

PART IV

Generation

20 Generative models: an overview

20.1 Introduction

A **generative model** is a joint probability distribution $p(\mathbf{x})$, for $\mathbf{x} \in \mathcal{X}$. In some cases, the model may be conditioned on inputs or covariates $\mathbf{c} \in \mathcal{C}$, which gives rise to a **conditional generative model** of the form $p(\mathbf{x}|\mathbf{c})$.

There are many kinds of generative model. We give a brief summary in Section 20.2, and go into more detail in subsequent chapters. See also [Tom22] for a recent book on this topic that goes into more depth.

20.2 Types of generative model

There are many kinds of generative model, some of which we list in Table 20.1. At a high level, we can distinguish between **deep generative models (DGM)** — which use deep neural network to learn a complex mapping from a single latent vector \mathbf{z} to the observed data \mathbf{x} — and more “classical” **probabilistic graphical models (PGM)**, that map a set of interconnected latent variables $\mathbf{z}_1, \dots, \mathbf{z}_L$ to the observed variables $\mathbf{x}_1, \dots, \mathbf{x}_D$ using simpler, often linear, mappings. Of course, many hybrids are possible. For example, PGMs can use neural networks, and DGMs can use structured state spaces. We discuss PGMs in general terms in Chapter 4, and give examples in Chapter 28, Chapter 29, Chapter 30. In this part of the book, we mostly focus on DGMs.

The main kinds of DGM are: **variational autoencoders (VAE)**, **autoregressive models (ARM)** models, **normalizing flows**, **diffusion models**, **energy based models (EBM)**, and **generative adversarial networks (GAN)**. We can categorize these models in terms of the following criteria (see Figure 20.1 for a visual summary):

- **Density**: does the model support pointwise evaluation of the probability density function $p(\mathbf{x})$, and if so, is this fast or slow, exact, approximate or a bound, etc? For **implicit models**, such as GANs, there is no well-defined density $p(\mathbf{x})$. For other models, we can only compute a lower bound on the density (VAEs), or an approximation to the density (EBMs, UPGMs).
- **Sampling**: does the model support generating new samples, $\mathbf{x} \sim p(\mathbf{x})$, and if so, is this fast or slow, exact or approximate? Directed PGMs, VAEs, and GANs all support fast sampling. However, undirected PGMs, EBMs, ARM, diffusion, and flows are slow for sampling.
- **Training**: what kind of method is used for parameter estimation? For some models (such as AR, flows and directed PGMs), we can perform exact maximum likelihood estimation (MLE), although

Subtypes
of DGM

Model	Chapter	Density	Sampling	Training	Latents	Architecture
PGM-D	Section 4.2	Exact, fast	Fast	MLE	Optional	Sparse DAG
PGM-U	Section 4.3	Approx, slow	Slow	MLE-A	Optional	Sparse graph
VAE	Chapter 21	LB, fast	Fast	MLE-LB	\mathbb{R}^L	Encoder-Decoder
ARM	Chapter 22	Exact, fast	Slow	MLE	None	Sequential
Flows	Chapter 23	Exact, slow/fast	Slow	MLE	\mathbb{R}^D	Invertible
EBM	Chapter 24	Approx, slow	Slow	MLE-A	Optional	Discriminative
Diffusion	Chapter 25	LB	Slow	MLE-LB	\mathbb{R}^D	Encoder-Decoder
GAN	Chapter 26	NA	Fast	Min-max	\mathbb{R}^L	Generator-Discriminator

Table 20.1: Characteristics of common kinds of generative model. Here D is the dimensionality of the observed \mathbf{x} , and L is the dimensionality of the latent \mathbf{z} , if present. (We usually assume $L \ll D$, although overcomplete representations can have $L \gg D$.) Abbreviations: Approx = approximate, ARM = autoregressive model, EBM = energy based model, GAN = generative adversarial network, MLE = maximum likelihood estimation, MLE-A = MLE (approximate), MLE-LB = MLE (lower bound), NA = not available, PGM = probabilistic graphical model, PGM-D = directed PGM, PGM-U = undirected PGM, VAE = variational autoencoder.

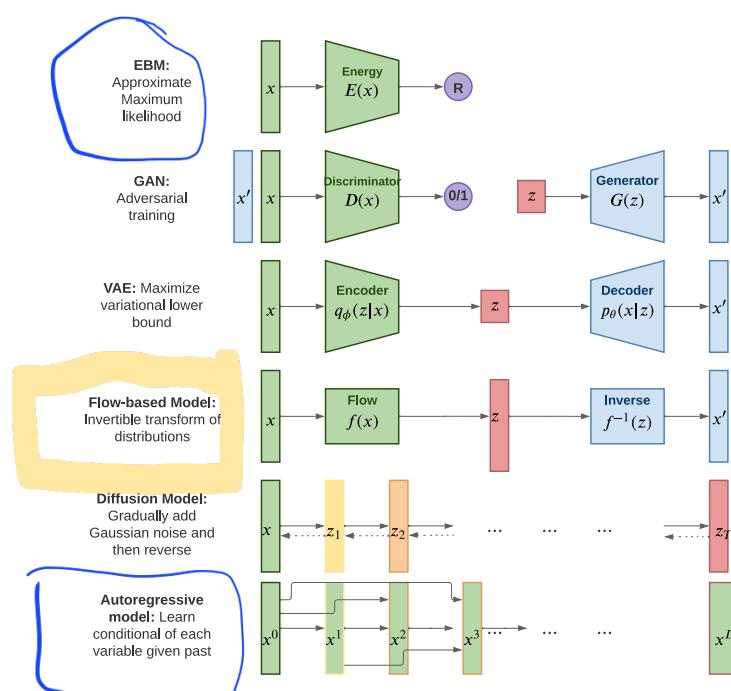


Figure 20.1: Summary of various kinds of deep generative model. Here \mathbf{x} is the observed data, \mathbf{z} is the latent code, and \mathbf{x}' is a sample from the model. AR models do not have a latent code \mathbf{z} . For diffusion models and flow models, the size of \mathbf{z} is the same as \mathbf{x} . For AR models, x^d is the d 'th dimension of \mathbf{x} . R represents real-valued output, 0/1 represents binary output. Adapted from Figure 1 of [Wen21].

the objective is usually non-convex, so we can only reach a local optimum. For other models, we cannot tractably compute the likelihood. In the case of VAEs, we maximize a lower bound on the likelihood; in the case of EBMs and UGMs, we maximize an approximation to the likelihood. For GANs we have to use min-max training, which can be unstable, and there is no clear objective function to monitor.

- Latents: does the model use a latent vector \mathbf{z} to generate \mathbf{x} or not, and if so, is it the same size as \mathbf{x} or is it a potentially compressed representation? For example, ARMs do not use latents; flows and diffusion use latents, but they are not compressed.¹ Graphical models, including EBMs, may or may not use latents.
- Architecture: what kind of neural network should we use, and are there restrictions? For flows, we are restricted to using invertible neural networks where each layer has a tractable Jacobian. For EBMs, we can use any model we like. The other models have different restrictions.

20.3 Goals of generative modeling

There are several different kinds of tasks that we can use generative models for, as we discuss below.

20.3.1 Generating data

One of the main goals of generative models is to generate (create) new data samples. This is sometimes called **generative AI** (see e.g., [GBGM23] for a recent survey). For example, if we fit a model $p(\mathbf{x})$ to images of faces, we can sample new faces from it, as illustrated in Figure 25.10.² Similar methods can be used to create samples of text, audio, etc. When this technology is abused to make fake content, they are called **deep fakes** (see e.g., [Ngu+19]). Generative models can also be used to create **synthetic data** for training discriminative models (see e.g., [Wil+20; Jor+22]).

To control what is generated, it is useful to use a **conditional generative model** of the form $p(\mathbf{x}|\mathbf{c})$. Here are some examples:

- \mathbf{c} = text prompt, \mathbf{x} = image. This is a **text-to-image** model (see Figure 20.2, Figure 20.3 and Figure 22.6 for examples).
- \mathbf{c} = image, \mathbf{x} = text. This is an **image-to-text** model, which is useful for **image captioning**.
- \mathbf{c} = image, \mathbf{x} = image. This is an **image-to-image** model, and can be used for image colorization, inpainting, uncropping, JPEG artefact restoration, etc. See Figure 20.4 for examples.
- \mathbf{c} = sequence of sounds, \mathbf{x} = sequence of words. This is a **speech-to-text** model, which is useful for **automatic speech recognition (ASR)**.
- \mathbf{c} = sequence of English words, \mathbf{x} = sequence of French words. This is a **sequence-to-sequence** model, which is useful for **machine translation**.

¹ Flow models define a latent vector \mathbf{z} that has the same size as \mathbf{x} , although the internal deterministic computation may use vectors that are larger or smaller than the input (see e.g., the DenseFlow paper [GGS21]).

² These images were made with a technique called score-based generative modeling (Section 25.3), although similar results can be obtained using many other techniques. See for example <https://this-person-does-not-exist.com/en> which shows results from a GAN model (Chapter 26).

1
2
3
4
5
6
7
8
9
10
11



12 (a) *Teddy bears swimming at the
13 Olympics 400m Butterfly event.*



12 (b) *A cute corgi lives in a house
13 made out of sushi.*

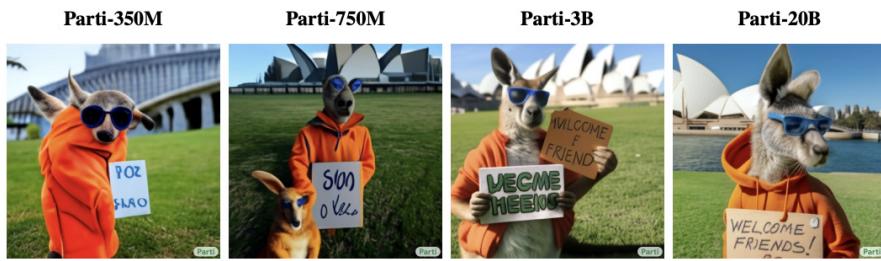


12 (c) *A cute sloth holding a small trea-
13 sure chest. A bright golden glow is
coming from the chest.*

14 *Figure 20.2: Some 1024×1024 images generated from text prompts by the Imagen diffusion model (Sec-
15 tion 25.6.4). From Figure 1 of [Sah+22b]. Used with kind permission of William Chan.*

16

17



26 A portrait photo of a kangaroo wearing an orange hoodie and blue sunglasses standing on the grass
27 in front of the Sydney Opera House holding a sign on the chest that says Welcome Friends!
28

29 *Figure 20.3: Some images generated from the Parti transformer model (Section 22.4.2) in response to a
30 text prompt. We show results from models of increasing size (350M, 750M, 3B, 20B). Multiple samples are
31 generated, and the highest ranked one is shown. From Figure 10 of [Yu+22]. Used with kind permission of
32 Jiahui Yu.*

33

34

35

36

- 37 • \mathbf{c} = initial prompt, \mathbf{x} = continuation of the text. This is another sequence-to-sequence model,
38 which is useful for automatic **text generation** (see Figure 22.5 for an example).

39

40 Note that, in the conditional case, we sometimes denote the inputs by \mathbf{x} and the outputs by \mathbf{y} . In
41 this case the model has the familiar form $p(\mathbf{y}|\mathbf{x})$. In the special case that \mathbf{y} denotes a low dimensional
42 quantity, such as a integer class label, $y \in \{1, \dots, C\}$, we get a predictive (discriminative) model.
43 The main difference between a discriminative model and a conditional generative model is this: in a
44 discriminative model, we assume there is one correct output, whereas in a conditional generative
45 model, we assume there may be multiple correct outputs. This makes it harder to evaluate generative
46 models, as we discuss in Section 20.4.

47

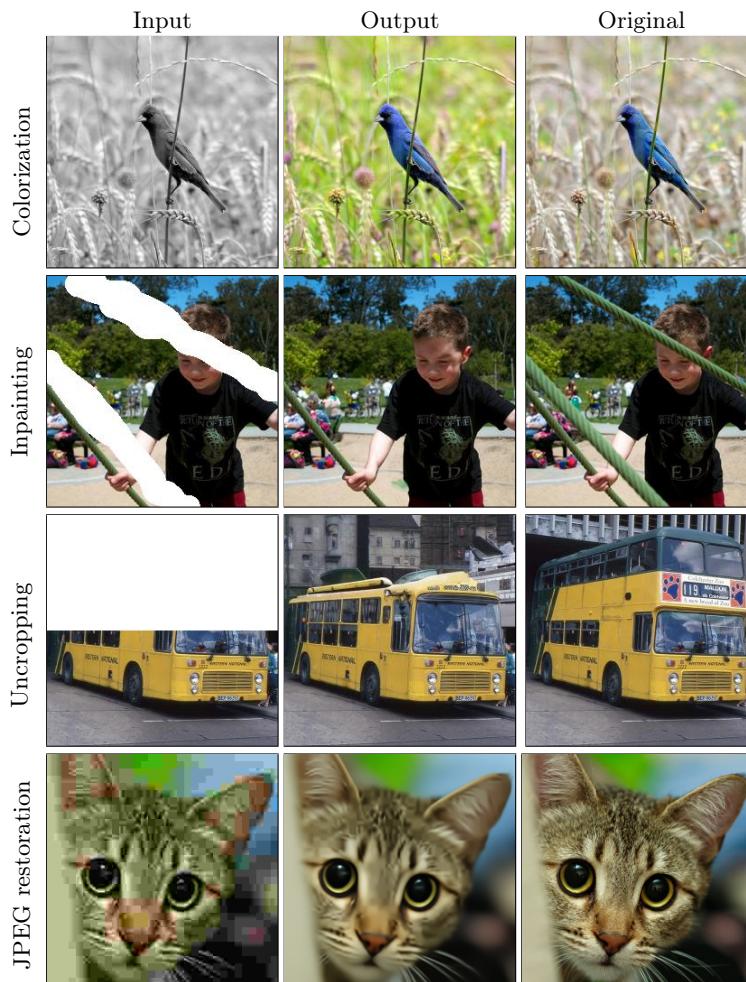
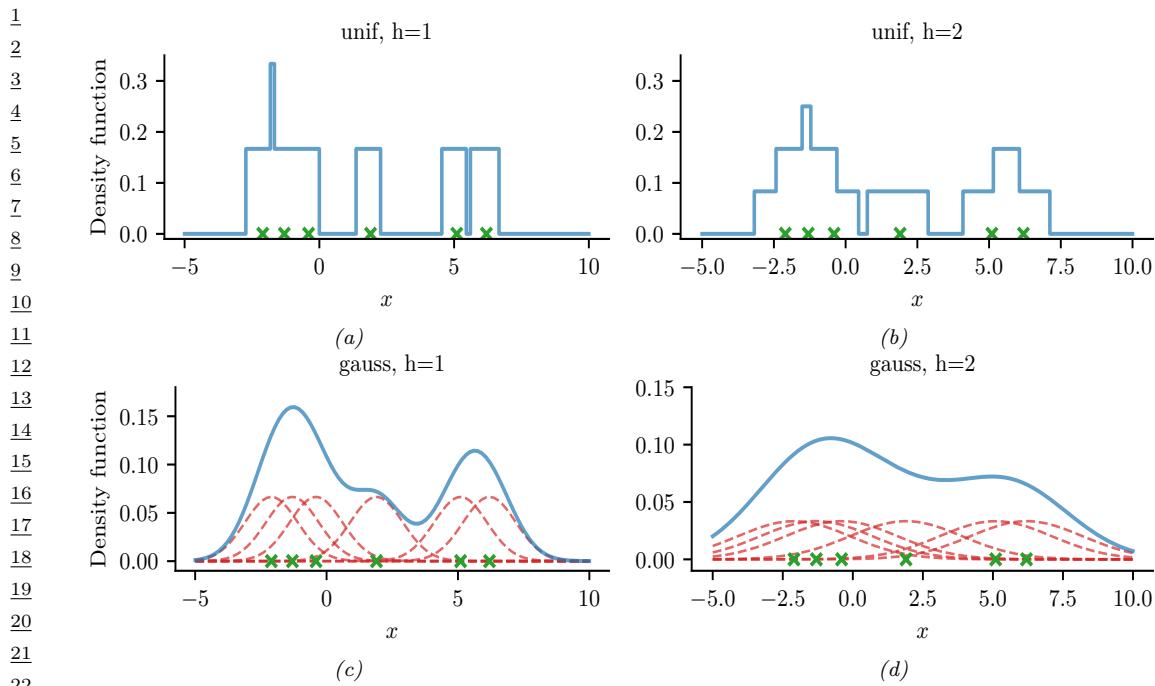


Figure 20.4: Illustration of some image-to-image tasks using the Palette conditional diffusion model (Section 25.6.4). From Figure 1 of [Sah+22a]. Used with kind permission of Chitwan Saharia.

20.3.2 Density estimation

The task of **density estimation** refers to evaluating the probability of an observed data vector, i.e., computing $p(\mathbf{x})$. This can be useful for outlier detection (Section 19.3.2), data compression (Section 5.4), generative classifiers, model comparison, etc.

A simple approach to this problem, which works in low dimensions, is to use **kernel density**



23 *Figure 20.5: A nonparametric (Parzen) density estimator in 1d estimated from 6 datapoints, denoted by x .*
 24 *Top row: uniform kernel. Bottom row: Gaussian kernel. Left column: bandwidth parameter $h = 1$. Right*
 25 *column: bandwidth parameter $h = 2$. Adapted from http://en.wikipedia.org/wiki/Kernel_density_estimation. Generated by `parzen_window_demo.ipynb`.*

29 **estimation** or **KDE**, which has the form

$$31 \quad p(\mathbf{x}|\mathcal{D}) = \frac{1}{N} \sum_{n=1}^N \mathcal{K}_h(\mathbf{x} - \mathbf{x}_n) \quad (20.1)$$

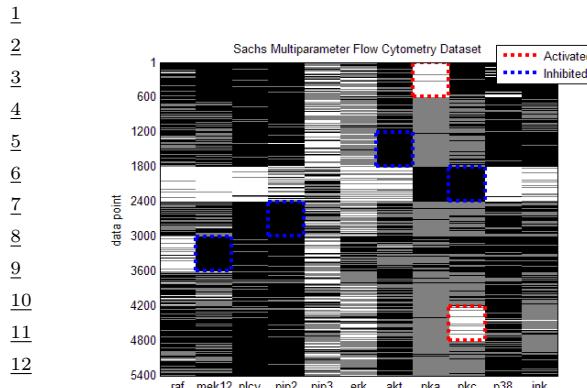
34 Here $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ is the data, and \mathcal{K}_h is a density kernel with **bandwidth** h , which is a
 35 function $\mathcal{K} : \mathbb{R} \rightarrow \mathbb{R}_+$ such that $\int \mathcal{K}(x)dx = 1$ and $\int x\mathcal{K}(x)dx = 0$. We give a 1d example of this in
 36 Figure 20.5: in the top row, we use a uniform (boxcar) kernel, and in the bottom row, we use a
 37 Gaussian kernel.

38 In higher dimensions, KDE suffers from the **curse of dimensionality** (see e.g., [AHK01]), and
 39 we need to use parametric density models $p_\theta(\mathbf{x})$ of some kind.

