \right]. \quad (13.23)$$

$$\mathbf{H}_{k+1} = \mathbf{a} \left[\mathbf{H}'_k \cdot \mathbf{Softmask}[\mathbf{S}, \mathbf{A} + \mathbf{I}] \right], \quad (13.24)$$

$$\mathbf{h}_n \leftarrow \mathbf{f} \left[\mathbf{x}_n, \mathbf{x}_{m \in \mathbf{ne}[n]}, \mathbf{e}_{e \in \mathbf{nee}[n]}, \mathbf{h}_{m \in \mathbf{ne}[n]}, \phi \right], \quad (13.25)$$

$$\mathbf{h}_{k+1}^{(n)} = \mathbf{mlp} \left[(1 + \epsilon_k) \mathbf{h}_k^{(n)} + \sum_{m \in \mathbf{ne}[n]} \mathbf{h}_k^{(m)} \right]. \quad (13.26)$$

$$\mathbf{A}_1 = \begin{bmatrix} 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 1 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 & 0 & 0 \end{bmatrix} \quad \text{and} \quad \mathbf{A}_2 = \begin{bmatrix} 0 & 0 & 1 & 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 1 & 1 \\ 0 & 1 & 0 & 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 & 1 & 1 & 0 \end{bmatrix}.$$

$$\text{sig}[\beta_K + \boldsymbol{\omega}_K \mathbf{H}_K \mathbf{1}] = \text{sig}[\beta_K + \boldsymbol{\omega}_K \mathbf{H}_K \mathbf{P} \mathbf{1}], \quad (13.27)$$

$$\begin{aligned} \mathbf{H}_{k+1} &= \text{GraphLayer}[\mathbf{H}_k, \mathbf{A}] \\ &= \mathbf{a} \left[\beta_k \mathbf{1}^T + \boldsymbol{\Omega}_k \begin{bmatrix} \mathbf{H}_k \\ \mathbf{H}_k \mathbf{A} \end{bmatrix} \right], \end{aligned} \quad (13.28)$$

$$\text{GraphLayer}[\mathbf{H}_k, \mathbf{A}] \mathbf{P} = \text{GraphLayer}[\mathbf{H}_k \mathbf{P}, \mathbf{P}^T \mathbf{A} \mathbf{P}], \quad (13.29)$$

$$\mathbf{agg}[n] = \frac{1}{1 + |\text{ne}[n]|} \left(\mathbf{h}_n + \sum_{m \in \text{ne}[n]} \mathbf{h}_m \right). \quad (13.30)$$

Chapter 14

Unsupervised learning

$$L[\phi] = - \sum_{i=1}^I \log \left[Pr(\mathbf{x}_i | \phi) \right]. \quad (14.1)$$

$$IS = \exp \left[\frac{1}{I} \sum_{i=1}^I D_{KL} \left[Pr(y | \mathbf{x}_i^*) || Pr(y) \right] \right], \quad (14.2)$$

$$Pr(y) = \frac{1}{I} \sum_{i=1}^I Pr(y | \mathbf{x}_i^*). \quad (14.3)$$

Chapter 15

Generative adversarial networks

$$x_j^* = g[z_j, \theta] = z_j + \theta, \quad (15.1)$$

$$\hat{\phi} = \operatorname{argmin}_{\phi} \left[\sum_i -(1 - y_i) \log [1 - \operatorname{sig}[f[\mathbf{x}_i, \phi]]] - y_i \log [\operatorname{sig}[f[\mathbf{x}_i, \phi]]] \right], \quad (15.2)$$

$$\hat{\phi} = \operatorname{argmin}_{\phi} \left[\sum_j -\log [1 - \operatorname{sig}[f[\mathbf{x}_j^*, \phi]]] - \sum_i \log [\operatorname{sig}[f[\mathbf{x}_i, \phi]]] \right], \quad (15.3)$$

$$\hat{\theta} = \operatorname{argmax}_{\theta} \left[\min_{\phi} \left[\sum_j -\log [1 - \operatorname{sig}[f[\mathbf{g}[\mathbf{z}_j, \theta], \phi]]] - \sum_i \log [\operatorname{sig}[f[\mathbf{x}_i, \phi]]] \right] \right]. \quad (15.4)$$

$$\begin{aligned} L[\phi] &= \sum_j -\log [1 - \operatorname{sig}[f[\mathbf{g}[\mathbf{z}_j, \theta], \phi]]] - \sum_i \log [\operatorname{sig}[f[\mathbf{x}_i, \phi]]] \\ L[\theta] &= \sum_j \log [1 - \operatorname{sig}[f[\mathbf{g}[\mathbf{z}_j, \theta], \phi]]], \end{aligned} \quad (15.5)$$

$$\begin{aligned} L[\phi] &= -\frac{1}{J} \sum_{j=1}^J \left(\log [1 - \operatorname{sig}[f[\mathbf{x}_j^*, \phi]]] \right) - \frac{1}{I} \sum_{i=1}^I \left(\log [\operatorname{sig}[f[\mathbf{x}_i, \phi]]] \right) \\ &\approx -\mathbb{E}_{\mathbf{x}^*} \left[\log [1 - \operatorname{sig}[f[\mathbf{x}^*, \phi]]] \right] - \mathbb{E}_{\mathbf{x}} \left[\log [\operatorname{sig}[f[\mathbf{x}, \phi]]] \right] \\ &= -\int Pr(\mathbf{x}^*) \log [1 - \operatorname{sig}[f[\mathbf{x}^*, \phi]]] d\mathbf{x}^* - \int Pr(\mathbf{x}) \log [\operatorname{sig}[f[\mathbf{x}, \phi]]] d\mathbf{x}, \end{aligned} \quad (15.6)$$

$$Pr(\text{real}|\tilde{\mathbf{x}}) = \text{sig}[f[\tilde{\mathbf{x}}, \phi]] = \frac{Pr(\tilde{\mathbf{x}}|\text{real})}{Pr(\tilde{\mathbf{x}}|\text{generated}) + Pr(\tilde{\mathbf{x}}|\text{real})} = \frac{Pr(\mathbf{x})}{Pr(\mathbf{x}^*) + Pr(\mathbf{x})}, \quad (15.7)$$

$$\begin{aligned} L[\phi] &= - \int Pr(\mathbf{x}^*) \log[1 - \text{sig}[f[\mathbf{x}^*, \phi]]] d\mathbf{x}^* - \int Pr(\mathbf{x}) \log[\text{sig}[f[\mathbf{x}, \phi]]] d\mathbf{x} \\ &= - \int Pr(\mathbf{x}^*) \log\left[1 - \frac{Pr(\mathbf{x})}{Pr(\mathbf{x}^*) + Pr(\mathbf{x})}\right] d\mathbf{x}^* - \int Pr(\mathbf{x}) \log\left[\frac{Pr(\mathbf{x})}{Pr(\mathbf{x}^*) + Pr(\mathbf{x})}\right] d\mathbf{x} \\ &= - \int Pr(\mathbf{x}^*) \log\left[\frac{Pr(\mathbf{x}^*)}{Pr(\mathbf{x}^*) + Pr(\mathbf{x})}\right] d\mathbf{x}^* - \int Pr(\mathbf{x}) \log\left[\frac{Pr(\mathbf{x})}{Pr(\mathbf{x}^*) + Pr(\mathbf{x})}\right] d\mathbf{x}. \end{aligned} \quad (15.8)$$

$$\begin{aligned} D_{JS}[Pr(\mathbf{x}^*) || Pr(\mathbf{x})] & \\ &= \frac{1}{2} D_{KL}\left[Pr(\mathbf{x}^*) \left\| \frac{Pr(\mathbf{x}^*) + Pr(\mathbf{x})}{2} \right\| \right] + \frac{1}{2} D_{KL}\left[Pr(\mathbf{x}) \left\| \frac{Pr(\mathbf{x}^*) + Pr(\mathbf{x})}{2} \right\| \right] \\ &= \frac{1}{2} \int \underbrace{Pr(\mathbf{x}^*) \log\left[\frac{2Pr(\mathbf{x}^*)}{Pr(\mathbf{x}^*) + Pr(\mathbf{x})}\right] d\mathbf{x}^*}_{\text{quality}} + \frac{1}{2} \int \underbrace{Pr(\mathbf{x}) \log\left[\frac{2Pr(\mathbf{x})}{Pr(\mathbf{x}^*) + Pr(\mathbf{x})}\right] d\mathbf{x}}_{\text{coverage}}. \end{aligned} \quad (15.9)$$

$$D_w[Pr(x)||q(x)] = \min_{\mathbf{P}} \left[\sum_{i,j} P_{ij} \cdot |i - j| \right], \quad (15.10)$$

$$\begin{aligned} \sum_j P_{ij} &= Pr(x = i) && \text{initial distribution of } Pr(x) \\ \sum_i P_{ij} &= q(x = j) && \text{initial distribution of } q(x) \\ P_{ij} &\geq 0 && \text{non-negative masses.} \end{aligned} \quad (15.11)$$

$$D_w[Pr(x)||q(x)] = \max_{\mathbf{f}} \left[\sum_i Pr(x = i) f_i - \sum_j q(x = j) f_j \right], \quad (15.12)$$

$$|f_{i+1} - f_i| < 1. \quad (15.13)$$

$$D_w[Pr(\mathbf{x}), q(\mathbf{x})] = \min_{\pi[\bullet, \bullet]} \left[\iint \pi(\mathbf{x}_1, \mathbf{x}_2) \cdot \|\mathbf{x}_1 - \mathbf{x}_2\| d\mathbf{x}_1 d\mathbf{x}_2 \right], \quad (15.14)$$

$$D_w[Pr(\mathbf{x}), q(\mathbf{x})] = \max_{f[\mathbf{x}]} \left[\int Pr(\mathbf{x}) f[\mathbf{x}] d\mathbf{x} - \int q(\mathbf{x}) f[\mathbf{x}] d\mathbf{x} \right], \quad (15.15)$$

$$\begin{aligned}
L[\phi] &= \sum_j f[\mathbf{x}_j^*, \phi] - \sum_i f[\mathbf{x}_i, \phi] \\
&= \sum_j f[\mathbf{g}[\mathbf{z}_j, \boldsymbol{\theta}], \phi] - \sum_i f[\mathbf{x}_i, \phi],
\end{aligned} \tag{15.16}$$

$$\left| \frac{\partial f[\mathbf{x}, \phi]}{\partial \mathbf{x}} \right| < 1. \tag{15.17}$$

$$\mathbf{b} = [Pr(x=1), Pr(x=2), Pr(x=3), Pr(x=4), q(x=1), q(x=2), q(x=3), q(x=4)]^T \tag{15.18}$$

$$Pr(z) = \begin{cases} 0 & z < 0 \\ 1 & 0 \leq z \leq 1, \\ 0 & z > 1 \end{cases}, \quad \text{and} \quad Pr(z) = \begin{cases} 0 & z < a \\ 1 & a \leq z \leq a+1, \\ 0 & z > a \end{cases}. \tag{15.19}$$

$$D_{kl} = \log \left[\frac{\sigma_2}{\sigma_1} \right] + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}, \tag{15.20}$$

$$D_w = (\mu_1 - \mu_2)^2 + \sigma_1 + \sigma_2 - 2\sqrt{\sigma_1\sigma_2}, \tag{15.21}$$

Chapter 16

Normalizing flows

$$Pr(x|\phi) = \left| \frac{\partial \mathbf{f}[z, \phi]}{\partial z} \right|^{-1} \cdot Pr(z), \quad (16.1)$$

$$\begin{aligned} \hat{\phi} &= \operatorname{argmax}_{\phi} \left[\prod_{i=1}^I Pr(x_i|\phi) \right] \\ &= \operatorname{argmin}_{\phi} \left[\sum_{i=1}^I -\log [Pr(x_i|\phi)] \right] \\ &= \operatorname{argmin}_{\phi} \left[\sum_{i=1}^I \log \left[\left| \frac{\partial \mathbf{f}[z_i, \phi]}{\partial z_i} \right| \right] - \log [Pr(z_i)] \right], \end{aligned} \quad (16.2)$$

$$Pr(\mathbf{x}|\phi) = \left| \frac{\partial \mathbf{f}[\mathbf{z}, \phi]}{\partial \mathbf{z}} \right|^{-1} \cdot Pr(\mathbf{z}), \quad (16.3)$$

$$\mathbf{x} = \mathbf{f}[\mathbf{z}, \phi] = \mathbf{f}_K \left[\mathbf{f}_{K-1} \left[\dots \mathbf{f}_2 [\mathbf{f}_1 [\mathbf{z}, \phi_1], \phi_2], \dots \phi_{K-1} \right], \phi_K \right]. \quad (16.4)$$

$$\mathbf{z} = \mathbf{f}^{-1}[\mathbf{x}, \phi] = \mathbf{f}_1^{-1} \left[\mathbf{f}_2^{-1} \left[\dots \mathbf{f}_{K-1}^{-1} [\mathbf{f}_K^{-1}[\mathbf{x}, \phi_K], \phi_{K-1}], \dots \phi_2 \right], \phi_1 \right]. \quad (16.5)$$

$$\frac{\partial \mathbf{f}[\mathbf{z}, \phi]}{\partial \mathbf{z}} = \frac{\partial \mathbf{f}_K[\mathbf{f}_{K-1}, \phi_K]}{\partial \mathbf{f}_{K-1}} \cdot \frac{\partial \mathbf{f}_{K-1}[\mathbf{f}_{K-2}, \phi_{K-1}]}{\partial \mathbf{f}_{K-2}} \dots \frac{\partial \mathbf{f}_2[\mathbf{f}_1, \phi_2]}{\partial \mathbf{f}_1} \cdot \frac{\partial \mathbf{f}_1[\mathbf{z}, \phi_1]}{\partial \mathbf{z}}, \quad (16.6)$$

$$\left| \frac{\partial \mathbf{f}[\mathbf{z}, \phi]}{\partial \mathbf{z}} \right| = \left| \frac{\partial \mathbf{f}_K[\mathbf{f}_{K-1}, \phi_K]}{\partial \mathbf{f}_{K-1}} \right| \cdot \left| \frac{\partial \mathbf{f}_{K-1}[\mathbf{f}_{K-2}, \phi_{K-1}]}{\partial \mathbf{f}_{K-2}} \right| \dots \left| \frac{\partial \mathbf{f}_2[\mathbf{f}_1, \phi_2]}{\partial \mathbf{f}_1} \right| \cdot \left| \frac{\partial \mathbf{f}_1[\mathbf{z}, \phi_1]}{\partial \mathbf{z}} \right|. \quad (16.7)$$

$$\begin{aligned}
\hat{\phi} &= \operatorname{argmax}_{\phi} \left[\prod_{i=1}^I Pr(\mathbf{z}_i) \cdot \left| \frac{\partial \mathbf{f}[\mathbf{z}_i, \phi]}{\partial \mathbf{z}_i} \right|^{-1} \right] \\
&= \operatorname{argmin}_{\phi} \left[\sum_{i=1}^I \log \left[\left| \frac{\partial \mathbf{f}[\mathbf{z}_i, \phi]}{\partial \mathbf{z}_i} \right| \right] - \log[Pr(\mathbf{z}_i)] \right],
\end{aligned} \tag{16.8}$$

$$\mathbf{\Omega} = \mathbf{PL}(\mathbf{U} + \mathbf{D}), \tag{16.9}$$

$$\mathbf{f}[\mathbf{h}] = \left[\mathbf{f}[h_1, \phi], \mathbf{f}[h_2, \phi], \dots, \mathbf{f}[h_D, \phi] \right]^T. \tag{16.10}$$

$$\left| \frac{\partial \mathbf{f}[\mathbf{h}]}{\partial \mathbf{h}} \right| = \prod_{d=1}^D \left| \frac{\partial \mathbf{f}[h_d]}{\partial h_d} \right|. \tag{16.11}$$

$$\mathbf{f}[h, \phi] = \left(\sum_{k=1}^{b-1} \phi_k \right) + (hK - b + 1)\phi_b, \tag{16.12}$$

$$\begin{aligned}
\mathbf{h}'_1 &= \mathbf{h}_1 \\
\mathbf{h}'_2 &= \mathbf{g}[\mathbf{h}_2, \phi[\mathbf{h}_1]].
\end{aligned} \tag{16.13}$$

$$\begin{aligned}
\mathbf{h}_1 &= \mathbf{h}'_1 \\
\mathbf{h}_2 &= \mathbf{g}^{-1}[\mathbf{h}'_2, \phi[\mathbf{h}_1]].
\end{aligned} \tag{16.14}$$

$$h'_d = \mathbf{g}[h_d, \phi[\mathbf{h}_{1:d-1}]]. \tag{16.15}$$

$$\begin{aligned}
h'_1 &= \mathbf{g}[h_1, \phi] \\
h'_2 &= \mathbf{g}[h_2, \phi[h_1]] \\
h'_3 &= \mathbf{g}[h_3, \phi[h_{1:2}]] \\
h'_4 &= \mathbf{g}[h_4, \phi[h_{1:3}]].
\end{aligned} \tag{16.16}$$

$$\begin{aligned}
h_1 &= g^{-1}[h'_1, \phi] \\
h_2 &= g^{-1}[h'_2, \phi[h_1]] \\
h_3 &= g^{-1}[h'_3, \phi[h_{1:2}]] \\
h_4 &= g^{-1}[h'_4, \phi[h_{1:3}]].
\end{aligned} \tag{16.17}$$

$$\begin{aligned}
\mathbf{h}'_1 &= \mathbf{h}_1 + \mathbf{f}_1[\mathbf{h}_2, \phi_1] \\
\mathbf{h}'_2 &= \mathbf{h}_2 + \mathbf{f}_2[\mathbf{h}'_1, \phi_2],
\end{aligned} \tag{16.18}$$

$$\begin{aligned}
\mathbf{h}_2 &= \mathbf{h}'_2 - \mathbf{f}_2[\mathbf{h}'_1, \phi_2] \\
\mathbf{h}_1 &= \mathbf{h}'_1 - \mathbf{f}_1[\mathbf{h}_2, \phi_1].
\end{aligned} \tag{16.19}$$

$$\text{dist}[\mathbf{f}[z'], \mathbf{f}[z]] < \beta \cdot \text{dist}[z', z] \quad \forall z, z', \tag{16.20}$$

$$y = z + \mathbf{f}[z] \tag{16.21}$$

$$\begin{aligned}
\log \left[\left\| \mathbf{I} + \frac{\partial \mathbf{f}[\mathbf{h}, \phi]}{\partial \mathbf{h}} \right\| \right] &= \text{trace} \left[\log \left[\mathbf{I} + \frac{\partial \mathbf{f}[\mathbf{h}, \phi]}{\partial \mathbf{h}} \right] \right] \\
&= \sum_{k=1}^{\infty} \frac{(-1)^{k-1}}{k} \text{trace} \left[\frac{\partial \mathbf{f}[\mathbf{h}, \phi]}{\partial \mathbf{h}} \right]^k,
\end{aligned} \tag{16.22}$$