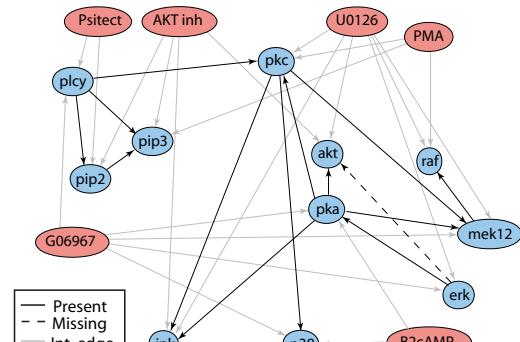
41 20.3.3 Imputation

43 The task of **imputation** refers to “filling in” missing values of a data vector or data matrix. For
 44 example, suppose \mathbf{X} is an $N \times D$ matrix of data (think of a spreadsheet) in which some entries, call
 45 them \mathbf{X}_m , may be missing, while the rest, \mathbf{X}_o , are observed. A simple way to fill in the missing
 46 data is to use the mean value of each feature, $\mathbb{E}[x_d]$; this is called **mean value imputation**, and is
 47

	Data sample	Variables			Missing values replaced by means						
		A	B	C	A	B	C				
1	1	6	6	NA	2	6	7.5				
2	2	NA	6	0	9	6	0				
3	3	NA	6	NA	9	6	7.5				
4	4	10	10	10	10	10	10				
5	5	10	10	10	10	10	10				
6	6	10	10	10	10	10	10				
7	Average		9	8	7.5						
8							9				
9							8				
10							7.5				
11											
12											
13											
14	Figure 20.6: Missing data imputation using the mean of each column.										
15											
16											
17											
18	illustrated in Figure 20.6. However, this ignores dependencies between the variables within each row,										
19	and does not return any measure of uncertainty.										
20	We can generalize this by fitting a generative model to the observed data, $p(\mathbf{X}_o)$, and then										
21	computing samples from $p(\mathbf{X}_m \mathbf{X}_o)$. This is called multiple imputation . A generative model can										
22	be used to fill in more complex data types, such as in-painting occluded pixels in an image (see										
23	Figure 20.4).										
24	See Section 3.11 for a more general discussion of missing data.										
25											
26	20.3.4 Structure discovery										
27											
28	Some kinds of generative models have latent variables \mathbf{z} , which are assumed to be the “causes”										
29	that generated the observed data \mathbf{x} . We can use Bayes’ rule to invert the model to compute										
30	$p(\mathbf{z} \mathbf{x}) \propto p(\mathbf{z})p(\mathbf{x} \mathbf{z})$. This can be useful for discovering latent, low-dimensional patterns in the data.										
31	For example, suppose we perturb various proteins in a cell and measure the resulting phosphorylation										
32	state using a technique known as flow cytometry, as in [Sac+05]. An example of such a dataset is										
33	shown in Figure 20.7(a). Each row represents a data sample $\mathbf{x}_n \sim p(\cdot \mathbf{a}_n, \mathbf{z})$, where $\mathbf{x} \in \mathbb{R}^{11}$ is a										
34	vector of outputs (phosphorylations), $\mathbf{a} \in \{0, 1\}^6$ is a vector of input actions (perturbations) and \mathbf{z} is										
35	the unknown cellular signaling network structure. We can infer the graph structure $p(\mathbf{z} \mathcal{D})$ using										
36	graphical model structure learning techniques (see Section 30.3). In particular, we can use the										
37	dynamic programming method described in [EM07] to get the result is shown in Figure 20.7(b). Here										
38	we plot the median graph, which includes all edges for which $p(z_{ij} = 1 \mathcal{D}) > 0.5$. (For a more recent										
39	approach to this problem, see e.g., [Bro+20b].)										
40											
41	20.3.5 Latent space interpolation										
42											
43	One of the most interesting abilities of certain latent variable models is the ability to generate samples										
44	that have certain desired properties by interpolating between existing datapoints in latent space.										
45	To explain how this works, let \mathbf{x}_1 and \mathbf{x}_2 be two inputs (e.g., images), and let $\mathbf{z}_1 = e(\mathbf{x}_1)$ and										
46	$\mathbf{z}_2 = e(\mathbf{x}_2)$ be their latent encodings. (The method used for computing these will depend on the										
47											

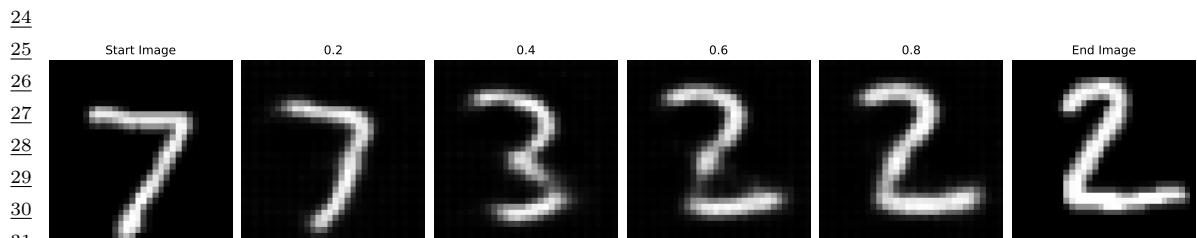


(a)



(b)

16 Figure 20.7: (a) A design matrix consisting of 5400 datapoints (rows) measuring the state (using flow
17 cytometry) of 11 proteins (columns) under different experimental conditions. The data has been discretized
18 into 3 states: low (black), medium (grey), and high (white). Some proteins were explicitly controlled using
19 activating or inhibiting chemicals. (b) A directed graphical model representing dependencies between various
20 proteins (blue circles) and various experimental interventions (pink ovals), which was inferred from this data.
21 We plot all edges for which $p(G_{ij} = 1 | \mathcal{D}) > 0.5$. Dotted edges are believed to exist in nature but were not
22 discovered by the algorithm (1 false negative). Solid edges are true positives. The light colored edges represent
23 the effects of intervention. From Figure 6d of [EM07].



32
33 Figure 20.8: Interpolation between two MNIST images in the latent space of a β -VAE (with $\beta = 0.5$).
34 Generated by [mnist_vae_ae_comparison.ipynb](#).



44 Figure 20.9: Interpolation between two CelebA images in the latent space of a β -VAE (with $\beta = 0.5$).
45 Generated by [celeba_vae_ae_comparison.ipynb](#).



Figure 20.10: Arithmetic in the latent space of a β -VAE (with $\beta = 0.5$). The first column is an input image, with embedding \mathbf{z} . Subsequent columns show the decoding of $\mathbf{z} + s\Delta$, where $s \in \{-2, -1, 0, 1, 2\}$ and $\Delta = \bar{\mathbf{z}}^+ - \bar{\mathbf{z}}^-$ is the difference in the average embeddings of images with or without a certain attribute (here, wearing sunglasses). Generated by `celeba_vae_ae_comparison.ipynb`.

type of model; we discuss the details in later chapters.) We can regard \mathbf{z}_1 and \mathbf{z}_2 as two “anchors” in latent space. We can now generate new images that interpolate between these points by computing $\mathbf{z} = \lambda\mathbf{z}_1 + (1 - \lambda)\mathbf{z}_2$, where $0 \leq \lambda \leq 1$, and then decoding by computing $\mathbf{x}' = d(\mathbf{z})$, where $d()$ is the decoder. This is called **latent space interpolation**, and will generate data that combines semantic features from both \mathbf{x}_1 and \mathbf{x}_2 . (The justification for taking a linear interpolation is that the learned manifold often has approximately zero curvature, as shown in [SKTF18]. However, sometimes it is better to use nonlinear interpolation [Whi16; MB21; Fad+20].)

We can see an example of this process in Figure 20.8, where we use a β -VAE model (Section 21.3.1) fit to the MNIST dataset. We see that the model is able to produce plausible interpolations between the digit 7 and the digit 2. As a more interesting example, we can fit a β -VAE to the **CelebA** dataset [Liu+15].³ The results are shown in Figure 20.9, and look reasonable. (We can get much better quality if we use a larger model trained on more data for a longer amount of time.)

It is also possible to perform interpolation in the latent space of text models, as illustrated in Figure 21.7.

20.3.6 Latent space arithmetic

In some cases, we can go beyond interpolation, and can perform **latent space arithmetic**, in which we can increase or decrease the amount of a desired “semantic factor of variation”. This was first shown in the **word2vec** model [Mik+13], but it also is possible in other latent variable models. For example, consider our VAE model fit to the CelebA dataset, which has faces of celebrities and some corresponding attributes. Let \mathbf{X}_i^+ be a set of images which have attribute i , and \mathbf{X}_i^- be a set of images which do not have this attribute. Let \mathbf{Z}_i^+ and \mathbf{Z}_i^- be the corresponding embeddings, and $\bar{\mathbf{z}}_i^+$ and $\bar{\mathbf{z}}_i^-$ be the average of these embeddings. We define the offset vector as $\Delta_i = \bar{\mathbf{z}}_i^+ - \bar{\mathbf{z}}_i^-$. If we add some positive multiple of Δ_i to a new point \mathbf{z} , we increase the amount of the attribute i ; if we subtract some multiple of Δ_i , we decrease the amount of the attribute i [Whi16].

We give an example of this in Figure 20.10. We consider the attribute of wearing sunglasses. The j 'th reconstruction is computed using $\hat{\mathbf{x}}_j = d(\mathbf{z} + s_j\Delta)$, where $\mathbf{z} = e(\mathbf{x})$ is the encoding of the original image, and s_j is a scale factor. When $s_j > 0$ we add sunglasses to the face. When $s_j < 0$ we

³3. CelebA contains about 200k images of famous celebrities. The images are also annotated with 40 attributes. We reduce the resolution of the images to 64×64 , as is conventional.

1 remove sunglasses; but this also has the side effect of making the face look younger and more female,
2 possibly a result of dataset bias.
3

4

5 **20.3.7 Generative design** 6

7 Another interesting use case for (deep) generative models is **generative design**, in which we use
8 the model to generate candidate objects, such as molecules, which have desired properties (see
9 e.g., [RNA22]). One approach is to fit a VAE to unlabeled samples, and then to perform Bayesian
10 optimization (Section 6.6) in its latent space, as discussed in Section 21.3.5.2.
11

12

13 **20.3.8 Model-based reinforcement learning**

14 We discuss reinforcement learning (RL) in Chapter 35. The main success stories of RL to date have
15 been in computer games, where simulators exist and data is abundant. However, in other areas, such
16 as robotics, data is expensive to acquire. In this case, it can be useful to learn a generative “**world**
17 **model**”, so the agent can do planning and learning “in its head”. See Section 35.4 for more details.
18

19

20 **20.3.9 Representation learning**

21 **Representation learning** refers to learning (possibly uninterpretable) latent factors \mathbf{z} that generate
22 the observed data \mathbf{x} . The primary goal is for these features to be used in “**downstream**” supervised
23 tasks. This is discussed in Chapter 32.
24

25

26 **20.3.10 Data compression**

27 Models which can assign high probability to frequently occurring data vectors (e.g., images, sentences),
28 and low probability to rare vectors, can be used for **data compression**, since we can assign shorter
29 codes to the more common items. Indeed, the optimal coding length for a vector \mathbf{x} from some
30 stochastic source $p(\mathbf{x})$ is $l(\mathbf{x}) = -\log p(\mathbf{x})$, as proved by Shannon. See Section 5.4 for details.
31

32

33 **20.4 Evaluating generative models**

34

35 *This section is written by Mihaela Rosca, Shakir Mohamed, and Balaji Lakshminarayanan.*

36 Evaluating generative models requires metrics which capture

37

38 • **sample quality** — are samples generated by the model a part of the data distribution?
39

40

41 • **sample diversity** — are samples from the model distribution capturing all modes of the data
distribution?
42

43

44 • **generalization** — is the model generalizing beyond the training data?
45

46

47 There is no known metric which meets all these requirements, but various metrics have been proposed
to capture different aspects of the learned distribution, some of which we discuss below.

20.4.1 Likelihood-based evaluation

A standard way to measure how close a model q is to a true distribution p is in terms of the KL divergence (Section 5.1):

$$D_{\text{KL}}(p \parallel q) = \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} = -\mathbb{H}(p) + \mathbb{H}_{ce}(p, q) \quad (20.2)$$

where $\mathbb{H}(p)$ is a constant, and $\mathbb{H}_{ce}(p, q)$ is the cross entropy. If we approximate $p(\mathbf{x})$ by the empirical distribution, we can evaluate the cross entropy in terms of the empirical **negative log likelihood** on the dataset:

$$\text{NLL} = -\frac{1}{N} \sum_{n=1}^N \log q(\mathbf{x}_n) \quad (20.3)$$

Usually we care about negative log likelihood on a held-out test set.⁴

20.4.1.1 Computing the log-likelihood

For models of discrete data, such as language models, it is easy to compute the (negative) log likelihood. However, it is common to measure performance using a quantity called **perplexity**, which is defined as 2^H , where $H = \text{NLL}$ is the cross entropy or negative log likelihood.

For image and audio models, one complication is that the model is usually a continuous distribution $p(\mathbf{x}) \geq 0$ but the data is usually discrete (e.g., $\mathbf{x} \in \{0, \dots, 255\}^D$ if we use one byte per pixel). Consequently the average log likelihood can be arbitrary large, since the pdf can be bigger than 1. To avoid this it is standard practice to use **uniform dequantization** [TOB16], in which we add uniform random noise to the discrete data, and then treat it as continuous-valued data. This gives a lower bound on the average log likelihood of the discrete model on the original data.

To see this, let \mathbf{z} be a continuous latent variable, and \mathbf{x} be a vector of binary observations computed by rounding, so $p(\mathbf{x}|\mathbf{z}) = \delta(\mathbf{x} - \text{round}(\mathbf{z}))$, computed elementwise. We have $p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$. Let $q(\mathbf{z}|\mathbf{x})$ be a probabilistic inverse of \mathbf{x} , that is, it has support only on values where $p(\mathbf{x}|\mathbf{z}) = 1$. In this case, Jensen's inequality gives

$$\log p(\mathbf{x}) \geq \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z}) + \log p(\mathbf{z}) - \log q(\mathbf{z}|\mathbf{x})] \quad (20.4)$$

$$= \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{z}) - \log q(\mathbf{z}|\mathbf{x})] \quad (20.5)$$

Thus if we model the density of $\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})$, which is a dequantized version of \mathbf{x} , we will get a lower bound on $p(\mathbf{x})$.

20.4.1.2 Likelihood can be hard to compute

Unfortunately, for many models, computing the likelihood can be computationally expensive, since it requires knowing the normalization constant of the probability model. One solution is to use variational inference (Chapter 10), which provides a way to efficiently compute lower (and sometimes

⁴ In some applications, we report **bits per dimension**, which is the log likelihood using log base 2, divided by the dimensionality of \mathbf{x} . To compute this metric, recall that $\log_2 L = \frac{\log_e L}{\log_e 2}$, and hence $bpd = \text{NLL} \log_e(2) \frac{1}{\|\mathbf{x}\|}$.

1 upper) bounds on the log likelihood. Another solution is to use annealed importance sampling
2 (Section 11.5.4.1), which provides a way to estimate the log likelihood using Monte Carlo sampling.
3 However, in the case of implicit generative models, such as GANs (Chapter 26), the likelihood is not
4 even defined, so we need to find evaluation metrics that do not rely on likelihood.
5

6 20.4.1.3 Likelihood is not related to sample quality

7 A more subtle concern with likelihood is that it is often uncorrelated with the perceptual quality of
8 the samples, at least for real-valued data, such as images and sound. In particular, a model can have
9 great log-likelihood but create poor samples and vice versa.
10

11 To see why a model can have good likelihoods but create bad samples, consider the following
12 argument from [TOB16]. Suppose q_0 is a density model for D -dimensional data \mathbf{x} which performs
13 arbitrarily well as judged by average log-likelihood, and suppose q_1 is a bad model, such as white
14 noise. Now consider samples generated from the mixture model
15

$$\underline{16} \quad q_2(\mathbf{x}) = 0.01q_0(\mathbf{x}) + 0.99q_1(\mathbf{x}) \quad (20.6)$$

17 Clearly 99% of the samples will be poor. However, the log-likelihood per pixel will hardly change
18 between q_2 and q_0 if D is large, since
19

$$\underline{21} \quad \log q_2(\mathbf{x}) = \log[0.01q_0(\mathbf{x}) + 0.99q_1(\mathbf{x})] \geq \log[0.01q_0(\mathbf{x})] = \log q_0(\mathbf{x}) - 2 \quad (20.7)$$

22 For high-dimensional data, $|\log q_0(\mathbf{x})| \sim D \gg 100$, so $\log q_2(\mathbf{x}) \approx \log q_0(\mathbf{x})$, and hence mixing in the
23 poor sampler does not significantly impact the log likelihood.
24

25 Now consider a case where the model has good samples but bad likelihoods. To achieve this,
26 suppose q is a GMM centered on the training images:
27

$$\underline{28} \quad q(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \mathcal{N}(\mathbf{x} | \mathbf{x}_n, \epsilon^2 \mathbf{I}) \quad (20.8)$$

29 If ϵ is small enough that the Gaussian noise is imperceptible, then samples from this model will look
30 good, since they correspond to the training set of real images. But this model will almost certainly
31 have poor likelihood on the test set due to overfitting. (In this case we say the model has effectively
32 just memorized the training set.)
33

34 20.4.2 Distances and divergences in feature space

35 Due to the challenges associated with comparing distributions in high dimensional spaces, and the
36 desire to compare distributions in a semantically meaningful way, it is common to use domain-specific
37 **perceptual distance metrics**, that measure how similar data vectors are to each other or to the
38 training data. However, most metrics used to evaluate generative models do not directly compare
39 raw data (e.g., pixels) but use a neural network to obtain features from the raw data and compare
40 the feature distribution obtained from model samples with the feature distribution obtained from
41 the dataset. The neural network used to obtain features can be trained solely for the purpose of
42 evaluation, or can be pretrained; a common choice is to use a pretrained classifier (see e.g., [Sal+16;
43 Heu+17b; Bin+18; Kyn+19; SSG18a]).
44

The **Inception score** [Sal+16] measures the average KL divergence between the marginal distribution of class labels obtained from the samples $p_{\theta}(y) = \int p_{\text{disc}}(y|\mathbf{x})p_{\theta}(\mathbf{x})d\mathbf{x}$ (where the integral is approximated by sampling images \mathbf{x} from a fixed dataset) and the distribution $p(y|\mathbf{x})$ induced by samples from the model, $\mathbf{x} \sim p_{\theta}(\mathbf{x})$. (The term comes from the “Inception” model [Sze+15b] that is often used to define $p_{\text{disc}}(y|\mathbf{x})$.) This leads to the following score:

$$\text{IS} = \exp [\mathbb{E}_{p_{\theta}(\mathbf{x})} D_{\text{KL}}(p_{\text{disc}}(Y|\mathbf{x}) \parallel p_{\theta}(Y))] \quad (20.9)$$

To understand this, let us rewrite the log score as follows:

$$\log(\text{IS}) = \mathbb{H}(p_{\theta}(Y)) - \mathbb{E}_{p_{\theta}(\mathbf{x})} [\mathbb{H}(p_{\text{disc}}(Y|\mathbf{x}))] \quad (20.10)$$

Thus we see that a high scoring model will be equally likely to generate samples from all classes, thus maximizing the entropy of $p_{\theta}(Y)$, while also ensuring that each individual sample is easy to classify, thus minimizing the entropy of $p_{\text{disc}}(Y|\mathbf{x})$.

The Inception score solely relies on class labels, and thus does not measure overfitting or sample diversity outside the predefined dataset classes. For example, a model which generates one perfect example per class would get a perfect Inception score, despite not capturing the variety of examples inside a class, as shown in Figure 20.11a. To address this drawback, the **Fréchet Inception distance** or **FID** score [Heu+17b] measures the Fréchet distance between two Gaussian distributions on sets of features of a pre-trained classifier. One Gaussian is obtained by passing model samples through a pretrained classifier, and the other by passing dataset samples through the same classifier. If we assume that the mean and covariance obtained from model features are μ_m and Σ_m and those from the data are μ_d and Σ_d , then the FID is

$$\text{FID} = \|\mu_m - \mu_d\|_2^2 + \text{tr}(\Sigma_d + \Sigma_m - 2(\Sigma_d \Sigma_m)^{1/2}) \quad (20.11)$$

Since it uses features instead of class logits, the Fréchet distance captures more than modes captured by class labels, as shown in Figure 20.11b. Unlike the Inception score, a lower score is better since we want the two distributions to be as close as possible.