$$\begin{aligned}
\text{trace}[\mathbf{A}] &= \text{trace}[\mathbf{A} \mathbb{E}[\boldsymbol{\epsilon} \boldsymbol{\epsilon}^T]] \\
&= \text{trace}[\mathbb{E}[\mathbf{A} \boldsymbol{\epsilon} \boldsymbol{\epsilon}^T]] \\
&= \mathbb{E}[\text{trace}[\mathbf{A} \boldsymbol{\epsilon} \boldsymbol{\epsilon}^T]] \\
&= \mathbb{E}[\text{trace}[\boldsymbol{\epsilon}^T \mathbf{A} \boldsymbol{\epsilon}]] \\
&= \mathbb{E}[\boldsymbol{\epsilon}^T \mathbf{A} \boldsymbol{\epsilon}],
\end{aligned} \tag{16.23}$$

$$\begin{aligned}
\text{trace}[\mathbf{A}] &= \mathbb{E}[\boldsymbol{\epsilon}^T \mathbf{A} \boldsymbol{\epsilon}] \\
&\approx \frac{1}{I} \sum_{i=1}^I \boldsymbol{\epsilon}_i^T \mathbf{A} \boldsymbol{\epsilon}_i.
\end{aligned} \tag{16.24}$$

$$\hat{\phi} = \operatorname{argmin}_{\phi} \left[\text{KL} \left[\frac{1}{I} \sum_{i=1}^I \delta[\mathbf{x} - \mathbf{f}[\mathbf{z}_i, \phi]] \middle| \middle| q(\mathbf{x}) \right] \right]. \quad (16.25)$$

$$\hat{\phi} = \operatorname{argmin}_{\phi} \left[\text{KL} \left[\frac{1}{I} \sum_{i=1}^I \delta[\mathbf{x} - \mathbf{x}_i] \middle| \middle| Pr(\mathbf{x}_i, \phi) \right] \right]. \quad (16.26)$$

$$Pr(z) = \frac{1}{\sqrt{2\pi}} \exp \left[\frac{-z^2}{2} \right], \quad (16.27)$$

$$x = \mathbf{f}[z] = \frac{1}{1 + \exp[-z]}. \quad (16.28)$$

$$\mathbf{\Omega}_1 = \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & -5 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 2 \end{bmatrix} \quad \mathbf{\Omega}_2 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 2 & 4 & 0 & 0 \\ 1 & -1 & 2 & 0 \\ 4 & -2 & -2 & 1 \end{bmatrix}. \quad (16.29)$$

$$Pr(\mathbf{x}) = Pr(\mathbf{z}) \cdot \left| \frac{\partial \mathbf{f}[\mathbf{z}]}{\partial \mathbf{z}} \right|^{-1}. \quad (16.30)$$

$$\text{LReLU}[z] = \begin{cases} 0.1z & z < 0 \\ z & z \geq 0 \end{cases}. \quad (16.31)$$

$$\mathbf{f}[\mathbf{z}] = \left[\text{LReLU}[z_1], \text{LReLU}[z_2], \dots, \text{LReLU}[z_D] \right]^T. \quad (16.32)$$

$$\mathbf{h}' = \mathbf{f}[h, \phi] = \sqrt{[Kh - b + 1] \phi_b} + \sum_{k=1}^{b-1} \sqrt{\phi_k}, \quad (16.33)$$

Chapter 17

Variational autoencoders

$$Pr(\mathbf{x}) = \int Pr(\mathbf{x}, \mathbf{z}) d\mathbf{z}. \quad (17.1)$$

$$Pr(\mathbf{x}) = \int Pr(\mathbf{x}|\mathbf{z}) Pr(\mathbf{z}) d\mathbf{z}. \quad (17.2)$$

$$\begin{aligned} Pr(z = n) &= \lambda_n \\ Pr(x|z = n) &= \text{Norm}_x[\mu_n, \sigma_n^2]. \end{aligned} \quad (17.3)$$

$$\begin{aligned} Pr(x) &= \sum_{n=1}^N Pr(x, z = n) \\ &= \sum_{n=1}^N Pr(x|z = n) \cdot Pr(z = n) \\ &= \sum_{n=1}^N \lambda_n \cdot \text{Norm}_x[\mu_n, \sigma_n^2]. \end{aligned} \quad (17.4)$$

$$Pr(\mathbf{z}) = \text{Norm}_{\mathbf{z}}[\mathbf{0}, \mathbf{I}]. \quad (17.5)$$

$$Pr(\mathbf{x}|\mathbf{z}, \phi) = \text{Norm}_{\mathbf{x}}[\mathbf{f}[\mathbf{z}, \phi], \sigma^2 \mathbf{I}]. \quad (17.6)$$

$$\begin{aligned}
Pr(\mathbf{x}|\phi) &= \int Pr(\mathbf{x}, \mathbf{z}|\phi) d\mathbf{z} \\
&= \int Pr(\mathbf{x}|\mathbf{z}, \phi) \cdot Pr(\mathbf{z}) d\mathbf{z} \\
&= \int \text{Norm}_{\mathbf{x}}[\mathbf{f}[\mathbf{z}, \phi], \sigma^2 \mathbf{I}] \cdot \text{Norm}_{\mathbf{z}}[\mathbf{0}, \mathbf{I}] d\mathbf{z}.
\end{aligned} \tag{17.7}$$

$$\hat{\phi} = \underset{\phi}{\operatorname{argmax}} \left[\sum_{i=1}^I \log [Pr(\mathbf{x}_i|\phi)] \right], \tag{17.8}$$

$$Pr(\mathbf{x}_i|\phi) = \int \text{Norm}_{\mathbf{x}_i}[\mathbf{f}[\mathbf{z}, \phi], \sigma^2 \mathbf{I}] \cdot \text{Norm}_{\mathbf{z}}[\mathbf{0}, \mathbf{I}] d\mathbf{z}. \tag{17.9}$$

$$g[\mathbb{E}[y]] \geq \mathbb{E}[g[y]]. \tag{17.10}$$

$$\log [\mathbb{E}[y]] \geq \mathbb{E}[\log[y]], \tag{17.11}$$

$$\log \left[\int Pr(y) y dy \right] \geq \int Pr(y) \log[y] dy. \tag{17.12}$$

$$\log \left[\int Pr(y) h[y] dy \right] \geq \int Pr(y) \log[h[y]] dy. \tag{17.13}$$

$$\begin{aligned}
\log[Pr(\mathbf{x}|\phi)] &= \log \left[\int Pr(\mathbf{x}, \mathbf{z}|\phi) d\mathbf{z} \right] \\
&= \log \left[\int q(\mathbf{z}) \frac{Pr(\mathbf{x}, \mathbf{z}|\phi)}{q(\mathbf{z})} d\mathbf{z} \right],
\end{aligned} \tag{17.14}$$

$$\log \left[\int q(\mathbf{z}) \frac{Pr(\mathbf{x}, \mathbf{z}|\phi)}{q(\mathbf{z})} d\mathbf{z} \right] \geq \int q(\mathbf{z}) \log \left[\frac{Pr(\mathbf{x}, \mathbf{z}|\phi)}{q(\mathbf{z})} \right] d\mathbf{z}, \tag{17.15}$$

$$\text{ELBO}[\theta, \phi] = \int q(\mathbf{z}|\theta) \log \left[\frac{Pr(\mathbf{x}, \mathbf{z}|\phi)}{q(\mathbf{z}|\theta)} \right] d\mathbf{z}. \tag{17.16}$$

$$\begin{aligned}
\text{ELBO}[\boldsymbol{\theta}, \phi] &= \int q(\mathbf{z}|\boldsymbol{\theta}) \log \left[\frac{Pr(\mathbf{x}, \mathbf{z}|\phi)}{q(\mathbf{z}|\boldsymbol{\theta})} \right] d\mathbf{z} \\
&= \int q(\mathbf{z}|\boldsymbol{\theta}) \log \left[\frac{Pr(\mathbf{z}|\mathbf{x}, \phi) Pr(\mathbf{x}|\phi)}{q(\mathbf{z}|\boldsymbol{\theta})} \right] d\mathbf{z} \\
&= \int q(\mathbf{z}|\boldsymbol{\theta}) \log [Pr(\mathbf{x}|\phi)] d\mathbf{z} + \int q(\mathbf{z}|\boldsymbol{\theta}) \log \left[\frac{Pr(\mathbf{z}|\mathbf{x}, \phi)}{q(\mathbf{z}|\boldsymbol{\theta})} \right] d\mathbf{z} \\
&= \log [Pr(\mathbf{x}|\phi)] + \int q(\mathbf{z}|\boldsymbol{\theta}) \log \left[\frac{Pr(\mathbf{z}|\mathbf{x}, \phi)}{q(\mathbf{z}|\boldsymbol{\theta})} \right] d\mathbf{z} \\
&= \log [Pr(\mathbf{x}|\phi)] - D_{KL} [q(\mathbf{z}|\boldsymbol{\theta}) \parallel Pr(\mathbf{z}|\mathbf{x}, \phi)]. \tag{17.17}
\end{aligned}$$

$$\begin{aligned}
\text{ELBO}[\boldsymbol{\theta}, \phi] &= \int q(\mathbf{z}|\boldsymbol{\theta}) \log \left[\frac{Pr(\mathbf{x}, \mathbf{z}|\phi)}{q(\mathbf{z}|\boldsymbol{\theta})} \right] d\mathbf{z} \\
&= \int q(\mathbf{z}|\boldsymbol{\theta}) \log \left[\frac{Pr(\mathbf{x}|\mathbf{z}, \phi) Pr(\mathbf{z})}{q(\mathbf{z}|\boldsymbol{\theta})} \right] d\mathbf{z} \\
&= \int q(\mathbf{z}|\boldsymbol{\theta}) \log [Pr(\mathbf{x}|\mathbf{z}, \phi)] d\mathbf{z} + \int q(\mathbf{z}|\boldsymbol{\theta}) \log \left[\frac{Pr(\mathbf{z})}{q(\mathbf{z}|\boldsymbol{\theta})} \right] d\mathbf{z} \\
&= \int q(\mathbf{z}|\boldsymbol{\theta}) \log [Pr(\mathbf{x}|\mathbf{z}, \phi)] d\mathbf{z} - D_{KL} [q(\mathbf{z}|\boldsymbol{\theta}) \parallel Pr(\mathbf{z})], \tag{17.18}
\end{aligned}$$

$$Pr(\mathbf{z}|\mathbf{x}, \phi) = \frac{Pr(\mathbf{x}|\mathbf{z}, \phi) Pr(\mathbf{z})}{Pr(\mathbf{x}|\phi)}, \tag{17.19}$$

$$q(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}) = \text{Norm}_{\mathbf{z}} [\mathbf{g}_{\boldsymbol{\mu}}[\mathbf{x}, \boldsymbol{\theta}], \mathbf{g}_{\boldsymbol{\Sigma}}[\mathbf{x}, \boldsymbol{\theta}]], \tag{17.20}$$

$$\text{ELBO}[\boldsymbol{\theta}, \phi] = \int q(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}) \log [Pr(\mathbf{x}|\mathbf{z}, \phi)] d\mathbf{z} - D_{KL} [q(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}) \parallel Pr(\mathbf{z})], \tag{17.21}$$

$$\mathbb{E}_{\mathbf{z}} [\mathbf{a}[\mathbf{z}]] = \int \mathbf{a}[\mathbf{z}] q(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}) d\mathbf{z} \approx \frac{1}{N} \sum_{n=1}^N \mathbf{a}[\mathbf{z}_n^*], \tag{17.22}$$

$$\text{ELBO}[\boldsymbol{\theta}, \phi] \approx \log [Pr(\mathbf{x}|\mathbf{z}^*, \phi)] - D_{KL} [q(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}) \parallel Pr(\mathbf{z})]. \tag{17.23}$$

$$D_{KL} [q(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}) \parallel Pr(\mathbf{z})] = \frac{1}{2} \left(\text{Tr}[\boldsymbol{\Sigma}] + \boldsymbol{\mu}^T \boldsymbol{\mu} - D_{\mathbf{z}} - \log [\det[\boldsymbol{\Sigma}]] \right). \tag{17.24}$$

$$\mathbf{z}^* = \boldsymbol{\mu} + \boldsymbol{\Sigma}^{1/2} \boldsymbol{\epsilon}^*, \quad (17.25)$$

$$\begin{aligned} Pr(\mathbf{x}) &= \int Pr(\mathbf{x}|\mathbf{z})Pr(\mathbf{z})d\mathbf{z} \\ &= \mathbb{E}_{\mathbf{z}}[Pr(\mathbf{x}|\mathbf{z})] \\ &= \mathbb{E}_{\mathbf{z}}[\text{Norm}_{\mathbf{x}}[\mathbf{f}[\mathbf{z}, \boldsymbol{\phi}], \sigma^2 \mathbf{I}]]. \end{aligned} \quad (17.26)$$

$$Pr(\mathbf{x}) \approx \frac{1}{N} \sum_{n=1}^N Pr(\mathbf{x}|\mathbf{z}_n). \quad (17.27)$$

$$\begin{aligned} Pr(\mathbf{x}) &= \int Pr(\mathbf{x}|\mathbf{z})Pr(\mathbf{z})d\mathbf{z} \\ &= \int \frac{Pr(\mathbf{x}|\mathbf{z})Pr(\mathbf{z})}{q(\mathbf{z})} q(\mathbf{z})d\mathbf{z} \\ &= \mathbb{E}_{q(\mathbf{z})} \left[\frac{Pr(\mathbf{x}|\mathbf{z})Pr(\mathbf{z})}{q(\mathbf{z})} \right] \\ &\approx \frac{1}{N} \sum_{n=1}^N \frac{Pr(\mathbf{x}|\mathbf{z}_n)Pr(\mathbf{z}_n)}{q(\mathbf{z}_n)}, \end{aligned} \quad (17.28)$$

$$L_{\text{new}} = -\text{ELBO}[\boldsymbol{\theta}, \boldsymbol{\phi}] + \lambda_1 \mathbb{E}_{Pr(\mathbf{x})} [\text{r}_1[q(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})]] + \lambda_2 \text{r}_2[q(\mathbf{z}|\boldsymbol{\theta})]. \quad (17.29)$$

$$\text{ELBO}[\boldsymbol{\theta}, \boldsymbol{\phi}] \approx \log[Pr(\mathbf{x}|\mathbf{z}^*, \boldsymbol{\phi})] - \beta \cdot \text{D}_{KL}[q(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}) \parallel Pr(\mathbf{z})], \quad (17.30)$$

$$g[\mathbb{E}[y]] \leq \mathbb{E}[g[y]]. \quad (17.31)$$

$$\text{D}_{KL}[q(\mathbf{z}|\mathbf{x}) \parallel Pr(\mathbf{z}|\mathbf{x}, \boldsymbol{\phi})] = \int q(\mathbf{z}|\mathbf{x}) \log \left[\frac{q(\mathbf{z}|\mathbf{x})}{Pr(\mathbf{z}|\mathbf{x}, \boldsymbol{\phi})} \right] d\mathbf{z}. \quad (17.32)$$

$$\frac{\partial}{\partial \boldsymbol{\phi}} \mathbb{E}_{Pr(x|\boldsymbol{\phi})} [\text{f}[x]], \quad (17.33)$$

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\phi}} \mathbb{E}_{Pr(x|\boldsymbol{\phi})} [\text{f}[x]] &= \mathbb{E}_{Pr(x|\boldsymbol{\phi})} \left[\text{f}[x] \frac{\partial}{\partial \boldsymbol{\phi}} \log[Pr(x|\boldsymbol{\phi})] \right] \\ &\approx \frac{1}{I} \sum_{i=1}^I \text{f}[x_i] \frac{\partial}{\partial \boldsymbol{\phi}} \log[Pr(x_i|\boldsymbol{\phi})]. \end{aligned} \quad (17.34)$$

Chapter 18

Diffusion models

$$\begin{aligned}\mathbf{z}_1 &= \sqrt{1 - \beta_1} \cdot \mathbf{x} + \sqrt{\beta_1} \cdot \boldsymbol{\epsilon}_1 \\ \mathbf{z}_t &= \sqrt{1 - \beta_t} \cdot \mathbf{z}_{t-1} + \sqrt{\beta_t} \cdot \boldsymbol{\epsilon}_t \quad \forall t \in \{2, \dots, T\},\end{aligned}\tag{18.1}$$

$$\begin{aligned}q(\mathbf{z}_1|\mathbf{x}) &= \text{Norm}_{\mathbf{z}_1} \left[\sqrt{1 - \beta_1} \mathbf{x}, \beta_1 \mathbf{I} \right] \\ q(\mathbf{z}_t|\mathbf{z}_{t-1}) &= \text{Norm}_{\mathbf{z}_t} \left[\sqrt{1 - \beta_t} \mathbf{z}_{t-1}, \beta_t \mathbf{I} \right] \quad \forall t \in \{2, \dots, T\}.\end{aligned}\tag{18.2}$$

$$q(\mathbf{z}_{1\dots T}|\mathbf{x}) = q(\mathbf{z}_1|\mathbf{x}) \prod_{t=2}^T q(\mathbf{z}_t|\mathbf{z}_{t-1}).\tag{18.3}$$

$$\begin{aligned}\mathbf{z}_1 &= \sqrt{1 - \beta_1} \cdot \mathbf{x} + \sqrt{\beta_1} \cdot \boldsymbol{\epsilon}_1 \\ \mathbf{z}_2 &= \sqrt{1 - \beta_2} \cdot \mathbf{z}_1 + \sqrt{\beta_2} \cdot \boldsymbol{\epsilon}_2.\end{aligned}\tag{18.4}$$

$$\begin{aligned}\mathbf{z}_2 &= \sqrt{1 - \beta_2} \left(\sqrt{1 - \beta_1} \cdot \mathbf{x} + \sqrt{\beta_1} \cdot \boldsymbol{\epsilon}_1 \right) + \sqrt{\beta_2} \cdot \boldsymbol{\epsilon}_2 \\ &= \sqrt{1 - \beta_2} \left(\sqrt{1 - \beta_1} \cdot \mathbf{x} + \sqrt{1 - (1 - \beta_1)} \cdot \boldsymbol{\epsilon}_1 \right) + \sqrt{\beta_2} \cdot \boldsymbol{\epsilon}_2 \\ &= \sqrt{(1 - \beta_2)(1 - \beta_1)} \cdot \mathbf{x} + \sqrt{1 - \beta_2 - (1 - \beta_2)(1 - \beta_1)} \cdot \boldsymbol{\epsilon}_1 + \sqrt{\beta_2} \cdot \boldsymbol{\epsilon}_2.\end{aligned}\tag{18.5}$$

$$\mathbf{z}_2 = \sqrt{(1 - \beta_2)(1 - \beta_1)} \cdot \mathbf{x} + \sqrt{1 - (1 - \beta_2)(1 - \beta_1)} \cdot \boldsymbol{\epsilon},\tag{18.6}$$

$$\mathbf{z}_t = \sqrt{\alpha_t} \cdot \mathbf{x} + \sqrt{1 - \alpha_t} \cdot \boldsymbol{\epsilon},\tag{18.7}$$

$$q(\mathbf{z}_t|\mathbf{x}) = \text{Norm}_{\mathbf{z}_t} \left[\sqrt{\alpha_t} \cdot \mathbf{x}, (1 - \alpha_t)\mathbf{I} \right]. \quad (18.8)$$

$$\begin{aligned} q(\mathbf{z}_t) &= \iint q(\mathbf{z}_{1\dots t}, \mathbf{x}) d\mathbf{z}_{1\dots t-1} d\mathbf{x} \\ &= \iint q(\mathbf{z}_{1\dots t}|\mathbf{x}) Pr(\mathbf{x}) d\mathbf{z}_{1\dots t-1} d\mathbf{x}, \end{aligned} \quad (18.9)$$