Unfortunately, the Fréchet distance has been shown to have a high bias, with results varying widely based on the number of samples used to compute the score. To mitigate this issue, the **kernel Inception distance** has been introduced [Bin+18], which measures the squared MMD (Section 2.7.3) between the features obtained from the data and features obtained from model samples.

20.4.3 Precision and recall metrics

Since the FID only measures the distance between the data and model distributions, it is difficult to use it as a diagnostic tool: a bad (high) FID can indicate that the model is not able to generate high quality data, or that it puts too much mass around the data distribution, or that the model only captures a subset of the data (e.g., in Figure 26.6). Trying to disentangle between these two failure modes has been the motivation to seek individual precision (sample quality) and recall (sample diversity) metrics in the context of generative models [LPO17; Kyn+19]. (The diversity question is especially important in the context of GANs, where mode collapse (Section 26.3.3) can be an issue.)

A common approach is to use nearest neighbors in the feature space of a pretrained classifier to

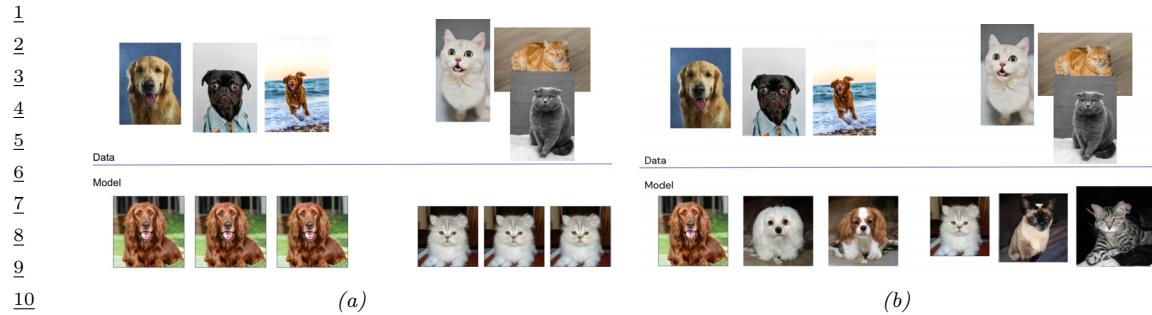


Figure 20.11: (a) Model samples with good (high) inception score are visually realistic. (b) Model samples with good (low) FID score are visually realistic and diverse.

define precision and recall [Kyn+19]. To formalize this, let us define

$$f_k(\phi, \Phi) = \begin{cases} 1 & \text{if } \exists \phi' \in \Phi \text{ s.t. } \|\phi - \phi'\|_2^2 \leq \|\phi' - \text{NN}_k(\phi', \Phi)\|_2^2 \\ 0 & \text{otherwise} \end{cases} \quad (20.12)$$

where Φ is a set of feature vectors and $\text{NN}_k(\phi', \Phi)$ is a function returning the k 'th nearest neighbor of ϕ' in Φ . We now define precision and recall as follows:

$$\text{precision}(\Phi_{model}, \Phi_{data}) = \frac{1}{|\Phi_{model}|} \sum_{\phi \in \Phi_{model}} f_k(\phi, \Phi_{data}); \quad (20.13)$$

$$\text{recall}(\Phi_{model}, \Phi_{data}) = \frac{1}{|\Phi_{data}|} \sum_{\phi \in \Phi_{data}} f_k(\phi, \Phi_{model}); \quad (20.14)$$

Precision and recall are always between 0 and 1. Intuitively, the precision metric measures whether samples are as close to data as data is to other data examples, while recall measures whether data is as close to model samples as model samples are to other samples. The parameter k controls how lenient the metrics will be — the higher k , the higher both precision and recall will be. As in classification, precision and recall in generative models can be used to construct a trade-off curve between different models which allows practitioners to make an informed decision regarding which model they want to use.

20.4.4 Statistical tests

Statistical tests have long been used to determine whether two sets of samples have been generated from the same distribution; these types of statistical tests are called **two sample tests**. Let us define the null hypothesis as the statement that both set of samples are from the same distribution. We then compute a statistic from the data and compare it to a threshold, and based on this we decide whether to reject the null hypothesis. In the context of evaluating implicit generative models such as GANs, statistics based on classifiers [Saj+18] and the MMD [Liu+20b] have been used. For use in scenarios with high dimensional input spaces, which are ubiquitous in the era of deep learning, two sample tests have been adapted to use learned features instead of raw data.

Like all other evaluation metrics for generative models, statistical tests have their own advantages and disadvantages: while users can specify Type 1 error — the chance they allow that the null hypothesis is wrongly rejected — statistical tests tend to be computationally expensive and thus cannot be used to monitor progress in training; hence they are best used to compare fully trained models.

20.4.5 Challenges with using pretrained classifiers

While popular and convenient, evaluation metrics that rely on pretrained classifiers (such as IS, FID, nearest neighbors in feature space, and statistical tests in feature space) have significant drawbacks. One might not have a pretrained classifier available for the dataset at hand, so classifiers trained on other datasets are used. Given the well known challenges with neural network generalization (see Section 17.4), the features of a classifier trained on images from one dataset might not be reliable enough to provide a fine grained signal of quality for samples obtained from a model trained on a different dataset. If the generative model is trained on the same dataset as the pre-trained classifier but the model is not capturing the data distribution perfectly, we are presenting the pre-trained classifier with out-of-distribution data and relying on its features to obtain score to evaluate our models. Far from being purely theoretical concerns, these issues have been studied extensively and have been shown to affect evaluation in practice [RV19; BS18].

20.4.6 Using model samples to train classifiers

Instead of using pretrained classifiers to evaluate samples, one can train a classifier on samples from conditional generative models, and then see how good these classifiers are at classifying data. For example, does adding synthetic (sampled) data to the real data help? This is closer to a reliable evaluation of generative model samples, since ultimately, the performance of generative models is dependent on the downstream task they are trained for. If used for semisupervised learning, one should assess how much adding samples to a classifier dataset helps with test accuracy. If used for model based reinforcement learning, one should assess how much the generative model helps with agent performance. For examples of this approach, see e.g., [SSM18; SSA18; RV19; SS20b; Jor+22].

20.4.7 Assessing overfitting

Many of the metrics discussed so far capture the sample quality and diversity, but do not capture overfitting to the training data. To capture overfitting, often a visual inspection is performed: a set of samples is generated from the model and for each sample its closest K nearest neighbors in the feature space of a pretrained classifier are obtained from the dataset. While this approach requires manually assessing samples, it is a simple way to test whether a model is simply memorizing the data. We show an example in Figure 20.12: since the model sample in the top left is quite different than its neighbors from the dataset (remaining images), we can conclude the sample is not simply memorised from the dataset. Similarly, sample diversity can be measured by approximating the support of the learned distribution by looking for similar samples in a large sample pool — as in the pigeonhole principle — but it is expensive and often requires manual human assessment[AZ17].

For likelihood-based models — such as variational autoencoders (Chapter 21), autoregressive models (Chapter 22), and normalizing flows (Chapter 23) — we can assess memorization by seeing



14 *Figure 20.12: Illustration of nearest neighbors in feature space: in the top left we have the query sample*
15 *generated using BigGAN, and the rest of the images are its nearest neighbors from the dataset. The nearest*
16 *neighbors search is done in the feature space of a pretrained classifier. From Figure 13 of [BDS18]. Used with*
17 *kind permission of Andy Brock.*

18

19

20 how much the log-likelihood of a model changes when a sample is included in the model's training
21 set or not [BW21].
22

23

24 20.4.8 Human evaluation

25 One approach to evaluate generative models is to use human evaluation, by presenting samples from
26 the model alongside samples from the data distribution, and ask human raters to compare the quality
27 of the samples [Zho+19b]. Human evaluation is a suitable metric if the model is used to create art or
28 other data for human display, or if reliable automated metrics are hard to obtain. However, human
29 evaluation can be difficult to standardize, hard to automate, and can be expensive or cumbersome to
30 set up.
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

21 Variational autoencoders

21.1 Introduction

In this chapter, we discuss generative models of the form

$$\mathbf{z} \sim p_{\theta}(\mathbf{z}) \tag{21.1}$$

$$\mathbf{x}|\mathbf{z} \sim \text{Expfam}(\mathbf{x}|d_{\theta}(\mathbf{z})) \tag{21.2}$$

where $p(\mathbf{z})$ is some kind of prior on the latent code \mathbf{z} , $d_{\theta}(\mathbf{z})$ is a deep neural network, known as the **decoder**, and $\text{Expfam}(\mathbf{x}|\boldsymbol{\eta})$ is an exponential family distribution, such as a Gaussian or product of Bernoullis. This is called a **deep latent variable model** or **DLVM**. When the prior is Gaussian (as is often the case), this model is called a **deep latent Gaussian model** or **DLGM**.

Posterior inference (i.e., computing $p_{\theta}(\mathbf{z}|\mathbf{x})$) is computationally intractable, as is computing the marginal likelihood

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p_{\theta}(\mathbf{z}) d\mathbf{z} \tag{21.3}$$

Hence we need to resort to approximate inference. For most of this chapter, we will use **amortized inference**, which we discussed in Section 10.1.5. This trains another model, $q_{\phi}(\mathbf{z}|\mathbf{x})$, called the **recognition network** or **inference network**, simultaneously with the generative model to do approximate posterior inference. This combination is called a **variational autoencoder** or **VAE** [KW14; RMW14b; KW19a], since it can be thought of as a probabilistic version of a deterministic autoencoder, discussed in Section 16.3.3.

In this chapter, we introduce the basic VAE, as well as some extensions. Note that the literature on VAE-like methods is vast¹, so we will only discuss a small subset of the ideas that have been explored.

21.2 VAE basics

In this section, we discuss the basics of variational autoencoders.

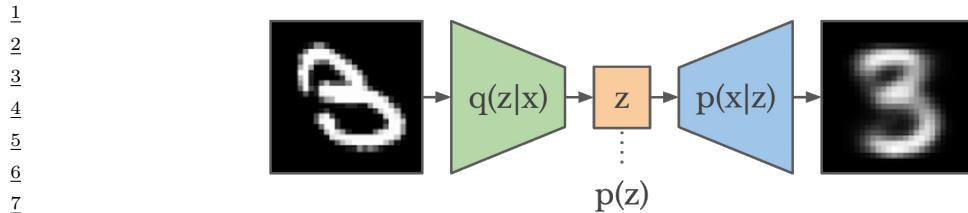


Figure 21.1: Schematic illustration of a VAE. From a figure in [Haf18]. Used with kind permission of Danijar Hafner.

21.2.1 Modeling assumptions

In the simplest setting, a VAE defines a generative model of the form

$$p_{\theta}(z, x) = p_{\theta}(z)p_{\theta}(x|z) \quad (21.4)$$

where $p_{\theta}(z)$ is usually a Gaussian, and $p_{\theta}(x|z)$ is usually a product of exponential family distributions (e.g., Gaussians or Bernoullis), with parameters computed by a neural network decoder, $d_{\theta}(z)$. For example, for binary observations, we can use

$$p_{\theta}(x|z) = \prod_{d=1}^D \text{Ber}(x_d | \sigma(d_{\theta}(z))) \quad (21.5)$$

In addition, a VAE fits a recognition model

$$q_{\phi}(z|x) = q(z|e_{\phi}(x)) \approx p_{\theta}(z|x) \quad (21.6)$$

to perform approximate posterior inference. Here $q_{\phi}(z|x)$ is usually a Gaussian, with parameters computed by a neural network encoder $e_{\phi}(x)$:

$$q_{\phi}(z|x) = \mathcal{N}(z|\mu, \text{diag}(\exp(\ell))) \quad (21.7)$$

$$(\mu, \ell) = e_{\phi}(x) \quad (21.8)$$

where $\ell = \log \sigma$. The model can be thought of as encoding the input x into a stochastic latent bottleneck z and then decoding it to approximately reconstruct the input, as shown in Figure 21.1.

The idea of training an inference network to “invert” a generative network, rather than running an optimization algorithm to infer the latent code, is called amortized inference, and is discussed in Section 10.1.5. This idea was first proposed in the **Helmholtz machine** [Day+95]. However, that paper did not present a single unified objective function for inference and generation, but instead used the wake-sleep (Section 10.6) method for training. By contrast, the VAE optimizes a variational lower bound on the log-likelihood, which means that convergence to a locally optimal MLE of the parameters is guaranteed.

We can use other approaches to fitting the DLGM (see e.g., [Hof17; DF19]). However, learning an inference network to fit the DLGM is often faster and can have some regularization benefits (see e.g., [KP20]).²

⁴⁵ 1. For example, the website <https://github.com/matthewvowels1/Awesome-VAEs> lists over 900 papers.

⁴⁶ 2. Combining a generative model with an inference model in this way results in what has been called a “monference”,

21.2.2 Model fitting

We can fit a VAE using amortized stochastic variational inference, as we discuss in Section 10.2.1.6. For example, suppose we use a VAE with a diagonal Bernoulli likelihood model, and a full covariance Gaussian as our variational posterior. Then we can use the methods discussed in Section 10.2.1.2 to derive the fitting algorithm. See Algorithm 21.1 for the corresponding pseudocode.

Algorithm 21.1: Fitting a VAE with Bernoulli likelihood and full covariance Gaussian posterior. Based on Algorithm 2 of [KW19a].

```

1 1 Initialize  $\theta, \phi$ 
2 2 repeat
3   3   Sample  $x \sim p_{\mathcal{D}}$ 
4   4   Sample  $\epsilon \sim q_0$ 
5   5    $(\mu, \log \sigma, L') = e_\phi(x)$ 
6   6    $M = np.triu(np.ones(K), -1)$ 
7   7    $L = M \odot L' + \text{diag}(\sigma)$ 
8   8    $z = L\epsilon + \mu$ 
9   9    $p_p = d_\theta(z)$ 
10  10   $\mathcal{L}_{\text{logqz}} = -\sum_{k=1}^K [\frac{1}{2}\epsilon_k^2 + \frac{1}{2}\log(2\pi) + \log \sigma_k]$  // from  $q_\phi(z|x)$  in Equation (10.47)
11  11   $\mathcal{L}_{\text{logpz}} = -\sum_{k=1}^K [\frac{1}{2}z_k^2 + \frac{1}{2}\log(2\pi)]$  // from  $p_\theta(z)$  in Equation (10.48)
12  12   $\mathcal{L}_{\text{logpx}} = -\sum_{d=1}^D [x_d \log p_d + (1-x_d) \log(1-p_d)]$  // from  $p_\theta(x|z)$ 
13  13   $\mathcal{L} = \mathcal{L}_{\text{logpx}} + \mathcal{L}_{\text{logpz}} - \mathcal{L}_{\text{logqz}}$ 
14  14  Update  $\theta := \theta - \eta \nabla_\theta \mathcal{L}$ 
15  15  Update  $\phi := \phi - \eta \nabla_\phi \mathcal{L}$ 
16  16 until converged

```

21.2.3 Comparison of VAEs and autoencoders

VAEs are very similar to deterministic autoencoders (AE). There are 2 main differences: in the AE, the objective is the log likelihood of the reconstruction without any KL term; and in addition, the encoding is deterministic, so the encoder network just needs to compute $\mathbb{E}[z|x]$ and not $\mathbb{V}[z|x]$. In view of these similarities, one can use the same codebase to implement both methods. However, it is natural to wonder what the benefits and potential drawbacks of the VAE are compared to the deterministic AE.

We shall answer this question by fitting both models to the CelebA dataset. Both models have the same convolutional structure with the following number of hidden channels per convolutional layer in the encoder: (32, 64, 128, 256, 512). The spatial size of each layer is as follows: (32, 16, 8, 4, 2). The final $2 \times 2 \times 512$ convolutional layer then gets reshaped and passed through a linear layer to generate the mean and (marginal) variance of the stochastic latent vector, which has size 256. The structure

i.e., model-inference hybrid. See the blog by Jacob Andreas, <http://blog.jacobandreas.net/monference.html>, for further discussion.

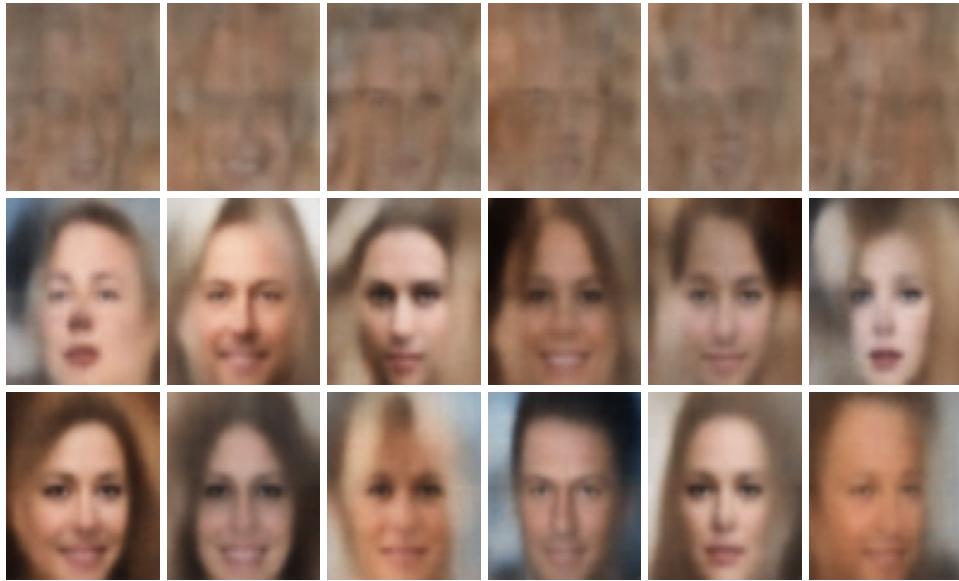


Figure 21.2: Illustration of unconditional image generation using (V)AEs trained on CelebA. Row 1: deterministic autoencoder. Row 2: β -VAE with $\beta = 0.5$. Row 3: VAE (with $\beta = 1$). Generated by celeba_vae_ae_comparison.ipynb.

of the decoder is the mirror image of the encoder. Each model is trained for 5 epochs with a batch size of 256, which takes about 20 minutes on a GPU.

The main advantage of a VAE over a deterministic autoencoder is that it defines a proper generative model, that can create sensible-looking novel images by decoding prior samples $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. By contrast, an autoencoder only knows how to decode latent codes derived from the training set, so does poorly when fed random inputs. This is illustrated in Figure 21.2.

We can also use both models to reconstruct a given input image. In Figure 21.3, we see that both AE and VAE can reconstruct the input images reasonably well, although the VAE reconstructions are somewhat blurry, for reasons we discuss in Section 21.3.1. We can reduce the amount of blurriness by scaling down the KL penalty term by a factor of β ; this is known as the β -VAE, and is discussed in more detail in Section 21.3.1.

21.2.4 VAEs optimize in an augmented space

In this section, we derive several alternative expressions for the ELBO which shed light on how VAEs work.