$$q(\mathbf{z}_t) = \int q(\mathbf{z}_t|\mathbf{x}) Pr(\mathbf{x}) d\mathbf{x}. \quad (18.10)$$

$$q(\mathbf{z}_{t-1}|\mathbf{z}_t) = \frac{q(\mathbf{z}_t|\mathbf{z}_{t-1})q(\mathbf{z}_{t-1})}{q(\mathbf{z}_t)}. \quad (18.11)$$

$$\begin{aligned} q(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{x}) &= \frac{q(\mathbf{z}_t|\mathbf{z}_{t-1}, \mathbf{x})q(\mathbf{z}_{t-1}|\mathbf{x})}{q(\mathbf{z}_t|\mathbf{x})} \\ &\propto q(\mathbf{z}_t|\mathbf{z}_{t-1})q(\mathbf{z}_{t-1}|\mathbf{x}) \\ &= \text{Norm}_{\mathbf{z}_t} \left[\sqrt{1 - \beta_t} \cdot \mathbf{z}_{t-1}, \beta_t \mathbf{I} \right] \text{Norm}_{\mathbf{z}_{t-1}} \left[\sqrt{\alpha_{t-1}} \cdot \mathbf{x}, (1 - \alpha_{t-1})\mathbf{I} \right] \\ &\propto \text{Norm}_{\mathbf{z}_{t-1}} \left[\frac{1}{\sqrt{1 - \beta_t}} \mathbf{z}_t, \frac{\beta_t}{1 - \beta_t} \mathbf{I} \right] \text{Norm}_{\mathbf{z}_{t-1}} \left[\sqrt{\alpha_{t-1}} \cdot \mathbf{x}, (1 - \alpha_{t-1})\mathbf{I} \right] \end{aligned} \quad (18.12)$$

$$\text{Norm}_{\mathbf{v}}[\mathbf{A}\mathbf{w}, \mathbf{B}] \propto \text{Norm}_{\mathbf{w}} \left[(\mathbf{A}^T \mathbf{B}^{-1} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{B}^{-1} \mathbf{v}, (\mathbf{A}^T \mathbf{B}^{-1} \mathbf{A})^{-1} \right], \quad (18.13)$$

$$\begin{aligned} \text{Norm}_{\mathbf{w}}[\mathbf{a}, \mathbf{A}] \cdot \text{Norm}_{\mathbf{w}}[\mathbf{b}, \mathbf{B}] &\propto \\ &\text{Norm}_{\mathbf{w}} \left[(\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1} (\mathbf{A}^{-1} \mathbf{a} + \mathbf{B}^{-1} \mathbf{b}), (\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1} \right], \end{aligned} \quad (18.14)$$

$$q(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{x}) = \text{Norm}_{\mathbf{z}_{t-1}} \left[\frac{(1 - \alpha_{t-1})}{1 - \alpha_t} \sqrt{1 - \beta_t} \mathbf{z}_t + \frac{\sqrt{\alpha_{t-1}} \beta_t}{1 - \alpha_t} \mathbf{x}, \frac{\beta_t (1 - \alpha_{t-1})}{1 - \alpha_t} \mathbf{I} \right]. \quad (18.15)$$

$$\begin{aligned} Pr(\mathbf{z}_T) &= \text{Norm}_{\mathbf{z}_T}[\mathbf{0}, \mathbf{I}] \\ Pr(\mathbf{z}_{t-1}|\mathbf{z}_t, \phi_t) &= \text{Norm}_{\mathbf{z}_{t-1}} \left[\mathbf{f}_t[\mathbf{z}_t, \phi_t], \sigma_t^2 \mathbf{I} \right] \\ Pr(\mathbf{x}|\mathbf{z}_1, \phi_1) &= \text{Norm}_{\mathbf{x}} \left[\mathbf{f}_1[\mathbf{z}_1, \phi_1], \sigma_1^2 \mathbf{I} \right], \end{aligned} \quad (18.16)$$

$$Pr(\mathbf{x}, \mathbf{z}_{1...T} | \phi_{1...T}) = Pr(\mathbf{x} | \mathbf{z}_1, \phi_1) \prod_{t=2}^T Pr(\mathbf{z}_{t-1} | \mathbf{z}_t, \phi_t) \cdot Pr(\mathbf{z}_T). \quad (18.17)$$

$$Pr(\mathbf{x} | \phi_{1...T}) = \int Pr(\mathbf{x}, \mathbf{z}_{1...T} | \phi_{1...T}) d\mathbf{z}_{1...T}. \quad (18.18)$$

$$\hat{\phi}_{1...T} = \operatorname{argmax}_{\phi_{1...T}} \left[\sum_{i=1}^I \log [Pr(\mathbf{x}_i | \phi_{1...T})] \right]. \quad (18.19)$$

$$\begin{aligned} \log [Pr(\mathbf{x} | \phi_{1...T})] &= \log \left[\int Pr(\mathbf{x}, \mathbf{z}_{1...T} | \phi_{1...T}) d\mathbf{z}_{1...T} \right] \\ &= \log \left[\int q(\mathbf{z}_{1...T} | \mathbf{x}) \frac{Pr(\mathbf{x}, \mathbf{z}_{1...T} | \phi_{1...T})}{q(\mathbf{z}_{1...T} | \mathbf{x})} d\mathbf{z}_{1...T} \right] \\ &\geq \int q(\mathbf{z}_{1...T} | \mathbf{x}) \log \left[\frac{Pr(\mathbf{x}, \mathbf{z}_{1...T} | \phi_{1...T})}{q(\mathbf{z}_{1...T} | \mathbf{x})} \right] d\mathbf{z}_{1...T}. \end{aligned} \quad (18.20)$$

$$\text{ELBO}[\phi_{1...T}] = \int q(\mathbf{z}_{1...T} | \mathbf{x}) \log \left[\frac{Pr(\mathbf{x}, \mathbf{z}_{1...T} | \phi_{1...T})}{q(\mathbf{z}_{1...T} | \mathbf{x})} \right] d\mathbf{z}_{1...T}. \quad (18.21)$$

$$\begin{aligned} \log \left[\frac{Pr(\mathbf{x}, \mathbf{z}_{1...T} | \phi_{1...T})}{q(\mathbf{z}_{1...T} | \mathbf{x})} \right] &= \log \left[\frac{Pr(\mathbf{x} | \mathbf{z}_1, \phi_1) \prod_{t=2}^T Pr(\mathbf{z}_{t-1} | \mathbf{z}_t, \phi_t) \cdot Pr(\mathbf{z}_T)}{q(\mathbf{z}_1 | \mathbf{x}) \prod_{t=2}^T q(\mathbf{z}_t | \mathbf{z}_{t-1})} \right] \\ &= \log \left[\frac{Pr(\mathbf{x} | \mathbf{z}_1, \phi_1)}{q(\mathbf{z}_1 | \mathbf{x})} \right] + \log \left[\frac{\prod_{t=2}^T Pr(\mathbf{z}_{t-1} | \mathbf{z}_t, \phi_t)}{\prod_{t=2}^T q(\mathbf{z}_t | \mathbf{z}_{t-1})} \right] + \log [Pr(\mathbf{z}_T)]. \end{aligned} \quad (18.22)$$

$$q(\mathbf{z}_t | \mathbf{z}_{t-1}) = q(\mathbf{z}_t | \mathbf{z}_{t-1}, \mathbf{x}) = \frac{q(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{x}) q(\mathbf{z}_t | \mathbf{x})}{q(\mathbf{z}_{t-1} | \mathbf{x})}, \quad (18.23)$$

$$\begin{aligned} &\log \left[\frac{Pr(\mathbf{x}, \mathbf{z}_{1...T} | \phi_{1...T})}{q(\mathbf{z}_{1...T} | \mathbf{x})} \right] \\ &= \log \left[\frac{Pr(\mathbf{x} | \mathbf{z}_1, \phi_1)}{q(\mathbf{z}_1 | \mathbf{x})} \right] + \log \left[\frac{\prod_{t=2}^T Pr(\mathbf{z}_{t-1} | \mathbf{z}_t, \phi_t) \cdot q(\mathbf{z}_{t-1} | \mathbf{x})}{\prod_{t=2}^T q(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{x}) \cdot q(\mathbf{z}_t | \mathbf{x})} \right] + \log [Pr(\mathbf{z}_T)] \\ &= \log [Pr(\mathbf{x} | \mathbf{z}_1, \phi_1)] + \log \left[\frac{\prod_{t=2}^T Pr(\mathbf{z}_{t-1} | \mathbf{z}_t, \phi_t)}{\prod_{t=2}^T q(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{x})} \right] + \log \left[\frac{Pr(\mathbf{z}_T)}{q(\mathbf{z}_T | \mathbf{x})} \right] \\ &\approx \log [Pr(\mathbf{x} | \mathbf{z}_1, \phi_1)] + \sum_{t=2}^T \log \left[\frac{Pr(\mathbf{z}_{t-1} | \mathbf{z}_t, \phi_t)}{q(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{x})} \right], \end{aligned} \quad (18.24)$$

$$\text{ELBO}[\phi_{1...T}] \quad (18.25)$$

$$\begin{aligned} &= \int q(\mathbf{z}_{1...T}|\mathbf{x}) \log \left[\frac{Pr(\mathbf{x}, \mathbf{z}_{1...T}|\phi_{1...T})}{q(\mathbf{z}_{1...T}|\mathbf{x})} \right] d\mathbf{z}_{1...T} \\ &\approx \int q(\mathbf{z}_{1...T}|\mathbf{x}) \left(\log [Pr(\mathbf{x}|\mathbf{z}_1, \phi_1)] + \sum_{t=2}^T \log \left[\frac{Pr(\mathbf{z}_{t-1}|\mathbf{z}_t, \phi_t)}{q(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{x})} \right] \right) d\mathbf{z}_{1...T} \\ &= \mathbb{E}_{q(\mathbf{z}_1|\mathbf{x})} [\log [Pr(\mathbf{x}|\mathbf{z}_1, \phi_1)]] - \sum_{t=2}^T \mathbb{E}_{q(\mathbf{z}_t|\mathbf{x})} \left[D_{KL} [q(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{x}) || Pr(\mathbf{z}_{t-1}|\mathbf{z}_t, \phi_t)] \right], \end{aligned}$$

$$Pr(\mathbf{x}|\mathbf{z}_1, \phi_1) = \text{Norm}_{\mathbf{x}} [\mathbf{f}_1[\mathbf{z}_1, \phi_1], \sigma_1^2 \mathbf{I}], \quad (18.26)$$

$$Pr(\mathbf{z}_{t-1}|\mathbf{z}_t, \phi_t) = \text{Norm}_{\mathbf{z}_{t-1}} [\mathbf{f}_t[\mathbf{z}_t, \phi_t], \sigma_t^2 \mathbf{I}] \quad (18.27)$$

$$q(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{x}) = \text{Norm}_{\mathbf{z}_{t-1}} \left[\frac{(1-\alpha_{t-1})}{1-\alpha_t} \sqrt{1-\beta_t} \mathbf{z}_t + \frac{\sqrt{\alpha_{t-1}} \beta_t}{1-\alpha_t} \mathbf{x}, \frac{\beta_t(1-\alpha_{t-1})}{1-\alpha_t} \mathbf{I} \right].$$

$$D_{KL} [q(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{x}) || Pr(\mathbf{z}_{t-1}|\mathbf{z}_t, \phi_t)] = \quad (18.28)$$

$$\frac{1}{2\sigma_t^2} \left\| \frac{(1-\alpha_{t-1})}{1-\alpha_t} \sqrt{1-\beta_t} \mathbf{z}_t + \frac{\sqrt{\alpha_{t-1}} \beta_t}{1-\alpha_t} \mathbf{x} - \mathbf{f}_t[\mathbf{z}_t, \phi_t] \right\|^2 + C.$$

$$\begin{aligned} L[\phi_{1...T}] &= \sum_{i=1}^I \overbrace{\left(-\log [\text{Norm}_{\mathbf{x}_i} [\mathbf{f}_1[\mathbf{z}_{i1}, \phi_1], \sigma_1^2 \mathbf{I}]] \right)}^{\text{reconstruction term}} \\ &\quad + \sum_{t=2}^T \frac{1}{2\sigma_t^2} \left\| \underbrace{\frac{1-\alpha_{t-1}}{1-\alpha_t} \sqrt{1-\beta_t} \mathbf{z}_{it} + \frac{\sqrt{\alpha_{t-1}} \beta_t}{1-\alpha_t} \mathbf{x}_i}_{\text{target, mean of } q(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{x})} - \underbrace{\mathbf{f}_t[\mathbf{z}_{it}, \phi_t]}_{\text{predicted } \mathbf{z}_{t-1}} \right\|^2 \end{aligned} \quad (18.29)$$

$$\mathbf{z}_t = \sqrt{\alpha_t} \cdot \mathbf{x} + \sqrt{1-\alpha_t} \cdot \epsilon. \quad (18.30)$$

$$\mathbf{x} = \frac{1}{\sqrt{\alpha_t}} \cdot \mathbf{z}_t - \frac{\sqrt{1-\alpha_t}}{\sqrt{\alpha_t}} \cdot \epsilon. \quad (18.31)$$

$$\begin{aligned}
& \frac{(1 - \alpha_{t-1})}{1 - \alpha_t} \sqrt{1 - \beta_t} \mathbf{z}_t + \frac{\sqrt{\alpha_{t-1}} \beta_t}{1 - \alpha_t} \mathbf{x} \\
&= \frac{(1 - \alpha_{t-1})}{1 - \alpha_t} \sqrt{1 - \beta_t} \mathbf{z}_t + \frac{\sqrt{\alpha_{t-1}} \beta_t}{1 - \alpha_t} \left(\frac{1}{\sqrt{\alpha_t}} \mathbf{z}_t - \frac{\sqrt{1 - \alpha_t}}{\sqrt{\alpha_t}} \boldsymbol{\epsilon} \right) \\
&= \frac{(1 - \alpha_{t-1})}{1 - \alpha_t} \sqrt{1 - \beta_t} \mathbf{z}_t + \frac{\beta_t}{1 - \alpha_t} \left(\frac{1}{\sqrt{1 - \beta_t}} \mathbf{z}_t - \frac{\sqrt{1 - \alpha_t}}{\sqrt{1 - \beta_t}} \boldsymbol{\epsilon} \right),
\end{aligned} \tag{18.32}$$

$$\begin{aligned}
& \frac{(1 - \alpha_{t-1})}{1 - \alpha_t} \sqrt{1 - \beta_t} \mathbf{z}_t + \frac{\sqrt{\alpha_{t-1}} \beta_t}{1 - \alpha_t} \mathbf{x} \\
&= \left(\frac{(1 - \alpha_{t-1}) \sqrt{1 - \beta_t}}{1 - \alpha_t} + \frac{\beta_t}{(1 - \alpha_t) \sqrt{1 - \beta_t}} \right) \mathbf{z}_t - \frac{\beta_t}{\sqrt{1 - \alpha_t} \sqrt{1 - \beta_t}} \boldsymbol{\epsilon} \\
&= \left(\frac{(1 - \alpha_{t-1})(1 - \beta_t)}{(1 - \alpha_t) \sqrt{1 - \beta_t}} + \frac{\beta_t}{(1 - \alpha_t) \sqrt{1 - \beta_t}} \right) \mathbf{z}_t - \frac{\beta_t}{\sqrt{1 - \alpha_t} \sqrt{1 - \beta_t}} \boldsymbol{\epsilon} \\
&= \frac{(1 - \alpha_{t-1})(1 - \beta_t) + \beta_t}{(1 - \alpha_t) \sqrt{1 - \beta_t}} \mathbf{z}_t - \frac{\beta_t}{\sqrt{1 - \alpha_t} \sqrt{1 - \beta_t}} \boldsymbol{\epsilon} \\
&= \frac{1 - \alpha_t}{(1 - \alpha_t) \sqrt{1 - \beta_t}} \mathbf{z}_t - \frac{\beta_t}{\sqrt{1 - \alpha_t} \sqrt{1 - \beta_t}} \boldsymbol{\epsilon} \\
&= \frac{1}{\sqrt{1 - \beta_t}} \mathbf{z}_t - \frac{\beta_t}{\sqrt{1 - \alpha_t} \sqrt{1 - \beta_t}} \boldsymbol{\epsilon},
\end{aligned} \tag{18.33}$$

$$\begin{aligned}
L[\phi_{1...T}] &= \sum_{i=1}^I \left(-\log \left[\text{Norm}_{\mathbf{x}_i} [\mathbf{f}_1[\mathbf{z}_{i1}, \phi_1], \sigma_1^2 \mathbf{I}] \right] \right. \\
&\quad \left. + \sum_{t=2}^T \frac{1}{2\sigma_t^2} \left\| \left(\frac{1}{\sqrt{1 - \beta_t}} \mathbf{z}_{it} - \frac{\beta_t}{\sqrt{1 - \alpha_t} \sqrt{1 - \beta_t}} \boldsymbol{\epsilon}_{it} \right) - \mathbf{f}_t[\mathbf{z}_{it}, \phi_t] \right\|^2 \right).
\end{aligned} \tag{18.34}$$

$$\mathbf{f}_t[\mathbf{z}_t, \phi_t] = \frac{1}{\sqrt{1 - \beta_t}} \mathbf{z}_t - \frac{\beta_t}{\sqrt{1 - \alpha_t} \sqrt{1 - \beta_t}} \mathbf{g}_t[\mathbf{z}_t, \phi_t]. \tag{18.35}$$

$$\begin{aligned}
L[\phi_{1...T}] &= \\
&\sum_{i=1}^I -\log \left[\text{Norm}_{\mathbf{x}_i} [\mathbf{f}_1[\mathbf{z}_{i1}, \phi_1], \sigma_1^2 \mathbf{I}] \right] + \sum_{t=2}^T \frac{\beta_t^2}{(1 - \alpha_t)(1 - \beta_t)2\sigma_t^2} \left\| \mathbf{g}_t[\mathbf{z}_{it}, \phi_t] - \boldsymbol{\epsilon}_{it} \right\|^2.
\end{aligned} \tag{18.36}$$

$$L[\phi_{1...T}] = \sum_{i=1}^I \frac{1}{2\sigma_1^2} \left\| \mathbf{x}_i - \mathbf{f}_1[\mathbf{z}_{i1}, \phi_1] \right\|^2 + \sum_{t=2}^T \frac{\beta_t^2}{(1 - \alpha_t)(1 - \beta_t)2\sigma_t^2} \left\| \mathbf{g}_t[\mathbf{z}_{it}, \phi_t] - \boldsymbol{\epsilon}_{it} \right\|^2 + \mathcal{C} \tag{18.37}$$

$$\frac{1}{2\sigma_1^2} \left\| \mathbf{x}_i - \mathbf{f}_1[\mathbf{z}_{i1}, \phi_1] \right\|^2 = \frac{1}{2\sigma_1^2} \left\| \frac{\beta_1}{\sqrt{1-\alpha_1}\sqrt{1-\beta_1}} \mathbf{g}_1[\mathbf{z}_{i1}, \phi_1] - \frac{\beta_1}{\sqrt{1-\alpha_1}\sqrt{1-\beta_1}} \boldsymbol{\epsilon}_{i1} \right\|^2 \quad (18.38)$$