First, let us define the joint generative distribution

$$p_{\theta}(\mathbf{x}, \mathbf{z}) = p_{\theta}(\mathbf{z})p_{\theta}(\mathbf{x}|\mathbf{z}) \quad (21.9)$$



Figure 21.3: Illustration of image reconstruction using (V)AEs trained and applied to CelebA. Row 1: original images. Row 2: deterministic autoencoder. Row 3: β -VAE with $\beta = 0.5$. Row 4: VAE (with $\beta = 1$). Generated by [celeba_vae_ae_comparison.ipynb](#).

from which we can derive the generative data marginal

$$p_{\theta}(\mathbf{x}) = \int_{\mathbf{z}} p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} \quad (21.10)$$

and the generative posterior

$$p_{\theta}(\mathbf{z}|\mathbf{x}) = p_{\theta}(\mathbf{x}, \mathbf{z}) / p_{\theta}(\mathbf{x}) \quad (21.11)$$

Let us also define the joint *inference* distribution

$$q_{\mathcal{D}, \phi}(\mathbf{z}, \mathbf{x}) = p_{\mathcal{D}}(\mathbf{x}) q_{\phi}(\mathbf{z}|\mathbf{x}) \quad (21.12)$$

where

$$p_{\mathcal{D}}(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \delta(\mathbf{x}_n - \mathbf{x}) \quad (21.13)$$

1 is the empirical distribution. From this we can derive the inference latent marginal, also called the
2 aggregated posterior:

3

$$\underline{q}_{\mathcal{D}, \phi}(\mathbf{z}) = \int_{\mathbf{x}} q_{\mathcal{D}, \phi}(\mathbf{x}, \mathbf{z}) d\mathbf{x} \quad (21.14)$$

4 and the inference likelihood
5

6

$$\underline{q}_{\mathcal{D}, \phi}(\mathbf{x}|\mathbf{z}) = q_{\mathcal{D}, \phi}(\mathbf{x}, \mathbf{z}) / q_{\mathcal{D}, \phi}(\mathbf{z}) \quad (21.15)$$

7 See Figure 21.4 for a visual illustration.

8 Having defined our terms, we can now derive various alternative versions of the ELBO, following
9 [ZSE19]. First note that the ELBO averaged over all the data is given by

10

$$\underline{L}(\boldsymbol{\theta}, \phi | \mathcal{D}) = \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})]] - \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \| p_{\boldsymbol{\theta}}(\mathbf{z}))] \quad (21.16)$$

11

$$= \mathbb{E}_{q_{\mathcal{D}, \phi}(\mathbf{x}, \mathbf{z})} [\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z}) + \log p_{\boldsymbol{\theta}}(\mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})] \quad (21.17)$$

12

$$= \mathbb{E}_{q_{\mathcal{D}, \phi}(\mathbf{x}, \mathbf{z})} \left[\log \frac{p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})}{q_{\mathcal{D}, \phi}(\mathbf{x}, \mathbf{z})} + \log p_{\mathcal{D}}(\mathbf{x}) \right] \quad (21.18)$$

13

$$= -D_{\text{KL}}(q_{\mathcal{D}, \phi}(\mathbf{x}, \mathbf{z}) \| p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})) + \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [\log p_{\mathcal{D}}(\mathbf{x})] \quad (21.19)$$

14 If we define $\stackrel{c}{=}$ to mean equal up to additive constants, we can rewrite the above as

15

$$\underline{L}(\boldsymbol{\theta}, \phi | \mathcal{D}) \stackrel{c}{=} -D_{\text{KL}}(q_{\phi}(\mathbf{x}, \mathbf{z}) \| p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})) \quad (21.20)$$

16

$$\stackrel{c}{=} -D_{\text{KL}}(p_{\mathcal{D}}(\mathbf{x}) \| p_{\boldsymbol{\theta}}(\mathbf{x})) - \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \| p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}))] \quad (21.21)$$

17 Thus maximizing the ELBO requires minimizing the two KL terms. The first KL term is minimized
18 by MLE, and the second KL term is minimized by fitting the true posterior. Thus if the posterior
19 family is limited, there may be a conflict between these objectives.

20 Finally, we note that the ELBO can also be written as

21

$$\underline{L}(\boldsymbol{\theta}, \phi | \mathcal{D}) \stackrel{c}{=} -D_{\text{KL}}(q_{\mathcal{D}, \phi}(\mathbf{z}) \| p_{\boldsymbol{\theta}}(\mathbf{z})) - \mathbb{E}_{q_{\mathcal{D}, \phi}(\mathbf{z})} [D_{\text{KL}}(q_{\phi}(\mathbf{x}|\mathbf{z}) \| p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z}))] \quad (21.22)$$

22 We see from Equation (21.22) that VAEs are trying to minimize the difference between the inference
23 marginal and generative prior, $D_{\text{KL}}(q_{\phi}(\mathbf{z}) \| p_{\boldsymbol{\theta}}(\mathbf{z}))$, while simultaneously minimizing reconstruction
24 error, $D_{\text{KL}}(q_{\phi}(\mathbf{x}|\mathbf{z}) \| p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z}))$. Since \mathbf{x} is typically of much higher dimensionality than \mathbf{z} , the latter
25 term usually dominates. Consequently, if there is a conflict between these two objectives (e.g., due to
26 limited modeling power), the VAE will favor reconstruction accuracy over posterior inference. Thus
27 the learned posterior may not be a very good approximation to the true posterior (see [ZSE19] for
28 further discussion).

29

30 21.3 VAE generalizations

31 In this section, we discuss some variants of the basic VAE model.

32

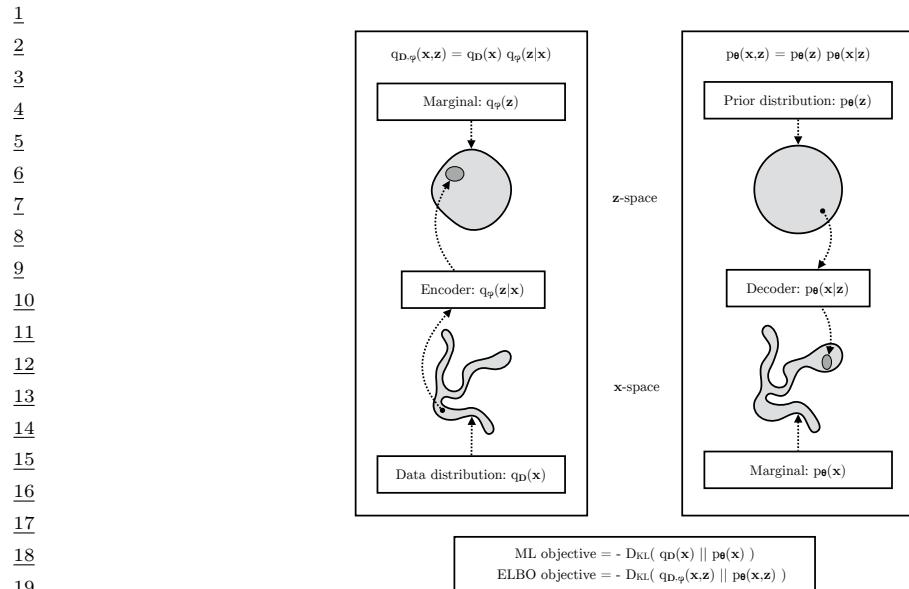


Figure 21.4: The maximum likelihood (ML) objective can be viewed as the minimization of $D_{KL}(p_D(x) || p_\theta(x))$. (Note: in the figure, $p_D(x)$ is denoted by $q_D(x)$.) The ELBO objective is minimization of $D_{KL}(q_{D,\phi}(x,z) || p_\theta(x,z))$, which upper bounds $D_{KL}(q_D(x) || p_\theta(x))$. From Figure 2.4 of [KW19a]. Used with kind permission of Durk Kingma.

21.3.1 β -VAE

It is often the case that VAEs generate somewhat blurry images, as illustrated in Figure 21.3, Figure 21.2 and Figure 20.9. This is not the case for models that optimize the exact likelihood, such as pixelCNNs (Section 22.3.2) and flow models (Chapter 23). To see why VAEs are different, consider the common case where the decoder is a Gaussian with fixed variance, so

$$\log p_\theta(x|z) = -\frac{1}{2\sigma^2} \|x - d_\theta(z)\|_2^2 + \text{const} \quad (21.23)$$

Let $e_\phi(x) = \mathbb{E}[q_\phi(z|x)]$ be the encoding of x , and $\mathcal{X}(z) = \{x : e_\phi(x) = z\}$ be the set of inputs that get mapped to z . For a fixed inference network, the optimal setting of the generator parameters, when using squared reconstruction loss, is to ensure $d_\theta(z) = \mathbb{E}[x : x \in \mathcal{X}(z)]$. Thus the decoder should predict the average of all inputs x that map to that z , resulting in blurry images.

We can solve this problem by increasing the expressive power of the posterior approximation (avoiding the merging of distinct inputs into the same latent code), or of the generator (by adding back information that is missing from the latent code), or both. However, an even simpler solution is to reduce the penalty on the KL term, making the model closer to a deterministic autoencoder:

$$\mathcal{L}_\beta(\theta, \phi|x) = \underbrace{-\mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x|z)]}_{\mathcal{L}_E} + \beta \underbrace{D_{KL}(q_\phi(z|x) || p_\theta(z))}_{\mathcal{L}_R} \quad (21.24)$$

where \mathcal{L}_E is the reconstruction error (negative log likelihood), and \mathcal{L}_R is the KL regularizer. This is

1 called the β -VAE objective [Hig+17a]. If we set $\beta = 1$, we recover the objective used in standard
2 VAEs; if we set $\beta = 0$, we recover the objective used in standard autoencoders.
3

4 By varying β from 0 to infinity, we can reach different points on the **rate distortion curve**, as
5 discussed in Section 5.4.2. These points make different tradeoffs between reconstruction error (distortion)
6 and how much information is stored in the latents about the input (rate of the corresponding
7 code). By using $\beta < 1$, we store more bits about each input, and hence can reconstruct images in a
8 less blurry way. If we use $\beta > 1$, we get a more compressed representation.
9

10 21.3.1.1 Disentangled representations

11 One advantage of using $\beta > 1$ is that it encourages the learning of a latent representation that is
12 “**disentangled**”. Intuitively this means that each latent dimension represents a different **factor of**
13 **variation** in the input. This is often formalized in terms of the total correlation (Section 5.3.5.1),
14 which is defined as follows:
15

$$\text{TC}(\mathbf{z}) = \sum_k \mathbb{H}(z_k) - \mathbb{H}(\mathbf{z}) = D_{\text{KL}} \left(p(\mathbf{z}) \parallel \prod_k p_k(z_k) \right) \quad (21.25)$$

16 This is zero iff the components of \mathbf{z} are all mutually independent, and hence disentangled. In [AS18],
17 they prove that using $\beta > 1$ will decrease the TC.
18

19 Unfortunately, in [Loc+18] they prove that nonlinear latent variable models are unidentifiable, and
20 therefore for any disentangled representation, there is an equivalent fully entangled representation
21 with exactly the same likelihood. Thus it is not possible to recover the correct latent representation
22 without choosing the appropriate inductive bias, via the encoder, decoder, prior, dataset, or learning
23 algorithm, i.e., merely adjusting β is not sufficient. See Section 32.4.1 for more discussion.
24

25

26 21.3.1.2 Connection with information bottleneck

27 In this section, we show that the β -VAE is an unsupervised version of the information bottleneck
28 (IB) objective from Section 5.6. If the input is \mathbf{x} , the hidden bottleneck is \mathbf{z} , and the target outputs
29 are $\tilde{\mathbf{x}}$, then the unsupervised IB objective becomes
30

$$\mathcal{L}_{\text{UIB}} = \beta \mathbb{I}(\mathbf{z}; \mathbf{x}) - \mathbb{I}(\mathbf{z}; \tilde{\mathbf{x}}) \quad (21.26)$$

$$= \beta \mathbb{E}_{p(\mathbf{x}, \mathbf{z})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{x})p(\mathbf{z})} \right] - \mathbb{E}_{p(\mathbf{z}, \tilde{\mathbf{x}})} \left[\log \frac{p(\mathbf{z}, \tilde{\mathbf{x}})}{p(\mathbf{z})p(\tilde{\mathbf{x}})} \right] \quad (21.27)$$

31 where

$$p(\mathbf{x}, \mathbf{z}) = p_{\mathcal{D}}(\mathbf{x})p(\mathbf{z}|\mathbf{x}) \quad (21.28)$$

$$p(\mathbf{z}, \tilde{\mathbf{x}}) = \int p_{\mathcal{D}}(\mathbf{x})p(\mathbf{z}|\mathbf{x})p(\tilde{\mathbf{x}}|\mathbf{z})d\mathbf{x} \quad (21.29)$$

32 Intuitively, the objective in Equation (21.26) means we should pick a representation \mathbf{z} that can
33 predict $\tilde{\mathbf{x}}$ reliably, while not memorizing too much information about the input \mathbf{x} . The tradeoff
34 parameter is controlled by β .
35

36

From Equation (5.181), we have the following variational upper bound on this unsupervised objective:

$$\mathcal{L}_{\text{UVIB}} = -\mathbb{E}_{q_{\mathcal{D}, \phi}(\mathbf{z}, \mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] + \beta \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p_{\theta}(\mathbf{z}))] \quad (21.30)$$

which matches Equation (21.24) when averaged over \mathbf{x} .

21.3.2 InfoVAE

In Section 21.2.4, we discussed some drawbacks of the standard ELBO objective for training VAEs, namely the tendency to ignore the latent code when the decoder is powerful (Section 21.4), and the tendency to learn a poor posterior approximation due to the mismatch between the KL terms in data space and latent space (Section 21.2.4). We can fix these problems to some degree by using a generalized objective of the following form:

$$\mathcal{L}(\theta, \phi|\mathbf{x}) = -\lambda D_{\text{KL}}(q_{\phi}(\mathbf{z}) \parallel p_{\theta}(\mathbf{z})) - \mathbb{E}_{q_{\phi}(\mathbf{z})} [D_{\text{KL}}(q_{\phi}(\mathbf{x}|\mathbf{z}) \parallel p_{\theta}(\mathbf{x}|\mathbf{z}))] + \alpha \mathbb{I}_q(\mathbf{x}; \mathbf{z}) \quad (21.31)$$

where $\alpha \geq 0$ controls how much we weight the mutual information $\mathbb{I}_q(\mathbf{x}; \mathbf{z})$ between \mathbf{x} and \mathbf{z} , and $\lambda \geq 0$ controls the tradeoff between \mathbf{z} -space KL and \mathbf{x} -space KL. This is called the **InfoVAE** objective [ZSE19]. If we set $\alpha = 0$ and $\lambda = 1$, we recover the standard ELBO, as shown in Equation (21.22).

Unfortunately, the objective in Equation (21.31) cannot be computed as written, because of the intractable MI term:

$$\mathbb{I}_q(\mathbf{x}; \mathbf{z}) = \mathbb{E}_{q_{\phi}(\mathbf{x}, \mathbf{z})} \left[\log \frac{q_{\phi}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{x})q_{\phi}(\mathbf{z})} \right] = -\mathbb{E}_{q_{\phi}(\mathbf{x}, \mathbf{z})} \left[\log \frac{q_{\phi}(\mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \quad (21.32)$$

However, using the fact that $q_{\phi}(\mathbf{x}|\mathbf{z}) = p_{\mathcal{D}}(\mathbf{x})q_{\phi}(\mathbf{z}|\mathbf{x})/q_{\phi}(\mathbf{z})$, we can rewrite the objective as follows:

$$\mathcal{L} = \mathbb{E}_{q_{\phi}(\mathbf{x}, \mathbf{z})} \left[-\lambda \log \frac{q_{\phi}(\mathbf{z})}{p_{\theta}(\mathbf{z})} - \log \frac{q_{\phi}(\mathbf{x}|\mathbf{z})}{p_{\theta}(\mathbf{x}|\mathbf{z})} - \alpha \log \frac{q_{\phi}(\mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \quad (21.33)$$

$$= \mathbb{E}_{q_{\phi}(\mathbf{x}, \mathbf{z})} \left[\log p_{\theta}(\mathbf{x}|\mathbf{z}) - \log \frac{q_{\phi}(\mathbf{z})^{\lambda+\alpha-1} p_{\mathcal{D}}(\mathbf{x})}{p_{\theta}(\mathbf{z})^{\lambda} q_{\phi}(\mathbf{z}|\mathbf{x})^{\alpha-1}} \right] \quad (21.34)$$

$$= \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]] - (1-\alpha) \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p_{\theta}(\mathbf{z}))] \\ - (\alpha + \lambda - 1) D_{\text{KL}}(q_{\phi}(\mathbf{z}) \parallel p_{\theta}(\mathbf{z})) - \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [\log p_{\mathcal{D}}(\mathbf{x})] \quad (21.35)$$

where the last term is a constant we can ignore. The first two terms can be optimized using the reparameterization trick. Unfortunately, the last term requires computing $q_{\phi}(\mathbf{z}) = \int_{\mathbf{x}} q_{\phi}(\mathbf{x}, \mathbf{z}) d\mathbf{x}$, which is intractable. Fortunately, we can easily sample from this distribution, by sampling $\mathbf{x} \sim p_{\mathcal{D}}(\mathbf{x})$ and $\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})$. Thus $q_{\phi}(\mathbf{z})$ is an **implicit probability model**, similar to a GAN (see Chapter 26).

As long as we use a strict divergence, meaning $D(q, p) = 0$ iff $q = p$, then one can show that this does not affect the optimality of the procedure. In particular, proposition 2 of [ZSE19] tells us the following:

Theorem 1. Let \mathcal{X} and \mathcal{Z} be continuous spaces, and $\alpha < 1$ (to bound the MI) and $\lambda > 0$. For any fixed value of $\mathbb{I}_q(\mathbf{x}; \mathbf{z})$, the approximate InfoVAE loss, with any strict divergence $D(q_{\phi}(\mathbf{z}), p_{\theta}(\mathbf{z}))$, is globally optimized if $p_{\theta}(\mathbf{x}) = p_{\mathcal{D}}(\mathbf{x})$ and $q_{\phi}(\mathbf{z}|\mathbf{x}) = p_{\theta}(\mathbf{z}|\mathbf{x})$.

1 **21.3.2.1 Connection with MMD VAE**

2 If we set $\alpha = 1$, the InfoVAE objective simplifies to

3

$$\underline{4} \quad \underline{\mathcal{L}} \stackrel{c}{=} \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]] - \lambda D_{\text{KL}}(q_{\phi}(\mathbf{z}) \parallel p_{\theta}(\mathbf{z})) \quad (21.36)$$

5

6 The **MMD VAE**³ replaces the KL divergence in the above term with the (squared) maximum mean
7 discrepancy or **MMD** divergence defined in Section 2.7.3. (This is valid based on the above theorem.)
8 The advantage of this approach over standard InfoVAE is that the resulting objective is tractable. In
9 particular, if we set $\lambda = 1$ and swap the sign we get

10

$$\underline{11} \quad \underline{\mathcal{L}} = \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [-\log p_{\theta}(\mathbf{x}|\mathbf{z})]] + \text{MMD}(q_{\phi}(\mathbf{z}), p_{\theta}(\mathbf{z})) \quad (21.37)$$

12

13 As we discuss in Section 2.7.3, we can compute the MMD as follows:

14

$$\underline{15} \quad \underline{\text{MMD}}(p, q) = \mathbb{E}_{p(\mathbf{z}), p(\mathbf{z}')} [\mathcal{K}(\mathbf{z}, \mathbf{z}')] + \mathbb{E}_{q(\mathbf{z}), q(\mathbf{z}')} [\mathcal{K}(\mathbf{z}, \mathbf{z}')] - 2\mathbb{E}_{p(\mathbf{z}), q(\mathbf{z}')} [\mathcal{K}(\mathbf{z}, \mathbf{z}')] \quad (21.38)$$

16

17 where $\mathcal{K}()$ is some kernel function, such as the RBF kernel, $\mathcal{K}(\mathbf{z}, \mathbf{z}') = \exp(-\frac{1}{2\sigma^2} \|\mathbf{z} - \mathbf{z}'\|_2^2)$. Intuitively
18 the MMD measures the similarity (in latent space) between samples from the prior and samples from
19 the aggregated posterior.

20 In practice, we can implement the MMD objective by using the posterior predicted mean $\mathbf{z}_n =$
21 $e_{\phi}(\mathbf{x}_n)$ for all B samples in the current minibatch, and comparing this to B random samples from
22 the $\mathcal{N}(\mathbf{0}, \mathbf{I})$ prior.

23 If we use a Gaussian decoder with fixed variance, the negative log likelihood is just a squared error
24 term:

25

$$\underline{26} \quad \underline{-\log p_{\theta}(\mathbf{x}|\mathbf{z})} = \|\mathbf{x} - d_{\theta}(\mathbf{z})\|_2^2 \quad (21.39)$$

27

28 Thus the entire model is deterministic, and just predicts the means in latent space and visible space.

29 **31 21.3.2.2 Connection with β -VAEs**

30 If we set $\alpha = 0$ and $\lambda = 1$, we get back the original ELBO. If $\lambda > 0$ is freely chosen, but we use
31 $\alpha = 1 - \lambda$, we get the β -VAE.

32 **36 21.3.2.3 Connection with adversarial autoencoders**

33 If we set $\alpha = 1$ and $\lambda = 1$, and D is chosen to be the Jensen-Shannon divergence (which can be
34 minimized by training a binary discriminator, as explained in Section 26.2.2), then we get a model
35 known as an **adversarial autoencoder** [Mak+15a].

36

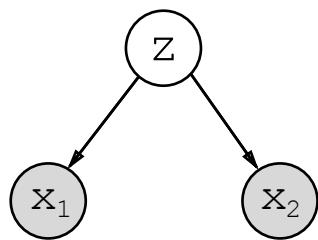
37 **41 21.3.3 Multimodal VAEs**

38 It is possible to extend VAEs to create joint distributions over different kinds of variables, such as
39 images and text. This is sometimes called a **multimodal VAE** or **MVAE**. Let us assume there are

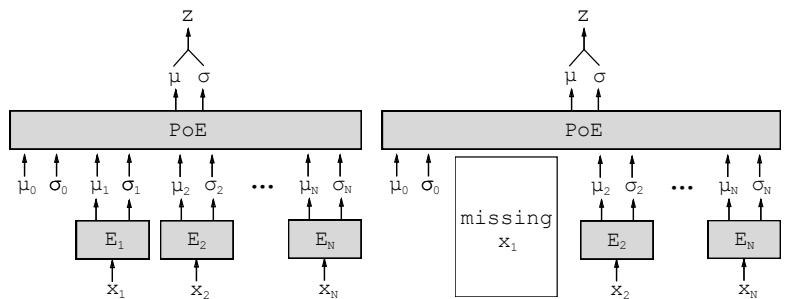
40

41 3. Proposed in <https://ermongroup.github.io/blog/a-tutorial-on-mmd-variational-autoencoders/>.

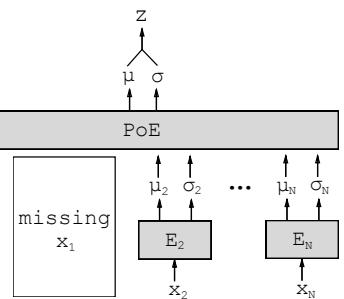
42



(a)



(b)



(c)

Figure 21.5: Illustration of multi-modal VAE. (a) The generative model with $N = 2$ modalities. (b) The product of experts (PoE) inference network is derived from N individual Gaussian experts E_i . μ_0 and σ_0 are parameters of the prior. (c) If a modality is missing, we omit its contribution to the posterior. From Figure 1 of [WG18]. Used with kind permission of Mike Wu.

M modalities. We assume they are conditionally independent given the latent code, and hence the generative model has the form

$$p_{\theta}(\mathbf{x}_1, \dots, \mathbf{x}_M, \mathbf{z}) = p(\mathbf{z}) \prod_{m=1}^M p_{\theta}(\mathbf{x}_m | \mathbf{z}) \quad (21.40)$$

where we treat $p(\mathbf{z})$ as a fixed prior. See Figure 21.5(a) for an illustration.

The standard ELBO is given by

$$\mathcal{L}(\theta, \phi | \mathbf{X}) = \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{X})} \left[\sum_m \log p_{\theta}(\mathbf{x}_m | \mathbf{z}) \right] - D_{\text{KL}}(q_{\phi}(\mathbf{z} | \mathbf{X}) \| p(\mathbf{z})) \quad (21.41)$$

where $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_M)$ is the observed data. However, the different likelihood terms $p(\mathbf{x}_m | \mathbf{z})$ may have different dynamic ranges (e.g., Gaussian pdf for pixels, and categorical pmf for text), so we introduce weight terms $\lambda_m \geq 0$ for each likelihood. In addition, let $\beta \geq 0$ control the amount of KL regularization. This gives us a weighted version of the ELBO, as follows:

$$\mathcal{L}(\theta, \phi | \mathbf{X}) = \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{X})} \left[\sum_m \lambda_m \log p_{\theta}(\mathbf{x}_m | \mathbf{z}) \right] - \beta D_{\text{KL}}(q_{\phi}(\mathbf{z} | \mathbf{X}) \| p(\mathbf{z})) \quad (21.42)$$

Often we don't have a lot of paired (aligned) data from all M modalities. For example, we may have a lot of images (modality 1), and a lot of text (modality 2), but very few (image, text) pairs. So it is useful to generalize the loss so it fits the marginal distributions of subsets of the features. Let $O_m = 1$ if modality m is observed (i.e., \mathbf{x}_m is known), and let $O_m = 0$ if it is missing or unobserved. Let $\mathbf{X} = \{\mathbf{x}_m : O_m = 1\}$ be the visible features. We now use the following objective:

$$\mathcal{L}(\theta, \phi | \mathbf{X}) = \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{X})} \left[\sum_{m:O_m=1} \lambda_m \log p_{\theta}(\mathbf{x}_m | \mathbf{z}) \right] - \beta D_{\text{KL}}(q_{\phi}(\mathbf{z} | \mathbf{X}) \| p(\mathbf{z})) \quad (21.43)$$

The key problem is how to compute the posterior $q_\phi(\mathbf{z}|\mathbf{X})$ given different subsets of features. In general this can be hard, since the inference network is a discriminative model that assumes all inputs are available. For example, if it is trained on (image, text) pairs, $q_\phi(\mathbf{z}|\mathbf{x}_1, \mathbf{x}_2)$, how can we compute the posterior just given an image, $q_\phi(\mathbf{z}|\mathbf{x}_1)$, or just given text, $q_\phi(\mathbf{z}|\mathbf{x}_2)$? (This issue arises in general with VAE when we have missing inputs.)

Fortunately, based on our conditional independence assumption between the modalities, we can compute the optimal form for $q_\phi(\mathbf{z}|\mathbf{X})$ given set of inputs by computing the exact posterior under the model, which is given by

$$p(\mathbf{z}|\mathbf{X}) = \frac{p(\mathbf{z})p(\mathbf{x}_1, \dots, \mathbf{x}_M|\mathbf{z})}{p(\mathbf{x}_1, \dots, \mathbf{x}_M)} = \frac{p(\mathbf{z})}{p(\mathbf{x}_1, \dots, \mathbf{x}_M)} \prod_{m=1}^M p(\mathbf{x}_m|\mathbf{z}) \quad (21.44)$$

$$= \frac{p(\mathbf{z})}{p(\mathbf{x}_1, \dots, \mathbf{x}_M)} \prod_{m=1}^M \frac{p(\mathbf{z}|\mathbf{x}_m)p(\mathbf{x}_m)}{p(\mathbf{z})} \quad (21.45)$$

$$\propto p(\mathbf{z}) \prod_{m=1}^M \frac{p(\mathbf{z}|\mathbf{x}_m)}{p(\mathbf{z})} \approx p(\mathbf{z}) \prod_{m=1}^M \tilde{q}(\mathbf{z}|\mathbf{x}_m) \quad (21.46)$$

This can be viewed as a product of experts (Section 24.1.1), where each $\tilde{q}(\mathbf{z}|\mathbf{x}_m)$ is an “expert” for the m ’th modality, and $p(\mathbf{z})$ is the prior. We can compute the above posterior for any subset of modalities for which we have data by modifying the product over m . If we use Gaussian distributions for the prior $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_0, \boldsymbol{\Lambda}_0^{-1})$ and marginal posterior ratio $\tilde{q}(\mathbf{z}|\mathbf{x}_m) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_m, \boldsymbol{\Lambda}_m^{-1})$, then we can compute the product of Gaussians using the result from Equation (2.154):

$$\prod_{m=0}^M \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_m, \boldsymbol{\Lambda}_m^{-1}) \propto \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma}), \quad \boldsymbol{\Sigma} = (\sum_m \boldsymbol{\Lambda}_m)^{-1}, \quad \boldsymbol{\mu} = \boldsymbol{\Sigma}(\sum_m \boldsymbol{\Lambda}_m \boldsymbol{\mu}_m) \quad (21.47)$$

Thus the overall posterior precision is the sum of individual expert posterior precisions, and the overall posterior mean is the precision weighted average of the individual expert posterior means. See Figure 21.5(b) for an illustration. For a linear Gaussian (factor analysis) model, we can ensure $q(\mathbf{z}|\mathbf{x}_m) = p(\mathbf{z}|\mathbf{x}_m)$, in which case the above solution is the exact posterior [WN18], but in general it will be an approximation.

We need to train the individual expert recognition models $q(\mathbf{z}|\mathbf{x}_m)$ as well as the joint model $q(\mathbf{z}|\mathbf{X})$, so the model knows what to do with fully observed as well as partially observed inputs at test time. In [Ved+18], they propose a somewhat complex “triple ELBO” objective. In [WG18], they propose the simpler approach of optimizing the ELBO for the fully observed feature vector, all the marginals, and a set of \mathcal{J} randomly chosen joint modalities:

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}|\mathbf{X}) = \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}|(\mathbf{x}_1, \dots, \mathbf{x}_M)) + \sum_{m=1}^M \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}|\mathbf{x}_m) + \sum_{j \in \mathcal{J}} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}|\mathbf{X}_j) \quad (21.48)$$

This generalizes nicely to the semi-supervised setting, in which we only have a few aligned (“labeled”) examples from the joint, but have many unaligned (“unlabeled”) examples from the individual marginals. See Figure 21.5(c) for an illustration.

Note that the above scheme can only handle the case of a fixed number of missingness patterns; we can generalize to allow for arbitrary missingness as discussed in [CNW20]. (See also Section 3.11 for a more general discussion of missing data.)

21.3.4 Semisupervised VAEs

In this section, we discuss how to extend VAEs to the **semi-supervised learning** setting in which we have both labeled data, $\mathcal{D}_L = \{(\mathbf{x}_n, y_n)\}$, and unlabeled data, $\mathcal{D}_U = \{(\mathbf{x}_n)\}$. We focus on the **M2** model, proposed in [Kin+14a].

The generative model has the following form:

$$p_{\theta}(\mathbf{x}, y) = p_{\theta}(y)p_{\theta}(\mathbf{x}|y) = p_{\theta}(y) \int p_{\theta}(\mathbf{x}|y, \mathbf{z})p_{\theta}(\mathbf{z})d\mathbf{z} \quad (21.49)$$

where \mathbf{z} is a latent variable, $p_{\theta}(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$ is the latent prior, $p_{\theta}(y) = \text{Cat}(y|\boldsymbol{\pi})$ the label prior, and $p_{\theta}(\mathbf{x}|y, \mathbf{z}) = p(\mathbf{x}|f_{\theta}(y, \mathbf{z}))$ is the likelihood, such as a Gaussian, with parameters computed by f (a deep neural network). The main innovation of this approach is to assume that data is generated according to both a latent class variable y as well as the continuous latent variable \mathbf{z} . The class variable y is observed for labeled data and unobserved for unlabeled data.

To compute the likelihood for the *labeled data*, $p_{\theta}(\mathbf{x}, y)$, we need to marginalize over \mathbf{z} , which we can do by using an inference network of the form

$$q_{\phi}(\mathbf{z}|y, \mathbf{x}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_{\phi}(y, \mathbf{x}), \text{diag}(\boldsymbol{\sigma}_{\phi}(y, \mathbf{x}))) \quad (21.50)$$

We then use the following variational lower bound

$$\log p_{\theta}(\mathbf{x}, y) \geq \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}, y)} [\log p_{\theta}(\mathbf{x}|y, \mathbf{z}) + \log p_{\theta}(y) + \log p_{\theta}(\mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x}, y)] = -\mathcal{L}(\mathbf{x}, y) \quad (21.51)$$

as is standard for VAEs (see Section 21.2). The only difference is that we observe two kinds of data: \mathbf{x} and y .

To compute the likelihood for the *unlabeled data*, $p_{\theta}(\mathbf{x})$, we need to marginalize over \mathbf{z} *and* y , which we can do by using an inference network of the form

$$q_{\phi}(\mathbf{z}, y|\mathbf{x}) = q_{\phi}(\mathbf{z}|\mathbf{x})q_{\phi}(y|\mathbf{x}) \quad (21.52)$$

$$q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_{\phi}(\mathbf{x}), \text{diag}(\boldsymbol{\sigma}_{\phi}(\mathbf{x}))) \quad (21.53)$$

$$q_{\phi}(y|\mathbf{x}) = \text{Cat}(y|\boldsymbol{\pi}_{\phi}(\mathbf{x})) \quad (21.54)$$

Note that $q_{\phi}(y|\mathbf{x})$ acts like a discriminative classifier, that imputes the missing labels. We then use the following variational lower bound:

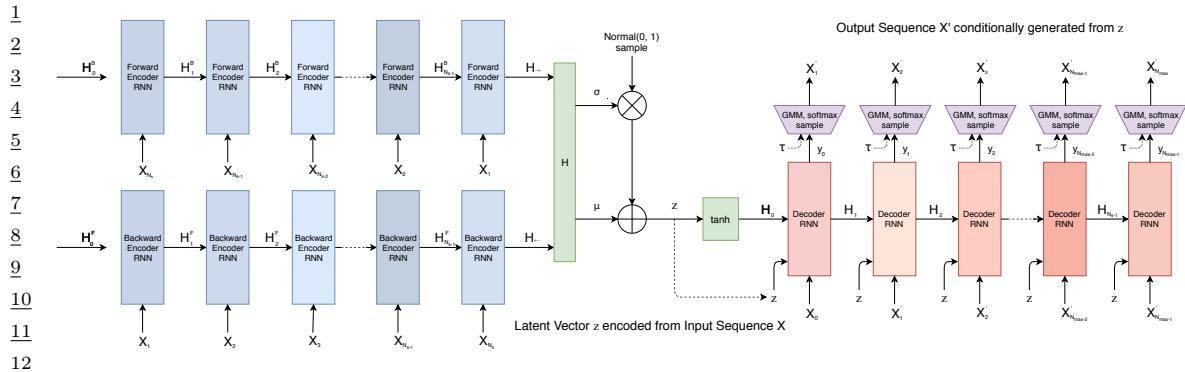
$$\log p_{\theta}(\mathbf{x}) \geq \mathbb{E}_{q_{\phi}(\mathbf{z}, y|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|y, \mathbf{z}) + \log p_{\theta}(y) + \log p_{\theta}(\mathbf{z}) - \log q_{\phi}(\mathbf{z}, y|\mathbf{x})] \quad (21.55)$$

$$= - \sum_y q_{\phi}(y|\mathbf{x}) \mathcal{L}(\mathbf{x}, y) + \mathbb{H}(q_{\phi}(y|\mathbf{x})) = -\mathcal{U}(\mathbf{x}) \quad (21.56)$$

Note that the discriminative classifier $q_{\phi}(y|\mathbf{x})$ is only used to compute the log-likelihood of the unlabeled data, which is undesirable. We can therefore add an extra classification loss on the supervised data, to get the following overall objective function:

$$\mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_L} [\mathcal{L}(\mathbf{x}, y)] + \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_U} [\mathcal{U}(\mathbf{x})] + \alpha \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_L} [-\log q_{\phi}(y|\mathbf{x})] \quad (21.57)$$

where \mathcal{D}_L is the labeled data, \mathcal{D}_U is the unlabeled data, and α is a hyperparameter that controls the relative weight of generative and discriminative learning.



13 Figure 21.6: Illustration of a VAE with a bidirectional RNN encoder and a unidirectional RNN decoder. The
14 output generator can use a GMM and/or softmax distribution. From Figure 2 of [HE18]. Used with kind
15 permission of David Ha.

1617

18 21.3.5 VAEs with sequential encoders/decoders

20 In this section, we discuss VAEs for sequential data, such as text and biosequences, in which the
21 data \mathbf{x} is a variable-length sequence, but we have a fixed-sized latent variable $\mathbf{z} \in \mathbb{R}^K$. (We consider
22 the more general case in which \mathbf{z} is a variable-length sequence of latents — known as **sequential**
23 **VAE** or **dynamic VAE** — in Section 29.13.) All we have to do is modify the decoder $p(\mathbf{x}|\mathbf{z})$ and
24 encoder $q(\mathbf{z}|\mathbf{x})$ to work with sequences.

25

26 21.3.5.1 Models

28 If we use an RNN for the encoder and decoder of a VAE, we get a model which is called a **VAE-RNN**,
29 as proposed in [Bow+16a]. In more detail, the generative model is $p(\mathbf{z}, \mathbf{x}_{1:T}) = p(\mathbf{z})\text{RNN}(\mathbf{x}_{1:T}|\mathbf{z})$,
30 where \mathbf{z} can be injected as the initial state of the RNN, or as an input to every time step. The
31 inference model is $q(\mathbf{z}|\mathbf{x}_{1:T}) = \mathcal{N}(\mathbf{z}|\mu(\mathbf{h}), \Sigma(\mathbf{h}))$, where $\mathbf{h} = [\mathbf{h}_T^\rightarrow, \mathbf{h}_1^\leftarrow]$ is the output of a bidirectional
32 RNN applied to $\mathbf{x}_{1:T}$. See Figure 21.6 for an illustration.

33 More recently, people have tried to combine transformers with VAEs. For example, in the **Optimus**
34 model of [Li+20], they use a BERT model for the encoder. In more detail, the encoder $q(\mathbf{z}|\mathbf{x})$ is
35 derived from the embedding vector associated with a dummy token corresponding to the “class label”
36 which is appended to the input sequence \mathbf{x} . The decoder is a standard autoregressive model (similar
37 to GPT), with one additional input, namely the latent vector \mathbf{z} . They consider two ways of injecting
38 the latent vector. The simplest approach is to add \mathbf{z} to the embedding layer of every token in the
39 decoding step, by defining $\mathbf{h}'_i = \mathbf{h}_i + \mathbf{W}\mathbf{z}$, where $\mathbf{h}_i \in \mathbb{R}^H$ is the original embedding for the i 'th
40 token, and $\mathbf{W} \in \mathbb{R}^{H \times K}$ is a decoding matrix, where K is the size of the latent vector. However, they
41 get better results in their experiments by letting all the layers of the decoder attend to the latent
42 code \mathbf{z} . An easy way to do this is to define the memory vector $\mathbf{h}_m = \mathbf{W}\mathbf{z}$, where $\mathbf{W} \in \mathbb{R}^{LH \times K}$,
43 where L is the number of layers in the decoder, and then to append $\mathbf{h}_m \in \mathbb{R}^{L \times H}$ to all the other
44 embeddings at each layer.

45 An alternative approach, known as **transformer VAE**, was proposed in [Gre20]. This model uses
46 a **funnel transformer** [Dai+20b] as the encoder, and the T5 [Raf+20a] conditional transformer for
47

he was silent for a long moment .
he was silent for a moment .
it was quiet for a moment .
it was dark and cold .
there was a pause .
it was my turn .

i went to the store to buy some groceries .
i store to buy some groceries .
i were to buy any groceries .
horses are to buy any groceries .
horses are to buy any animal .
horses the favorite any animal .
horses the favorite favorite animal .
horses are my favorite animal .

(a)

(b)

Figure 21.7: (a) Samples from the latent space of a VAE text model, as we interpolate between two sentences (on first and last line). Note that the intermediate sentences are grammatical, and semantically related to their neighbors. From Table 8 of [Bow+16b]. (b) Same as (a), but now using a deterministic autoencoder (with the same RNN encoder and decoder). From Table 1 of [Bow+16b]. Used with kind permission of Sam Bowman.

the decoder. In addition, it uses an MMD VAE (Section 21.3.2.1) to avoid posterior collapse.

21.3.5.2 Applications

In this section, we discuss some applications of VAEs to sequence data.