$$L[\phi_{1...T}] = \sum_{i=1}^I \sum_{t=1}^T \frac{\beta_t^2}{(1-\alpha_t)(1-\beta_t)2\sigma_t^2} \left\| \mathbf{g}_t[\mathbf{z}_{it}, \phi_t] - \boldsymbol{\epsilon}_{it} \right\|^2, \quad (18.39)$$

$$\begin{aligned} L[\phi_{1...T}] &= \sum_{i=1}^I \sum_{t=1}^T \left\| \mathbf{g}_t[\mathbf{z}_{it}, \phi_t] - \boldsymbol{\epsilon}_{it} \right\|^2 \\ &= \sum_{i=1}^I \sum_{t=1}^T \left\| \mathbf{g}_t \left[\sqrt{\alpha_t} \cdot \mathbf{x}_i + \sqrt{1-\alpha_t} \cdot \boldsymbol{\epsilon}_{it}, \phi_t \right] - \boldsymbol{\epsilon}_{it} \right\|^2, \end{aligned} \quad (18.40)$$

$$\mathbf{z}_{t-1} = \hat{\mathbf{z}}_{t-1} + \sigma_t^2 \frac{\partial \log[\Pr(c|\mathbf{z}_t)]}{\partial \mathbf{z}_t} + \sigma_t \boldsymbol{\epsilon}. \quad (18.41)$$

$$\mathbf{x}_t = \sqrt{1-\beta_t} \cdot \mathbf{x}_{t-1} + \sqrt{\beta_t} \cdot \boldsymbol{\epsilon}_t, \quad (18.42)$$

$$z = a \cdot \epsilon_1 + b \cdot \epsilon_2, \quad (18.43)$$

$$\begin{aligned} \mathbb{E}[z] &= 0 \\ \text{Var}[z] &= a^2 + b^2, \end{aligned} \quad (18.44)$$

$$\mathbf{z}_3 = \sqrt{(1-\beta_3)(1-\beta_2)(1-\beta_1)} \cdot \mathbf{x} + \sqrt{1-(1-\beta_3)(1-\beta_2)(1-\beta_1)} \cdot \boldsymbol{\epsilon}' \quad (18.45)$$

$$\text{Norm}_{\mathbf{v}}[\mathbf{A}\mathbf{w}, \mathbf{B}] \propto \text{Norm}_{\mathbf{w}} \left[(\mathbf{A}^T \mathbf{B}^{-1} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{B}^{-1} \mathbf{v}, (\mathbf{A}^T \mathbf{B}^{-1} \mathbf{A})^{-1} \right]. \quad (18.46)$$

$$\text{Norm}_{\mathbf{x}}[\mathbf{a}, \mathbf{A}] \text{Norm}_{\mathbf{x}}[\mathbf{b}, \mathbf{B}] \propto \text{Norm}_{\mathbf{x}} \left[(\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1} (\mathbf{A}^{-1} \mathbf{a} + \mathbf{B}^{-1} \mathbf{b}), (\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1} \right]. \quad (18.47)$$

$$D_{KL}\left[\text{Norm}_{\mathbf{w}}[\mathbf{a}, \mathbf{A}] \parallel \text{Norm}_{\mathbf{w}}[\mathbf{b}, \mathbf{B}]\right] = \frac{1}{2} \left(\text{tr} [\mathbf{B}^{-1} \mathbf{A}] - d + (\mathbf{a} - \mathbf{b})^T \mathbf{B}^{-1} (\mathbf{a} - \mathbf{b}) + \log \left[\frac{|\mathbf{B}|}{|\mathbf{A}|} \right] \right). \quad (18.48)$$

$$\sqrt{\frac{\alpha_t}{\alpha_{t-1}}} = \sqrt{1 - \beta_t}. \quad (18.49)$$

$$\frac{(1 - \alpha_{t-1})(1 - \beta_t) + \beta_t}{(1 - \alpha_t)\sqrt{1 - \beta_t}} = \frac{1}{\sqrt{1 - \beta_t}}. \quad (18.50)$$

Chapter 19

Reinforcement learning

$$G_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}. \quad (19.1)$$

$$v[s_t|\pi] = \mathbb{E}[G_t|s_t, \pi]. \quad (19.2)$$

$$q[s_t, a_t|\pi] = \mathbb{E}[G_t|s_t, a_t, \pi]. \quad (19.3)$$

$$v^*[s_t] = \max_{\pi} \left[\mathbb{E}[G_t|s_t, \pi] \right]. \quad (19.4)$$

$$q^*[s_t, a_t] = \max_{\pi} \left[\mathbb{E}[G_t|s_t, a_t, \pi] \right]. \quad (19.5)$$

$$\pi[a_t|s_t] \leftarrow \operatorname{argmax}_{a_t} [q^*[s_t, a_t]]. \quad (19.6)$$

$$v[s_t] = \sum_{a_t} \pi[a_t|s_t] q[s_t, a_t]. \quad (19.7)$$

$$q[s_t, a_t] = r[s_t, a_t] + \gamma \cdot \sum_{s_{t+1}} Pr(s_{t+1}|s_t, a_t) v[s_{t+1}]. \quad (19.8)$$

$$v[s_t] = \sum_{a_t} \pi[a_t|s_t] \left(r[s_t, a_t] + \gamma \cdot \sum_{s_{t+1}} Pr(s_{t+1}|s_t, a_t) v[s_{t+1}] \right). \quad (19.9)$$

$$q[s_t, a_t] = r[s_t, a_t] + \gamma \cdot \sum_{s_{t+1}} Pr(s_{t+1}|s_t, a_t) \left(\sum_{a_{t+1}} \pi[a_{t+1}|s_{t+1}] q[s_{t+1}, a_{t+1}] \right). \quad (19.10)$$

$$v[s_t] \leftarrow \sum_{a_t} \pi[a_t|s_t] \left(r[s_t, a_t] + \gamma \cdot \sum_{s_{t+1}} Pr(s_{t+1}|s_t, a_t) v[s_{t+1}] \right), \quad (19.11)$$

$$\pi[a_t|s_t] \leftarrow \operatorname{argmax}_{a_t} \left[r[s_t, a_t] + \gamma \cdot \sum_{s_{t+1}} Pr(s_{t+1}|s_t, a_t) v[s_{t+1}] \right]. \quad (19.12)$$

$$\pi[a|s] \leftarrow \operatorname{argmax}_a [q[s, a]]. \quad (19.13)$$

$$q[s_t, a_t] \leftarrow q[s_t, a_t] + \alpha \left(r[s_t, a_t] + \gamma \cdot q[s_{t+1}, a_{t+1}] - q[s_t, a_t] \right), \quad (19.14)$$

$$q[s_t, a_t] \leftarrow q[s_t, a_t] + \alpha \left(r[s_t, a_t] + \gamma \cdot \max_a [q[s_{t+1}, a]] - q[s_t, a_t] \right), \quad (19.15)$$

$$L[\phi] = \left(r[\mathbf{s}_t, a_t] + \gamma \cdot \max_a [q[\mathbf{s}_{t+1}, a, \phi]] - q[\mathbf{s}_t, a_t, \phi] \right)^2, \quad (19.16)$$

$$\phi \leftarrow \phi + \alpha \left(r[\mathbf{s}_t, a_t] + \gamma \cdot \max_a [q[\mathbf{s}_{t+1}, a, \phi]] - q[\mathbf{s}_t, a_t, \phi] \right) \frac{\partial q[\mathbf{s}_t, a_t, \phi]}{\partial \phi}. \quad (19.17)$$

$$\phi \leftarrow \phi + \alpha \left(r[\mathbf{s}_t, a_t] + \gamma \cdot \max_a [q[\mathbf{s}_{t+1}, a, \phi^-]] - q[\mathbf{s}_t, a_t, \phi] \right) \frac{\partial q[\mathbf{s}_t, a_t, \phi]}{\partial \phi}. \quad (19.18)$$

$$q[s_t, a_t] \leftarrow q[s_t, a_t] + \alpha \left(r[s_t, a_t] + \gamma \cdot \max_a [q[s_{t+1}, a]] - q[s_t, a_t] \right) \quad (19.19)$$

$$\begin{aligned} q_1[s_t, a_t] &\leftarrow q_1[s_t, a_t] + \alpha \left(r[s_t, a_t] + \gamma \cdot q_2 \left[s_{t+1}, \operatorname{argmax}_a [q_1[s_{t+1}, a]] \right] - q_1[s_t, a_t] \right) \\ q_2[s_t, a_t] &\leftarrow q_2[s_t, a_t] + \alpha \left(r[s_t, a_t] + \gamma \cdot q_1 \left[s_{t+1}, \operatorname{argmax}_a [q_2[s_{t+1}, a]] \right] - q_2[s_t, a_t] \right). \end{aligned} \quad (19.20)$$

$$\begin{aligned}
\phi_1 &\leftarrow \phi_1 + \alpha \left(r[\mathbf{s}_t, a_t] + \gamma \cdot q \left[\mathbf{s}_{t+1}, \arg\max_a [q[\mathbf{s}_{t+1}, a, \phi_1]], \phi_2 \right] - q[\mathbf{s}_t, a_t, \phi_1] \right) \frac{\partial q[\mathbf{s}_t, a_t, \phi_1]}{\partial \phi_1} \\
\phi_2 &\leftarrow \phi_2 + \alpha \left(r[\mathbf{s}_t, a_t] + \gamma \cdot q \left[\mathbf{s}_{t+1}, \arg\max_a [q[\mathbf{s}_{t+1}, a, \phi_2]], \phi_1 \right] - q[\mathbf{s}_t, a_t, \phi_2] \right) \frac{\partial q[\mathbf{s}_t, a_t, \phi_2]}{\partial \phi_2}.
\end{aligned} \tag{19.21}$$