Text

In [Bow+16b], they apply the VAE-RNN model to natural language sentences. (See also [MB16; SSB17] for related work.) Although this does not improve performance in terms of the standard perplexity measures (predicting the next word given the previous words), it does provide a way to infer a semantic representation of the sentence. This can then be used for latent space interpolation, as discussed in Section 20.3.5. The results of doing this with the VAE-RNN are illustrated in Figure 21.7a. (Similar results are shown in [Li+20], using a VAE-transformer.) By contrast, if we use a standard deterministic autoencoder, with the same RNN encoder and decoder networks, we learn a much less meaningful space, as illustrated in Figure 21.7b. The reason is that the deterministic autoencoder has “holes” in its latent space, which get decoded to nonsensical outputs.

However, because RNNs (and transformers) are powerful decoders, we need to address the problem of posterior collapse, which we discuss in Section 21.4. One common way to avoid this problem is to use KL annealing, but a more effective method is to use the InfoVAE method of Section 21.3.2, which includes adversarial autoencoders (used in [She+20] with an RNN decoder) and MMD autoencoders (used in [Gre20] with a transformer decoder).

Sketches

In [HE18], they apply the VAE-RNN model to generate sketches (line drawings) of various animals and hand-written characters. They call their model **sketch-rnn**. The training data records the sequence of (x, y) pen positions, as well as whether the pen was touching the paper or not. The emission model used a GMM for the real-valued location offsets, and a categorical softmax distribution for the discrete state.

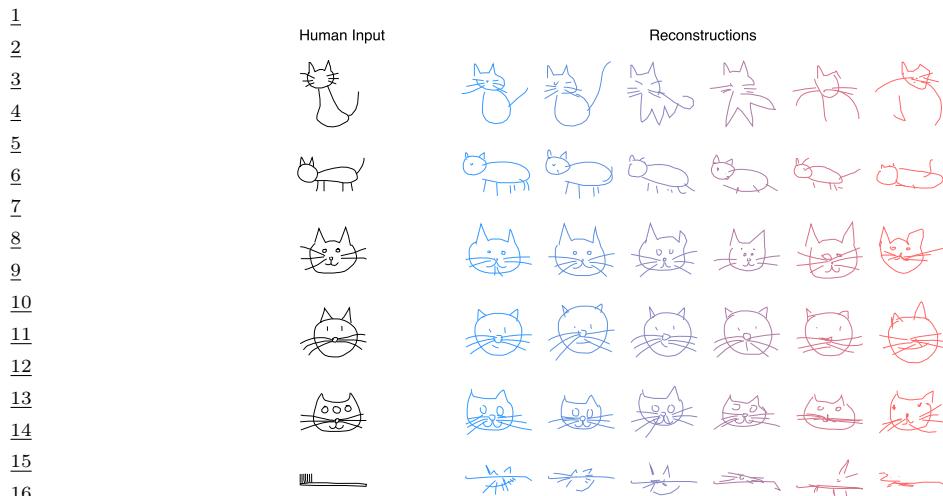


Figure 21.8: Conditional generation of cats from sketch-RNN model. We increase the temperature parameter from left to right. From Figure 5 of [HE18]. Used with kind permission of David Ha.

Figure 21.8 shows some samples from various class-conditional models. We vary the temperature parameter τ of the emission model to control the stochasticity of the generator. (More precisely, we multiply the GMM variances by τ , and divide the discrete probabilities by τ before renormalizing.) When the temperature is low, the model tries to reconstruct the input as closely as possible. However, when the input is untypical of the training set (e.g., a cat with three eyes, or a toothbrush), the reconstruction is “regularized” towards a canonical cat with two eyes, while still keeping some features of the input.

30 Molecular design

In [GB+18], they use VAE-RNNs to model molecular graph structure, represented as a string using the SMILES representation.⁴ It is also possible to learn a mapping from the latent space to some scalar quantity of interest, such as the solubility or drug efficacy of a molecule. We can then perform gradient-based optimization in the continuous latent space to try to generate new graphs which maximize this quantity. See Figure 21.9 for a sketch of this approach.

The main problem is to ensure that points in latent space decode to valid strings/molecules. There are various solutions to this, including using a **grammar VAE**, where the RNN decoder is replaced by a stochastic context free grammar. See [KPHL17] for details.

41 21.4 Avoiding posterior collapse

If the decoder $p_{\theta}(x|z)$ is sufficiently powerful (e.g., a pixel CNN, or an RNN for text), then the VAE does not need to use the latent code z for anything. This is called **posterior collapse** or **variational**

⁴⁶ 4. See https://en.wikipedia.org/wiki/Simplified_molecular_input_line-entry_system.

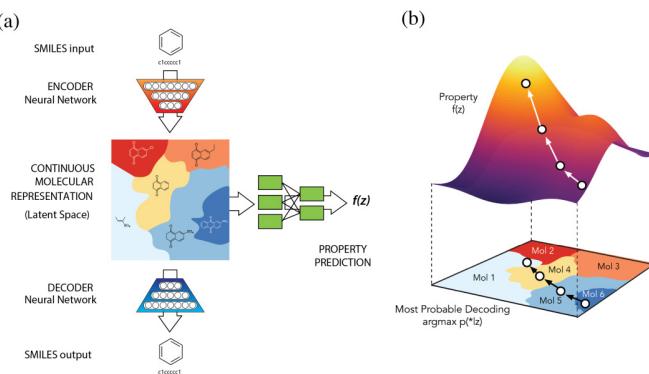


Figure 21.9: Application of VAE-RNN to molecule design. (a) The VAE-RNN model is trained on a sequence representation of molecules known as SMILES. We can fit an MLP to map from the latent space to properties of the molecule, such as its “fitness” $f(\mathbf{z})$. (b) We can perform gradient ascent in $f(\mathbf{z})$ space, and then decode the result to a new molecule with high fitness. From Figure 1 of [GB+18]. Used with kind permission of Rafael Gomez-Bombarelli.

overpruning (see e.g., [Che+17b; Ale+18; Hus17a; Phu+18; TT17; Yeu+17; Luc+19; DWW19; WBC21]). To see why this happens, consider Equation (21.21). If there exists a parameter setting for the generator θ^* such that $p_{\theta^*}(\mathbf{x}|\mathbf{z}) = p_{\mathcal{D}}(\mathbf{x})$ for every \mathbf{z} , then we can make $D_{\text{KL}}(p_{\mathcal{D}}(\mathbf{x}) \parallel p_{\theta}(\mathbf{x})) = 0$. Since the generator is independent of the latent code, we have $p_{\theta}(\mathbf{z}|\mathbf{x}) = p_{\theta}(\mathbf{z})$. The prior $p_{\theta}(\mathbf{z})$ is usually a simple distribution, such as a Gaussian, so we can find a setting of the inference parameters so that $q_{\phi^*}(\mathbf{z}|\mathbf{x}) = p_{\theta}(\mathbf{z})$, which ensures $D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p_{\theta}(\mathbf{z}|\mathbf{x})) = 0$. Thus we have successfully maximized the ELBO, but we have not learned any useful latent representation of the data, which is one of the goals of latent variable modeling.⁵ We discuss some solutions to posterior collapse below.

21.4.1 KL annealing

A common approach to solving this problem, proposed in [Bow+16a], is to use **KL annealing**, in which the KL penalty term in the ELBO is scaled by β , which is increased from 0.0 (corresponding to an autoencoder) to 1.0 (which corresponds to standard MLE training). (Note that, by contrast, the β -VAE model in Section 21.3.1 uses $\beta > 1$.)

KL annealing can work well, but requires tuning the schedule for β . A standard practice [Fu+19] is to use **cyclical annealing**, which repeats the process of increasing β multiple times. This ensures the progressive learning of more meaningful latent codes, by leveraging good representations learned in a previous cycle as a way to warmstart the optimization.

5. Note that [Luc+19; DWW20] show that posterior collapse can also happen in linear VAE models, where the ELBO corresponds to the exact marginal likelihood, so the problem is not only due to powerful (nonlinear) decoders, but is also related to spurious local maxima in the objective.

1 **21.4.2 Lower bounding the rate**

3 An alternative approach is to stick with the original unmodified ELBO objective, but to prevent the
4 rate (i.e., the $D_{\text{KL}}(q \parallel p)$ term) from collapsing to 0, by limiting the flexibility of q . For example,
5 [XD18; Dav+18] use a von Mises-Fisher (Section 2.2.5.3) prior and posterior, instead of a Gaussian,
6 and they constrain the posterior to have a fixed concentration, $q(\mathbf{z}|\mathbf{x}) = \text{vMF}(\mathbf{z}|\boldsymbol{\mu}(\mathbf{x}), \kappa)$. Here
7 the parameter κ controls the rate of the code. The δ -VAE method [Oor+19] uses a Gaussian
8 autoregressive prior and a diagonal Gaussian posterior. We can ensure the rate is at least δ by
9 adjusting the regression parameter of the AR prior.

10

11 **21.4.3 Free bits**

12 In this section, we discuss the method of **free bits** [Kin+16], which is another way of lower bounding
13 the rate. To explain this, consider a fully factorized posterior in which the KL penalty has the form

14

$$\mathcal{L}_R = \sum_i D_{\text{KL}}(q_{\phi}(z_i|\mathbf{x}) \parallel p_{\theta}(z_i)) \quad (21.58)$$

15 where z_i is the i 'th dimension of \mathbf{z} . We can replace this with a hinge loss, that will give up driving
16 down the KL for dimensions that are already beneath a target compression rate λ :

17

$$\mathcal{L}'_R = \sum_i \max(\lambda, D_{\text{KL}}(q_{\phi}(z_i|\mathbf{x}) \parallel p_{\theta}(z_i))) \quad (21.59)$$

18 Thus the bits where the KL is sufficiently small “are free”, since the model does not have to “pay” to
19 encode them according to the prior.

20

21 **21.4.4 Adding skip connections**

22 One reason for latent variable collapse is that the latent variables \mathbf{z} are not sufficiently “connected to”
23 the observed data \mathbf{x} . One simple solution is to modify the architecture of the generative model by
24 adding **skip connections**, similar to a residual network (Section 16.2.4), as shown in Figure 21.10.
25 This is called a **skip-VAE** [Die+19a].

26

27 **21.4.5 Improved variational inference**

28 The posterior collapse problem is caused in part by the poor approximation to the posterior. In
29 [He+19], they proposed to keep the model and VAE objective unchanged, but to more aggressively
30 update the inference network before each step of generative model fitting. This enables the inference
31 network to capture the current true posterior more faithfully, which will encourage the generator to
32 use the latent codes when it is useful to do so.

33 However, this only addresses the part of posterior collapse that is due to the amortization gap
34 [CLD18], rather than the more fundamental problem of variational pruning, in which the KL term
35 penalizes the model if its posterior deviates too far from the prior, which is often too simple to match
36 the aggregated posterior.

37 Another way to ameliorate variational pruning is to use lower bounds that are tighter than the
38 vanilla ELBO (Section 10.5.1), or more accurate posterior approximations (Section 10.4), or more
39 accurate (hierarchical) generative models (Section 21.5).

40

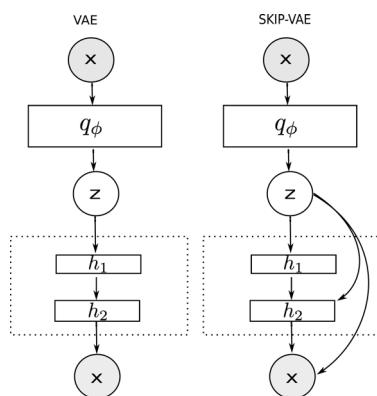


Figure 21.10: (a) VAE. (b) Skip-VAE. From Figure 1 of [Die+19a]. Used with kind permission of Adji Dieng.

21.4.6 Alternative objectives

An alternative to the above methods is to replace the ELBO objective with other objectives, such as the InfoVAE objective discussed in Section 21.3.2, which includes adversarial autoencoders and MMD autoencoders as special cases. The InfoVAE objective includes a term to explicitly enforce non-zero mutual information between \mathbf{x} and \mathbf{z} , which effectively solves the problem of posterior collapse.

21.5 VAEs with hierarchical structure

We define a **hierarchical VAE** or HVAE, with L stochastic layers, to be the following generative model:⁶

$$p_{\theta}(\mathbf{x}, \mathbf{z}_{1:L}) = p_{\theta}(\mathbf{z}_L) \left[\prod_{l=L-1}^1 p_{\theta}(\mathbf{z}_l | \mathbf{z}_{l+1}) \right] p_{\theta}(\mathbf{x} | \mathbf{z}_1) \quad (21.60)$$

We can improve on the above model by making it non-Markovian, i.e., letting each \mathbf{z}_l depend on all the higher level stochastic variables, $\mathbf{z}_{l+1:L}$, not just the preceding level, i.e.,

$$p_{\theta}(\mathbf{x}, \mathbf{z}) = p_{\theta}(\mathbf{z}_L) \left[\prod_{l=L-1}^1 p_{\theta}(\mathbf{z}_l | \mathbf{z}_{l+1:L}) \right] p_{\theta}(\mathbf{x} | \mathbf{z}_{1:L}) \quad (21.61)$$

Note that the likelihood is now $p_{\theta}(\mathbf{x} | \mathbf{z}_{1:L})$ instead of just $p_{\theta}(\mathbf{x} | \mathbf{z}_1)$. This is analogous to adding skip connections from all preceding variables to all their children. It is easy to implement this by using a deterministic “backbone” of residual connections, that accumulates all stochastic decisions, and propagates them down the chain, as illustrated in Figure 21.11(left). We discuss how to perform inference and learning in such models below.

6. There is a split in the literature about whether to label the top level as \mathbf{z}_L or \mathbf{z}_1 . We adopt the former convention, since we view lower numbered layers, such as \mathbf{z}_1 , as being “closer to the data”, and higher numbered layers, such as \mathbf{z}_L , as being “more abstract”.

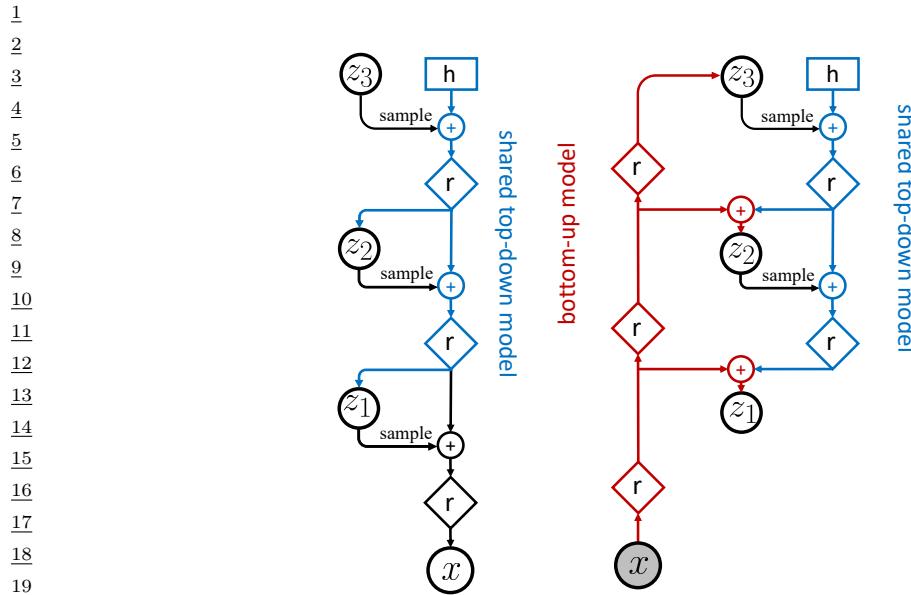


Figure 21.11: Hierarchical VAEs with 3 stochastic layers. Left: generative model. Right: inference network. Diamond is a residual network, \oplus is feature combination (e.g., concatenation), and h is a trainable parameter. We first do bottom-up inference, by propagating x up to z_3 to compute $z_3^s \sim q_\phi(z_3|x)$, and then we perform top-down inference by computing $z_2^s \sim q_\phi(z_2|x, z_3^s)$ and then $z_1^s \sim q_\phi(z_1|x, z_{2:3}^s)$. From Figure 2 of [VK20a]. Used with kind permission of Arash Vahdat.

21.5.1 Bottom-up vs top-down inference

To perform inference in a hierarchical VAE, we could use a **bottom-up inference model** of the form

$$q_\phi(z|x) = q_\phi(z_1|x) \prod_{l=2}^L q_\phi(z_l|x, z_{1:l-1}) \quad (21.62)$$

However, a better approach is to use a **top-down inference model** of the form

$$q_\phi(z|x) = q_\phi(z_L|x) \prod_{l=L-1}^1 q_\phi(z_l|x, z_{l+1:L}) \quad (21.63)$$

Inference for z_l combines bottom-up information from x with top-down information from higher layers, $z_{>l} = z_{l+1:L}$. See Figure 21.11(right) for an illustration.⁷

⁷ Note that it is also possible to have a stochastic bottom-up encoder and a stochastic top-down encoder, as discussed in the **BIVA** paper [Maa+19]. (BIVA stands for “bidirectional-inference variational autoencoder”.)

With the above model, the ELBO can be written as follows (using the chain rule for KL):

$$\underline{L}(\theta, \phi | \mathbf{x}) = \mathbb{E}_{q_\phi(\mathbf{z} | \mathbf{x})} [\log p_\theta(\mathbf{x} | \mathbf{z})] - D_{\text{KL}}(q_\phi(\mathbf{z}_L | \mathbf{x}) \| p_\theta(\mathbf{z}_L)) \quad (21.64)$$

$$-\sum_{l=L-1}^1 \mathbb{E}_{q_\phi(\mathbf{z}_{>l} | \mathbf{x})} [D_{\text{KL}}(q_\phi(\mathbf{z}_l | \mathbf{x}, \mathbf{z}_{>l}) \| p_\theta(\mathbf{z}_l | \mathbf{z}_{>l}))] \quad (21.65)$$

where

$$q_\phi(\mathbf{z}_{>l} | \mathbf{x}) = \prod_{i=l+1}^L q_\phi(\mathbf{z}_i | \mathbf{x}, \mathbf{z}_{>i}) \quad (21.66)$$

is the approximate posterior above layer l (i.e., the parents of \mathbf{z}_l).

The reason the top-down inference model is better is that it more closely approximates the true posterior of a given layer, which is given by

$$p_\theta(\mathbf{z}_l | \mathbf{x}, \mathbf{z}_{l+1:L}) \propto p_\theta(\mathbf{z}_l | \mathbf{z}_{l+1:L}) p_\theta(\mathbf{x} | \mathbf{z}_l, \mathbf{z}_{l+1:L}) \quad (21.67)$$

Thus the posterior combines the top-down prior term $p_\theta(\mathbf{z}_l | \mathbf{z}_{l+1:L})$ with the bottom-up likelihood term $p_\theta(\mathbf{x} | \mathbf{z}_l, \mathbf{z}_{l+1:L})$. We can approximate this posterior by defining

$$q_\phi(\mathbf{z}_l | \mathbf{x}, \mathbf{z}_{l+1:L}) \propto p_\theta(\mathbf{z}_l | \mathbf{z}_{l+1:L}) \tilde{q}_\phi(\mathbf{z}_l | \mathbf{x}, \mathbf{z}_{l+1:L}) \quad (21.68)$$

where $\tilde{q}_\phi(\mathbf{z}_l | \mathbf{x}, \mathbf{z}_{l+1:L})$ is a learned Gaussian approximation to the bottom-up likelihood. If both prior and likelihood are Gaussian, we can compute this product in closed form, as proposed in the **ladder network** paper [Sn+16; Søn+16].⁸ A more flexible approach is to let $q_\phi(\mathbf{z}_l | \mathbf{x}, \mathbf{z}_{l+1:L})$ be learned, but to force it to share some of its parameters with the learned prior $p_\theta(\mathbf{z}_l | \mathbf{z}_{l+1:L})$, as proposed in [Kin+16]. This reduces the number of parameters in the model, and ensures that the posterior and prior remain somewhat close.

21.5.2 Example: very deep VAE

There have been many papers exploring different kinds of HVAE models (see e.g., [Kin+16; Sn+16; Chi21a; VK20a; Maa+19]), and we do not have space to discuss them all. Here we focus on the “very deep VAE” or **VD-VAE** model of [Chi21a], since it is simple but yields state of the art results (at the time of writing).

The architecture is a simple convolutional VAE with bidirectional inference, as shown in Figure 21.12. For each layer, the prior and posterior are diagonal Gaussians. The author found that nearest-neighbor upsampling (in the decoder) worked much better than transposed convolution, and avoided posterior collapse. This enabled training with the vanilla VAE objective, without needing any of the tricks discussed in Section 21.5.4.

The low-resolution latents (at the top of the hierarchy) capture a lot of the global structure of each image; the remaining high-resolution latents are just used to fill in details, that make the image look more realistic, and improve the likelihood. This suggests the model could be useful for lossy

⁸ The term “ladder network” arises from the horizontal “rungs” in Figure 21.11(right). Note that a similar idea was independently proposed in [Sal16].

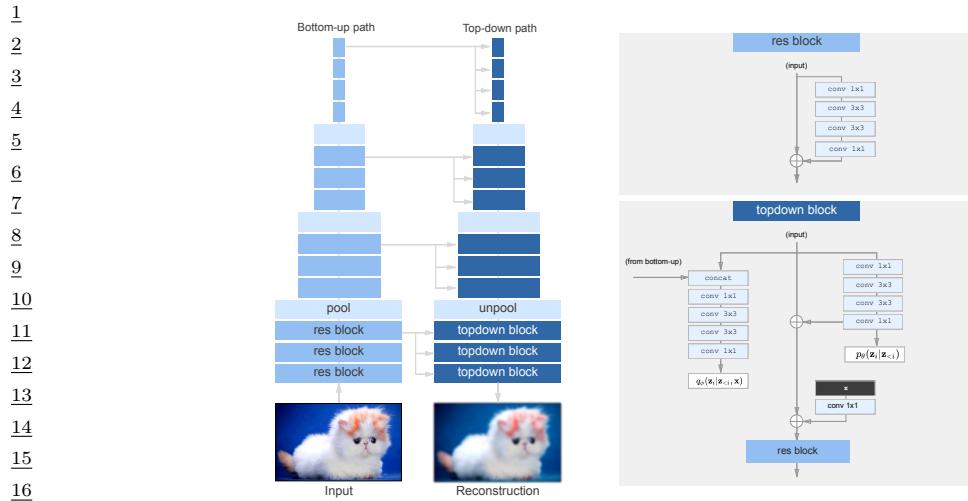


Figure 21.12: The top-down encoder used by the hierarchical VAE in [Chi21a]. Each convolution is preceded by the GELU nonlinearity. The model uses average pooling and nearest-neighbor upsampling for the pool and unpool layers. The posterior q_ϕ and prior p_θ are diagonal Gaussians. From Figure 3 of [Chi21a]. Used with kind permission of Rewon Child.



Figure 21.13: Samples from a VDVAE model (trained on FFHQ dataset) from different levels of the hierarchy. From Figure 1 of [Chi21a]. Used with kind permission of Rewon Child.

compression, since a lot of the low-level details can be drawn from the prior (i.e., “hallucinated”), rather than having to be sent by the encoder.

We can also use the model for unconditional sampling at multiple resolutions. This is illustrated in Figure 21.13, using a model with 78 stochastic layers trained on the FFHQ-256 dataset.⁹

21.5.3 Connection with autoregressive models

Until recently, most hierarchical VAEs only had a small number of stochastic layers. Consequently the images they generated have not looked as good, or had as high likelihoods, as images produced by other models, such as the autoregressive PixelCNN model (see Section 22.3.2). However, by endowing VAEs with many more stochastic layers, it is possible to outperform AR models in terms of

⁹ 9. This is a 256² version of the Flickr-Faces High Quality dataset from <https://github.com/NVlabs/ffhq-dataset>, which has 80k images at 1024² resolution.

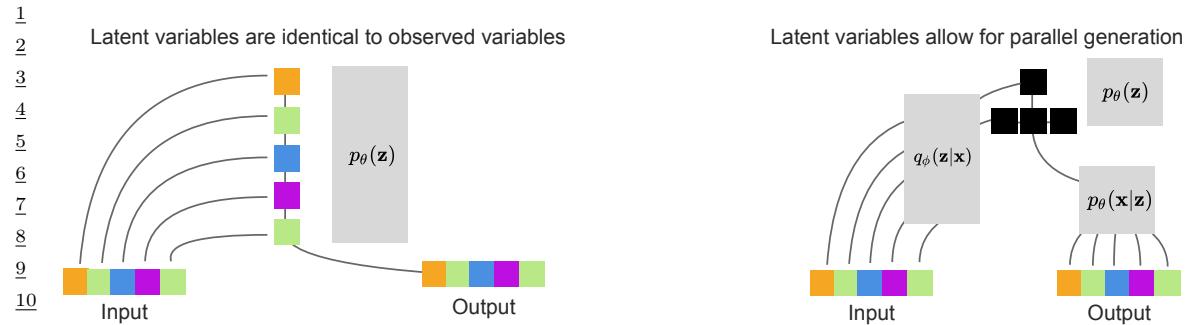


Figure 21.14: Left: a hierarchical VAE which emulates an autoregressive model using an identity encoder, autoregressive prior, and identity decoder. Right: a hierarchical VAE with a 2 layer hierarchical latent code. The bottom hidden nodes (black) are conditionally independent given the top layer. From Figure 2 of [Chi21a]. Used with kind permission of Rewon Child.

likelihood and sample quality, while using fewer parameters and much less computing power [Chi21a; VK20a; Maa+19].

To see why this is possible, note that we can represent any AR model as a degenerate VAE, as shown in Figure 21.14(left). The idea is simple: the encoder copies the input into latent space by setting $\mathbf{z}_{1:D} = \mathbf{x}_{1:D}$ (so $q_\phi(z_i = x_i | \mathbf{z}_{>i}, \mathbf{x}) = 1$), then the model learns an autoregressive prior $p_\theta(\mathbf{z}_{1:D}) = \prod_d p(z_d | \mathbf{z}_{1:d-1})$, and finally the likelihood function just copies the latent vector to output space, so $p_\theta(x_i = z_i | \mathbf{z}) = 1$. Since the encoder computes the exact (albeit degenerate) posterior, we have $q_\phi(\mathbf{z} | \mathbf{x}) = p_\theta(\mathbf{z} | \mathbf{x})$, so the ELBO is tight and reduces to the log likelihood,

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) = \log p_{\boldsymbol{\theta}}(\mathbf{z}) = \sum_d \log p_{\boldsymbol{\theta}}(x_d | \mathbf{x}_{< d}) \quad (21.69)$$

Thus we can emulate any AR model with a VAE providing it has at least D stochastic layers, where D is the dimensionality of the observed data.

34 In practice, data usually lives in a lower-dimensional manifold (see e.g., [DW19]), which can allow
35 for a much more compact latent code. For example, Figure 21.14(right) shows a hierarchical code
36 in which the latent factors at the lower level are conditionally independent given the higher level,
37 and hence can be generated in parallel. Such a tree-like structure can enable sample generation in
38 $O(\log D)$ time, whereas an autoregressive model always takes $O(D)$ time. (Recall that for an image
39 D is the number of pixels, so it grows quadratically with image resolution. For example, even a tiny
40 32×32 image has $D = 3072$.)

41 In addition to speed, hierarchical models also require many fewer parameters than “flat” models.
42 The typical architecture used for generating images is a **multi-scale** approach: the model starts from
43 a small, spatially arranged set of latent variables, and at each subsequent layer, the spatial resolution
44 is increased (usually by a factor of 2). This allows the high level to capture global, long-range
45 correlations (e.g., the symmetry of a face, or overall skin tone), while letting lower levels capture
46 fine-grained details.

1 **21.5.4 Variational pruning**

3 A common problem with hierarchical VAEs is that the higher level latent layers are often ignored, so
4 the model does not learn interesting high level semantics. This is caused by **variational pruning**.
5 This problem is analogous to the issue of latent variable collapse, which we discussed in Section 21.4.

6 A common heuristic to mitigate this problem is to use KL balancing coefficients [Che+17b], to
7 ensure that an equal amount of information is encoded in each layer. That is, we use the following
8 penalty:

$$\sum_{l=1}^L \gamma_l \mathbb{E}_{q_\phi(\mathbf{z}_{>l}|\mathbf{x})} [D_{\text{KL}}(q_\phi(\mathbf{z}_l|\mathbf{x}, \mathbf{z}_{>l}) \| p_\theta(\mathbf{z}_l|\mathbf{z}_{>l}))] \quad (21.70)$$

13 The balancing term γ_l is set to a small value when the KL penalty is small (on the current minibatch),
14 to encourage use of that layer, and is set to a large value when the KL term is large. (This is only
15 done during the “warm up period”.) Concretely, [VK20a] proposes to set the coefficients γ_l to be
16 proportional to the size of the layer, s_l , and the average KL loss:

$$\gamma_l \propto s_l \mathbb{E}_{\mathbf{x} \sim \mathcal{B}} [\mathbb{E}_{q_\phi(\mathbf{z}_{>l}|\mathbf{x})} [D_{\text{KL}}(q_\phi(\mathbf{z}_l|\mathbf{x}, \mathbf{z}_{>l}) \| p_\theta(\mathbf{z}_l|\mathbf{z}_{>l}))]] \quad (21.71)$$

19 where \mathcal{B} is the current minibatch.

20

21 **21.5.5 Other optimization difficulties**

23 A common problem when training (hierarchical) VAEs is that the loss can become unstable. The
24 main reason for this is that the KL term is unbounded (can become infinitely large). In [Chi21a], they
25 tackle the problem in two ways. First, ensure the initial random weights of the final convolutional
26 layer in each residual bottleneck block get scaled by $1/\sqrt{L}$. Second, skip an update step if the norm
27 of the gradient of the loss exceeds some threshold.

28 In the **Nouveau VAE** method of [VK20a], they use some more complicated measures to ensure
29 stability. First, they use batch normalization, but with various tweaks. Second, they use spectral
30 regularization for the encoder. Specifically they add the penalty $\beta \sum_i \lambda_i$, where λ_i is the largest
31 singular value of the i 'th convolutional layer (estimated using a single power iteration step), and
32 $\beta \geq 0$ is a tuning parameter. Third, they use inverse autoregressive flows (Section 23.2.4.3) in each
33 layer, instead of a diagonal Gaussian approximation. Fourth, they represent the posterior using a
34 residual representation. In particular, let us assume the prior for the i 'th variable in layer l is

$$p_\theta(z_l^i | \mathbf{z}_{>l}) = \mathcal{N}(z_l^i | \mu_i(\mathbf{z}_{>l}), \sigma_i(\mathbf{z}_{>l})) \quad (21.72)$$

37 They propose the following posterior approximation:

$$q_\phi(z_l^i | \mathbf{x}, \mathbf{z}_{>l}) = \mathcal{N}(z_l^i | \mu_i(\mathbf{z}_{>l}) + \Delta\mu_i(\mathbf{z}_{>l}, \mathbf{x}), \sigma_i(\mathbf{z}_{>l}) \cdot \Delta\sigma_i(\mathbf{z}_{>l}, \mathbf{x})) \quad (21.73)$$

40 where the Δ terms are the relative changes computed by the encoder. The corresponding KL penalty
41 reduces to the following (dropping the l subscript for brevity):

$$D_{\text{KL}}(q_\phi(z^i | \mathbf{x}, \mathbf{z}_{>l}) \| p_\theta(z^i | \mathbf{z}_{>l})) = \frac{1}{2} \left(\frac{\Delta\mu_i^2}{\sigma_i^2} + \Delta\sigma_i^2 - \log \Delta\sigma_i^2 - 1 \right) \quad (21.74)$$

45 So as long as σ_i is bounded from below, the KL term can be easily controlled just by adjusting the
46 encoder parameters.

47



Figure 21.15: Autoencoder for MNIST using 256 binary latents. Top row: input images. Middle row: reconstruction. Bottom row: latent code, reshaped to a 16×16 image. Generated by [quanzitized_autoencoder_mnist.ipynb](#).

21.6 Vector quantization VAE

In this section, we describe **VQ-VAE**, which stands for “vector quantized VAE” [OVK17; ROV19]. This is like a standard VAE except it uses a set of discrete latent variables.

21.6.1 Autoencoder with binary code

The simplest approach to the problem is to construct a standard VAE, but to add a discretization layer at the end of the encoder, $\mathbf{z}_e(\mathbf{x}) \in \{0, \dots, S - 1\}^K$, where S is the number of states, and K is the number of discrete latents. For example, we can binarize the latent vector (using $S = 2$) by clipping \mathbf{z} to lie in $\{0, 1\}^K$. This can be useful for data compression (see e.g., [BLS17]).

Suppose we assume the prior over the latent codes is uniform. Since the encoder is deterministic, the KL divergence reduces to a constant, equal to $\log K$. This avoids the problem with posterior collapse (Section 21.4). Unfortunately, the discontinuous quantization operation of the encoder prohibits the direct use of gradient based optimization. The solution proposed in [OVK17] is to use the straight-through estimator, which we discuss in Section 6.3.8. We show a simple example of this approach in Figure 21.15, where we use a Gaussian likelihood, so the loss function has the form

$$\mathcal{L} = \|\mathbf{x} - d(e(\mathbf{x}))\|_2^2 \quad (21.75)$$

where $e(\mathbf{x}) \in \{0, 1\}^K$ is the encoder, and $d(\mathbf{z}) \in \mathbb{R}^{28 \times 28}$ is the decoder.

21.6.2 VQ-VAE model

We can get a more expressive model by using a 3d tensor of discrete latents, $\mathbf{z} \in \mathbb{R}^{H \times W \times K}$, where K is the number of discrete values per latent variable. Rather than just binarizing the continuous vector $\mathbf{z}_e(\mathbf{x})_{ij}$, we compare it to a **codebook** of embedding vectors, $\{\mathbf{e}_k : k = 1 : K, \mathbf{e}_k \in \mathbb{R}^L\}$, and then set \mathbf{z}_{ij} to the index of the nearest codebook entry:

$$q(\mathbf{z}_{ij} = k | \mathbf{x}) = \begin{cases} 1 & \text{if } k = \operatorname{argmin}_{k'} \|\mathbf{z}_e(\mathbf{x})_{i,j,:} - \mathbf{e}_{k'}\|_2 \\ 0 & \text{otherwise} \end{cases} \quad (21.76)$$

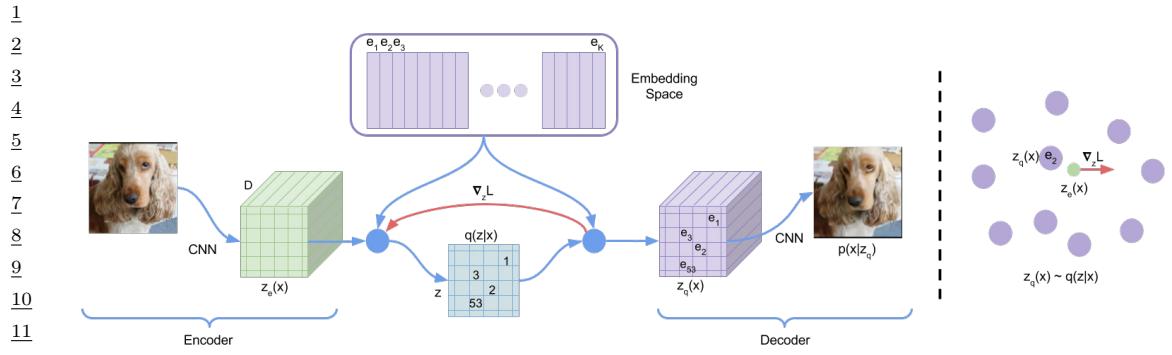


Figure 21.16: VQ-VAE architecture. From Figure 1 of [OVK17]. Used with kind permission of Aäron van den Oord.

When reconstructing the input we replace each discrete code index by the corresponding real-valued codebook vector:

$$(z_q)_{ij} = e_k \text{ where } z_{ij} = k \quad (21.77)$$

These values are then passed to the decoder, $p(\mathbf{x}|z_q)$, as usual. See Figure 21.16 for an illustration of the overall architecture. Note that although \mathbf{z}_q is generated from a discrete combination of codebook vectors, the use of a distributed code makes the model very expressive. For example, if we use a grid of 32×32 , with $K = 512$, then we can generate $512^{32 \times 32} = 2^{9216}$ distinct images, which is astronomically large.

To fit this model, we can minimize the negative log likelihood (reconstruction error) using the straight-through estimator, as before. This amounts to passing the gradients from the decoder input $\mathbf{z}_q(\mathbf{x})$ to the encoder output $\mathbf{z}_e(\mathbf{x})$, bypassing Equation (21.76), as shown by the red arrow in Figure 21.16. Unfortunately this means that the codebook entries will not get any learning signal. To solve this, the authors proposed to add an extra term to the loss, known as the **codebook loss**, that encourages the codebook entries e to match the output of the encoder. We treat the encoder $\mathbf{z}_e(\mathbf{x})$ as a fixed target, by adding a **stop gradient** operator to it; this ensures \mathbf{z}_e is treated normally in the forwards pass, but has zero gradient in the backwards pass. The modified loss (dropping the spatial indices i, j) becomes

$$\mathcal{L} = -\log p(\mathbf{x}|z_q(\mathbf{x})) + \|\text{sg}(\mathbf{z}_e(\mathbf{x})) - \mathbf{e}\|_2^2 \quad (21.78)$$

where \mathbf{e} refers to the codebook vector assigned to $\mathbf{z}_e(\mathbf{x})$, and sg is the stop gradient operator.