$$Pr(\boldsymbol{\tau}|\boldsymbol{\theta}) = Pr(\mathbf{s}_1) \prod_{t=1}^T \pi[a_t|\mathbf{s}_t, \boldsymbol{\theta}] Pr(\mathbf{s}_{t+1}|\mathbf{s}_t, a_t). \tag{19.22}$$

$$\boldsymbol{\theta} = \arg\max_{\boldsymbol{\theta}} \left[\mathbb{E}_{\boldsymbol{\tau}} [r[\boldsymbol{\tau}]] \right] = \arg\max_{\boldsymbol{\theta}} \left[\int Pr(\boldsymbol{\tau}|\boldsymbol{\theta}) r[\boldsymbol{\tau}] d\boldsymbol{\tau} \right], \tag{19.23}$$

$$\begin{aligned}
\boldsymbol{\theta} &\leftarrow \boldsymbol{\theta} + \alpha \cdot \frac{\partial}{\partial \boldsymbol{\theta}} \int Pr(\boldsymbol{\tau}|\boldsymbol{\theta}) r[\boldsymbol{\tau}] d\boldsymbol{\tau} \\
&= \boldsymbol{\theta} + \alpha \cdot \int \frac{\partial Pr(\boldsymbol{\tau}|\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} r[\boldsymbol{\tau}] d\boldsymbol{\tau}.
\end{aligned} \tag{19.24}$$

$$\begin{aligned}
\boldsymbol{\theta} &\leftarrow \boldsymbol{\theta} + \alpha \cdot \int \frac{\partial Pr(\boldsymbol{\tau}|\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} r[\boldsymbol{\tau}] d\boldsymbol{\tau} \\
&= \boldsymbol{\theta} + \alpha \cdot \int Pr(\boldsymbol{\tau}|\boldsymbol{\theta}) \frac{1}{Pr(\boldsymbol{\tau}|\boldsymbol{\theta})} \frac{\partial Pr(\boldsymbol{\tau}|\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} r[\boldsymbol{\tau}] d\boldsymbol{\tau} \\
&\approx \boldsymbol{\theta} + \alpha \cdot \frac{1}{I} \sum_{i=1}^I \frac{1}{Pr(\boldsymbol{\tau}_i|\boldsymbol{\theta})} \frac{\partial Pr(\boldsymbol{\tau}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} r[\boldsymbol{\tau}_i].
\end{aligned} \tag{19.25}$$

$$\frac{\partial \log[f[z]]}{\partial z} = \frac{1}{f[z]} \frac{\partial f[z]}{\partial z}, \tag{19.26}$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \cdot \frac{1}{I} \sum_{i=1}^I \frac{\partial \log[Pr(\boldsymbol{\tau}_i|\boldsymbol{\theta})]}{\partial \boldsymbol{\theta}} r[\boldsymbol{\tau}_i]. \tag{19.27}$$

$$\begin{aligned}
\log[Pr(\boldsymbol{\tau}|\boldsymbol{\theta})] &= \log \left[Pr(\mathbf{s}_1) \prod_{t=1}^T \pi[a_t|\mathbf{s}_t, \boldsymbol{\theta}] Pr(\mathbf{s}_{t+1}|\mathbf{s}_t, a_t) \right] \\
&= \log[Pr(\mathbf{s}_1)] + \sum_{t=1}^T \log[\pi[a_t|\mathbf{s}_t, \boldsymbol{\theta}]] + \sum_{t=1}^T \log[Pr(\mathbf{s}_{t+1}|\mathbf{s}_t, a_t)],
\end{aligned} \tag{19.28}$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \cdot \frac{1}{I} \sum_{i=1}^I \sum_{t=1}^T \frac{\partial \log[\pi[a_{it}|\mathbf{s}_{it}, \boldsymbol{\theta}]]}{\partial \boldsymbol{\theta}} r[\boldsymbol{\tau}_i], \quad (19.29)$$

$$r[\boldsymbol{\tau}_i] = \sum_{t=1}^T r_{i,t+1} = \sum_{k=1}^t r_{i,k+1} + \sum_{k=t}^T r_{i,k+1}, \quad (19.30)$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \cdot \frac{1}{I} \sum_{i=1}^I \sum_{t=1}^T \frac{\partial \log[\pi[a_{it}|\mathbf{s}_{it}, \boldsymbol{\theta}]]}{\partial \boldsymbol{\theta}} \sum_{k=t}^T r_{i,k+1}. \quad (19.31)$$

$$r[\boldsymbol{\tau}_{it}] = \sum_{k=t+1}^T \gamma^{k-t-1} r_{i,k+1}, \quad (19.32)$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \cdot \gamma^t \frac{\partial \log[\pi_{a_{it}}[\mathbf{s}_{it}, \boldsymbol{\theta}]]}{\partial \boldsymbol{\theta}} r[\boldsymbol{\tau}_{it}] \quad \forall i, t, \quad (19.33)$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \cdot \frac{1}{I} \sum_{i=1}^I \sum_{t=1}^T \frac{\partial \log[\pi_{a_{it}}[\mathbf{s}_{it}, \boldsymbol{\theta}]]}{\partial \boldsymbol{\theta}} (r[\boldsymbol{\tau}_{it}] - b). \quad (19.34)$$

$$\mathbb{E}_{\boldsymbol{\tau}} \left[\sum_{t=1}^T \frac{\partial \log[\pi_{a_{it}}[\mathbf{s}_{it}, \boldsymbol{\theta}]]}{\partial \boldsymbol{\theta}} \cdot b \right] = 0, \quad (19.35)$$

$$b = \sum_i \frac{\sum_{t=1}^T (\partial \log[\pi_{a_{it}}[\mathbf{s}_{it}, \boldsymbol{\theta}]] / \partial \boldsymbol{\theta})^2 r[\boldsymbol{\tau}_{it}]}{\sum_{t=1}^T (\partial \log[\pi_{a_{it}}[\mathbf{s}_{it}, \boldsymbol{\theta}]] / \partial \boldsymbol{\theta})^2}. \quad (19.36)$$

$$b = \frac{1}{I} \sum_i r[\boldsymbol{\tau}_i]. \quad (19.37)$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \cdot \frac{1}{I} \sum_{i=1}^I \sum_{t=1}^T \frac{\partial \log[\pi_{a_{it}}[\mathbf{s}_{it}, \boldsymbol{\theta}]]}{\partial \boldsymbol{\theta}} (r[\boldsymbol{\tau}_{it}] - b[\mathbf{s}_{it}]). \quad (19.38)$$

$$L[\phi] = \sum_{i=1}^I \sum_{t=1}^T \left(v[\mathbf{s}_{it}, \phi] - \sum_{j=t}^T r_{i,j+1} \right)^2. \quad (19.39)$$

$$r[\tau_{it}] \approx r_{i,t+1} + \gamma \cdot v[\mathbf{s}_{i,t+1}, \phi]. \quad (19.40)$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \cdot \frac{1}{I} \sum_{i=1}^I \sum_{t=1}^T \frac{\partial \log[Pr(a_{it}|\mathbf{s}_{it}, \boldsymbol{\theta})]}{\partial \boldsymbol{\theta}} \left(r_{i,t+1} + \gamma \cdot v[\mathbf{s}_{i,t+1}, \phi] - v[\mathbf{s}_{i,t}, \phi] \right). \quad (19.41)$$

$$L[\phi] = \sum_{i=1}^I \sum_{t=1}^T (r_{i,t+1} + \gamma \cdot v[\mathbf{s}_{i,t+1}, \phi] - v[\mathbf{s}_{i,t}, \phi])^2. \quad (19.42)$$

$$\pi'[a_t|s_t] \leftarrow \operatorname{argmax}_{a_t} \left[r[s_t, a_t] + \gamma \cdot \sum_{s_{t+1}} Pr(s_{t+1}|s_t, a_t) v[s_{t+1}|\pi] \right]. \quad (19.43)$$

$$\begin{aligned} v[s_t|\pi] &\leq q[s_t, \pi'[a_t|s_t]|\pi] \\ &= \mathbb{E}_{\pi'} [r_{t+1} + \gamma \cdot v[s_{t+1}|\pi]]. \end{aligned} \quad (19.44)$$

$$\pi[a|s] = \frac{\exp[q[s, a]/\tau]}{\sum_{a'} \exp[q[s, a']/\tau]}. \quad (19.45)$$

$$f[q[s, a]] = r[s, a] + \gamma \cdot \max_a [q[s', a]]. \quad (19.46)$$

$$\left\| f[q_1[s, a]] - f[q_2[s, a]] \right\|_{\infty} < \left\| q_1[s, a] - q_2[s, a] \right\|_{\infty} \quad \forall q_1, q_2. \quad (19.47)$$

$$\mathbb{E}_{\boldsymbol{\tau}} \left[\frac{\partial}{\partial \boldsymbol{\theta}} \log[Pr(\boldsymbol{\tau}|\boldsymbol{\theta})] b \right] = 0, \quad (19.48)$$

$$a' = a - c(b - \mu_b). \quad (19.49)$$

$$\mathbb{E}_{\boldsymbol{\tau}} \left[g[\boldsymbol{\theta}] (r[\boldsymbol{\tau}_t] - b) \right], \quad (19.50)$$

$$g[\theta] = \sum_{t=1}^T \frac{\partial \log[Pr(a_t | \mathbf{s}_t, \boldsymbol{\theta})]}{\partial \boldsymbol{\theta}}, \quad (19.51)$$

$$r[\boldsymbol{\tau}_t] = \sum_{k=t}^T r_k. \quad (19.52)$$

$$b = \frac{\mathbb{E}[g[\boldsymbol{\tau}]^2] r[\boldsymbol{\tau}]}{\mathbb{E}[g[\boldsymbol{\tau}]^2]}. \quad (19.53)$$

Chapter 20

Why does deep learning work?

Chapter 21

Ethics