An alternative way to update the codebook vectors is to use moving averages. To see how this works, first consider the batch setting. Let $\{\mathbf{z}_{i,1}, \dots, \mathbf{z}_{i,n_i}\}$ be the set of n_i outputs from the encoder that are closest to the dictionary item \mathbf{e}_i . We can update \mathbf{e}_i to minimize the MSE

$$\sum_{j=1}^{n_i} \|\mathbf{z}_{i,j} - \mathbf{e}_i\|_2^2 \quad (21.79)$$

1 which has the closed form update
2

$$\underline{3} \quad \underline{4} \quad \underline{5} \quad \mathbf{e}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} \mathbf{z}_{i,j} \quad (21.80)$$

6 This is like the M step of the EM algorithm when fitting the mean vectors of a GMM. In the minibatch
7 setting, we replace the above operations with an exponentially moving average, as follows:
8

$$\underline{9} \quad \underline{10} \quad N_i^t = \gamma N_i^{t-1} + (1 - \gamma) n_i^t \quad (21.81)$$

$$\underline{11} \quad \underline{12} \quad \mathbf{m}_i^t = \gamma \mathbf{m}_i^{t-1} + (1 - \gamma) \sum_j \mathbf{z}_{i,j}^t \quad (21.82)$$

$$\underline{13} \quad \underline{14} \quad \mathbf{e}_i^t = \frac{\mathbf{m}_i^t}{N_i^t} \quad (21.83)$$

16 The authors found $\gamma = 0.9$ to work well.

17 The above procedure will learn to update the codebook vectors so it matches the output of the
18 encoder. However, it is also important to ensure the encoder does not “change its mind” too often
19 about what codebook value to use. To prevent this, the authors propose to add a third term to
20 the loss, known as the **commitment loss**, that encourages the encoder output to be close to the
21 codebook values. Thus we get the final loss:

$$\underline{22} \quad \underline{23} \quad \mathcal{L} = -\log p(\mathbf{x}|\mathbf{z}_q(\mathbf{x})) + \|\text{sg}(\mathbf{z}_e(\mathbf{x})) - \mathbf{e}\|_2^2 + \beta \|\mathbf{z}_e(\mathbf{x}) - \text{sg}(\mathbf{e})\|_2^2 \quad (21.84)$$

24 The authors found $\beta = 0.25$ to work well, although of course the value depends on the scale of the
25 reconstruction loss (NLL) term. (A probabilistic interpretation of this loss can be found in [Hen+18].)
26 Overall, the decoder optimizes the first term only, the encoder optimizes the first and last term, and
27 the embeddings optimize the middle term.
28

29 30 21.6.3 Learning the prior

31 After training the VQ-VAE model, it is possible to learn a better prior, to match the aggregated
32 posterior. To do this, we just apply the encoder to a set of data, $\{\mathbf{x}_n\}$, thus converting them to
33 discrete sequences, $\{\mathbf{z}_n\}$. We can then learn a joint distribution $p(\mathbf{z})$ using any kind of sequence
34 model. In the original VQ-VAE paper [OVK17], they used the causal convolutional PixelCNN model
35 (Section 22.3.2). More recent work has used transformer decoders (Section 22.4). Samples from this
36 prior can then be decoded using the decoder part of the VQ-VAE model. We give some examples of
37 this in the sections below.
38

39 40 21.6.4 Hierarchical extension (VQ-VAE-2)

41 In [ROV19], they extend the original VQ-VAE model by using a hierarchical latent code. The model
42 is illustrated in Figure 21.17. They applied this to images of size $256 \times 256 \times 3$. The first latent layer
43 maps this to a quantized representation of size 64×64 , and the second latent layer maps this to a
44 quantized representation of size 32×32 . This hierarchical scheme allows the top level to focus on
45 high level semantics of the image, leaving fine visual details, such as texture, to the lower level. (See
46 Section 21.5 for more discussion of hierarchical VAEs.)
47

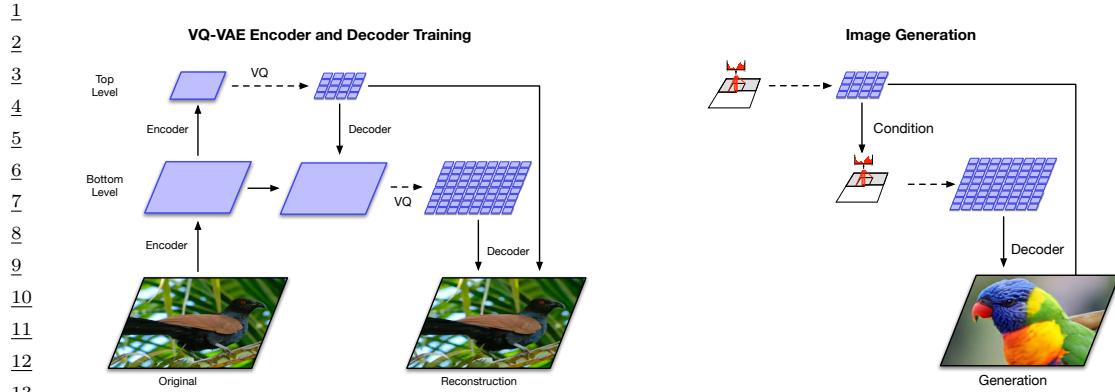


Figure 21.17: Hierarchical extension of VQ-VAE. (a) Encoder and decoder architecture. (b) Combining a Pixel-CNN prior with the decoder. From Figure 2 of [ROV19]. Used with kind permission of Aaron van den Oord.

After fitting the VQ-VAE, they learn a prior over the top level code using a PixelCNN model augmented with self-attention (Section 16.2.7) to capture long-range dependencies. (This hybrid model is known as PixelSNAIL [Che+17c].) For the lower level prior, they just use standard PixelCNN, since attention would be too expensive. Samples from the model can then be decoded using the VQ-VAE decoder, as shown in Figure 21.17.

21.6.5 Discrete VAE

In VQ-VAE, we use a one-hot encoding for the latents, $q(z = k|\mathbf{x}) = 1$ iff $k = \operatorname{argmin}_k \|\mathbf{z}_e(\mathbf{x}) - \mathbf{e}_k\|_2$, and then set $\mathbf{z}_q = \mathbf{e}_k$. This does not capture any uncertainty in the latent code, and requires the use of the straight-through estimator for training.

Various other approaches to fitting VAEs with discrete latent codes have been investigated. In the DALL-E paper (Section 22.4.2), they use a fairly simple method, based on using the Gumbel-softmax relaxation for the discrete variables (see Section 6.3.6). In brief, let $q(z = k|\mathbf{x})$ be the probability that the input \mathbf{x} is assigned to codebook entry k . We can exactly sample $w_k \sim q(z = k|\mathbf{x})$ from this by computing $w_k = \operatorname{argmax}_k g_k + \log q(z = k|\mathbf{x})$, where each g_k is from a Gumbel distribution. We can now “relax” this by using a softmax with temperature $\tau > 0$ and computing

$$w_k = \frac{\exp(\frac{g_k + \log q(z=k|\mathbf{x})}{\tau})}{\sum_{j=1}^K \exp(\frac{g_j + \log q(z=j|\mathbf{x})}{\tau})} \quad (21.85)$$

We now set the latent code to be a weighted sum of the codebook vectors:

$$\mathbf{z}_q = \sum_{k=1}^K w_k \mathbf{e}_k \quad (21.86)$$

In the limit that $\tau \rightarrow 0$, the distribution over weights \mathbf{w} converges to a one-hot distribution, in which case \mathbf{z} becomes equal to one of the codebook entries. But for finite τ , we “fill in” the space between the vectors.

47

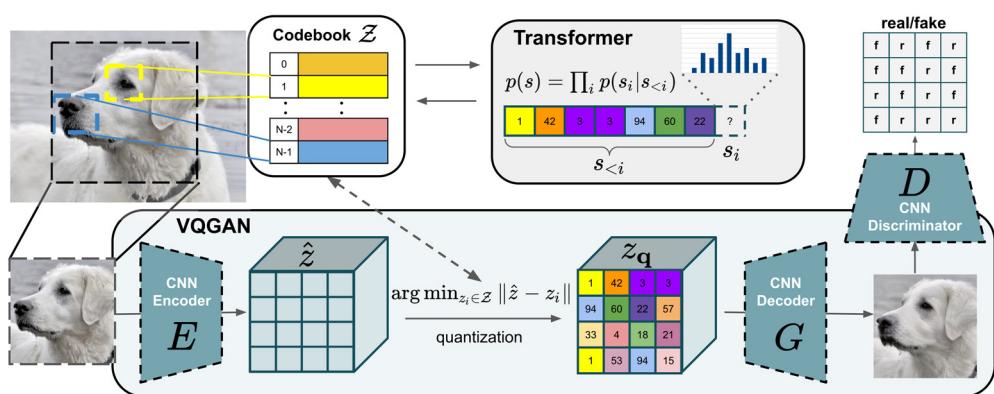


Figure 21.18: Illustration of the VQ-GAN. From Figure 2 of [ERO21]. Used with kind permission of Patrick Esser.

This allows us to express the ELBO in the usual differentiable way:

$$\mathcal{L} = -\mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z})] + \beta D_{\text{KL}}(q(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z})) \quad (21.87)$$

where $\beta > 0$ controls the amount of regularization. (Unlike VQ-VAE, the KL term is not a constant, because the encoder is stochastic.) Furthermore, since the Gumbel noise variables are sampled from a distribution that is independent of the encoder parameters, we can use the reparameterization trick (Section 6.3.5) to optimize this.

21.6.6 VQ-GAN

One drawback of VQ-VAE is that it uses mean squared error in its reconstruction loss, which can result in blurry samples. In the **VQ-GAN** paper [ERO21], they replace this with a (patch-wise) GAN loss (see Chapter 26), together with a perceptual loss; this results in much higher visual fidelity. In addition, they use a transformer (see Section 16.3.5) to model the prior on the latent codes. See Figure 21.18 for a visualization of the overall model. In [Yu+21], they replace the CNN encoder and decoder of the VQ-GAN model with transformers, yielding improved results; they call this **VIM** (vector-quantized image modeling).

22 Autoregressive models

22.1 Introduction

By the chain rule of probability, we can write any joint distribution over T variables as follows:

$$p(\mathbf{x}_{1:T}) = p(\mathbf{x}_1)p(\mathbf{x}_2|\mathbf{x}_1)p(\mathbf{x}_3|\mathbf{x}_2, \mathbf{x}_1)p(\mathbf{x}_4|\mathbf{x}_3, \mathbf{x}_2, \mathbf{x}_1) \dots = \prod_{t=1}^T p(\mathbf{x}_t|\mathbf{x}_{1:t-1}) \quad (22.1)$$

where $\mathbf{x}_t \in \mathcal{X}$ is the t 'th observation, and we define $p(\mathbf{x}_1|\mathbf{x}_{1:0}) = p(x_1)$ as the initial state distribution. This is called an **autoregressive model** or **ARM**. This corresponds to a fully connected DAG, in which each node depends on all its predecessors in the ordering, as shown in Figure 22.1. The models can also be conditioned on arbitrary inputs or context \mathbf{c} , in order to define $p(\mathbf{x}|\mathbf{c})$, although we omit this for notational brevity.

We could of course also factorize the joint distribution “backwards” in time, using

$$p(\mathbf{x}_{1:T}) = \prod_{t=T}^1 p(\mathbf{x}_t|\mathbf{x}_{t+1:T}) \quad (22.2)$$

However, this “anti-causal” direction is often harder to learn (see e.g., [PJS17]).

Although the decomposition in Equation (22.1) is general, each term in this expression (i.e., each conditional distribution $p(\mathbf{x}_t|\mathbf{x}_{1:t-1})$) becomes more and more complex, since it depends on an increasing number of arguments, which makes the terms slow to compute, and makes estimating their parameters more data hungry (see Section 2.6.3.2).

One approach to solving this intractability is to make the (first-order) **Markov assumption**, which gives rise to a **Markov model** $p(\mathbf{x}_t|\mathbf{x}_{1:t-1}) = p(\mathbf{x}_t|\mathbf{x}_{t-1})$, which we discuss in Section 2.6. (This is also called an auto-regressive model of order 1.) Unfortunately, the Markov assumption is very limiting. One way to relax it, and to make \mathbf{x}_t depend on all the past $\mathbf{x}_{1:t-1}$ without explicitly regressing on them, is to assume the past can be compressed into a **hidden state** \mathbf{z}_t . If \mathbf{z}_t is a deterministic function of the past observations $\mathbf{x}_{1:t-1}$, the resulting model is known as a **recurrent neural network**, discussed in Section 16.3.4. If \mathbf{z}_t is a stochastic function of the past hidden state, \mathbf{z}_{t-1} , the resulting model is known as a **hidden Markov model**, which we discuss in Section 29.2.

Another approach is to stay with the general AR model of Equation (22.1), but to use a restricted functional form, such as some kind of neural network, for the conditionals $p(\mathbf{x}_t|\mathbf{x}_{1:t-1})$. Thus rather than making conditional independence assumptions, or explicitly compressing the past into a sufficient

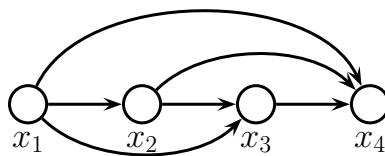


Figure 22.1: A fully-connected auto-regressive model.

11 statistic, we implicitly learn a compact mapping from the past to the future. In the sections below,
12 we discuss different functional forms for these conditional distributions.

13 The main advantage of such AR models is that it is easy to compute, and optimize, the exact
14 likelihood of each sequence (data vector). The main disadvantage is that generating samples is
15 inherently sequential, which can be slow. In addition, the method does not learn a compact latent
16 representation of the data.

18 22.2 Neural autoregressive density estimators (NADE)

21 A simple way to represent each conditional probability distribution $p(x_t | \mathbf{x}_{1:t-1})$ is to use a generalized
22 linear model, such as logistic regression, as proposed in [Fre98]. We can make the model be more
23 powerful by using a neural network. The resulting model is called the **neural auto-regressive**
24 **density estimator** or **NADE** model [LM11].

25 If we let $p(x_t | \mathbf{x}_{1:t-1})$ be a conditional mixture of Gaussians, we get a model known as **RNADE**
26 (“real-valued neural autoregressive density estimator”) of [UML13]. More precisely, this has the form

$$\text{p}(\mathbf{x}_t | \mathbf{x}_{1:t-1}) = \sum_{k=1}^K \pi_{t,k} \mathcal{N}(x_t | \mu_{t,k}, \sigma_{t,k}^2) \quad (22.3)$$

31 where the parameters are generated by a network, $(\boldsymbol{\mu}_t, \boldsymbol{\sigma}_t, \boldsymbol{\pi}_t) = f_t(\mathbf{x}_{1:t-1}; \boldsymbol{\theta}_t)$.

32 Rather than using separate neural networks, f_1, \dots, f_T , it is more efficient to create a single
33 network with T inputs and T outputs. This can be done using masking, resulting in a model called
34 the **MADE** (“masked autoencoder for density estimation”) model [Ger+15].

35 One disadvantage of NADE-type models is that they assume the variables have a natural linear
36 ordering. This makes sense for temporal or sequential data, but not for more general data types,
37 such as images or graphs. An orderless extension to NADE was proposed in [UML14; Uri+16].

40 22.3 Causal CNNs

42 One approach to representing the distribution $p(\mathbf{x}_t | \mathbf{x}_{1:t-1})$ is to try to identify patterns in the past
43 history that might be predictive of the value of x_t . If we assume these patterns can occur in any
44 location, it makes sense to use a **convolutional neural network** to detect them. However, we need
45 to make sure we only apply the convolutional mask to past inputs, not future ones. This can be done
46 using **masked convolution**, also called **causal convolution**. We discuss this in more detail below.

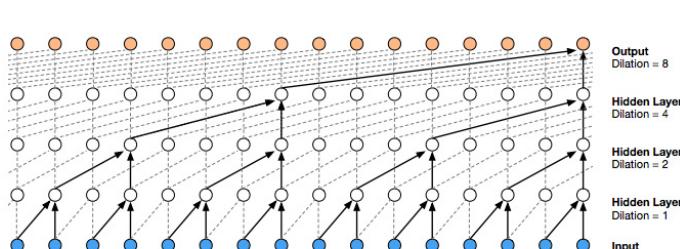


Figure 22.2: Illustration of the wavenet model using dilated (atrous) convolutions, with dilation factors of 1, 2, 4 and 8. From Figure 3 of [oor+16]. Used with kind permission of Aäron van den Oord.

22.3.1 1d causal CNN (convolutional Markov models)

Consider the following **convolutional Markov model** for 1d discrete sequences:

$$p(\mathbf{x}_{1:T}) = \prod_{t=1}^T p(x_t | \mathbf{x}_{1:t-1}; \boldsymbol{\theta}) = \prod_{t=1}^T \text{Cat}(x_t | \text{softmax}(\varphi(\sum_{\tau=1}^{t-k} \mathbf{w}^\top \mathbf{x}_{\tau:\tau+k}))) \quad (22.4)$$

where \mathbf{w} is the convolutional filter of size k , and we have assumed a single nonlinearity φ and categorical output, for notational simplicity. This is like regular 1d convolution except we “mask out” future inputs, so that x_t only depends on the past values. We can of course use deeper models, and we can condition on input features \mathbf{c} .

In order to capture long-range dependencies, we can use **dilated convolution** (see [Mur22, Sec 14.4.1]). This model has been successfully used to create a state of the art **text to speech** (TTS) synthesis system known as **wavenet** [oor+16]. See Figure 22.2 for an illustration.

The wavenet model is a conditional model, $p(\mathbf{x}|\mathbf{c})$, where \mathbf{c} is a set of linguistic features derived from an input sequence of words, and \mathbf{x} is raw audio. The **tacotron** system [Wan+17c] is a fully end-to-end approach, where the input is words rather than linguistic features.

Although wavenet produces high quality speech, it is too slow for use in production systems. However, it can be “distilled” into a parallel generative model [Oor+18], as we discuss in Section 23.2.4.3.

22.3.2 2d causal CNN (PixelCNN)

We can extend causal convolutions to 2d, to get an autoregressive model of the form

$$p(\mathbf{x}|\boldsymbol{\theta}) = \prod_{r=1}^R \prod_{c=1}^C p(x_{r,c} | f_{\boldsymbol{\theta}}(\mathbf{x}_{1:r-1,1:C}, \mathbf{x}_{r,1:c-1})) \quad (22.5)$$

where R is the number of rows, C is the number of columns, and we condition on all previously generated pixels in a **raster scan** order, as illustrated in Figure 22.3. This is called the **pixelCNN** model [Oor+16]. Naive sampling (generation) from this model takes $O(N)$ time, where $N = RC$ is the number of pixels, but [Ree+17] shows how to use a multiscale approach to reduce the complexity to $O(\log N)$.

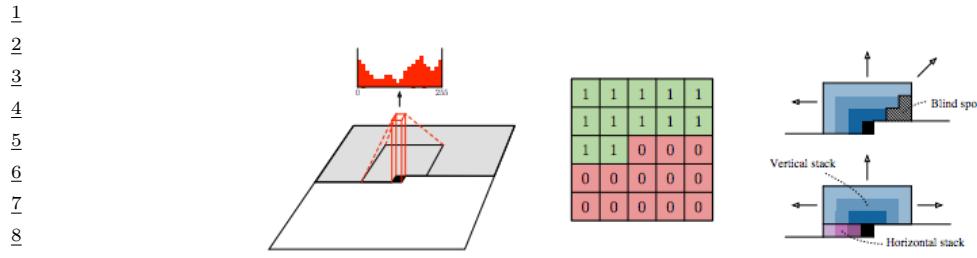


Figure 22.3: Illustration of causal 2d convolution in the PixelCNN model. The red histogram shows the empirical distribution over discretized values for a single pixel of a single RGB channel. The red and green 5×5 array shows the binary mask, which selects the top left context, in order to ensure the convolution is causal. The diagrams on the right illustrate how we can avoid blind spots by using a vertical context stack, that contains all previous rows, and a horizontal context stack, that just contains values from the current row. From Figure 1 of [Oor+16]. Used with kind permission of Aaron van den Oord.

Various extensions of this model have been proposed. The **pixelCNN++** model of [Sal+17c] improved the quality by using a mixture of logistic distributions, to capture the multimodality of $p(x_i | \mathbf{x}_{1:i-1})$. The **pixelRNN** of [OKK16] combined masked convolution with an RNN to get even longer range contextual dependencies. The **Subscale Pixel Network** of [MK19] proposed to generate the pixels such that the higher order bits are sampled before lower order bits, which allows high resolution details to be sampled conditioned on low resolution versions of the whole image, rather than just the top left corner.

22.4 Transformers

We introduced transformers in Section 16.3.5. They can be used for encoding sequences (as in BERT), or for decoding (generating) sequences. We can also combine the two, using an encoder-decoder combination, for conditional generation from $p(\mathbf{y}|\mathbf{c})$. Alternatively, we can define a joint sequence model $p(\mathbf{c}, \mathbf{y})$, where \mathbf{c} is the conditioning or context prompt, and then just condition the joint model, by giving it as the initial context.

The decoder (generator) works as follows. At each step t , the model applies masked (causal) self attention (Section 16.2.7) to the first t inputs, $\mathbf{y}_{1:t}$, to compute a set of attention weights, $\mathbf{a}_{1:t}$. From this it computes an activation vector $\mathbf{z}_t = \sum_{\tau=1}^t a_{t\tau} \mathbf{y}_\tau$. This is then passed through a feed-forward layer to compute $\mathbf{h}_t = \text{MLP}(\mathbf{z}_t)$. This process is repeated for each layer in the model. Finally the output is used to predict the next element in the sequence, $\mathbf{y}_{t+1} \sim \text{Cat}(\text{softmax}(\mathbf{W}\mathbf{h}_t))$.

At training time, all predictions can happen in parallel, since the target generated sequence is already available. That is, the t 'th output \mathbf{y}_t can be predicted given inputs $\mathbf{y}_{1:t-1}$, and this can be done for all t simultaneously. However, at test time, the model must be applied sequentially, so the output generated at $t+1$ is fed back into the model to predict $t+2$, etc. Note that the running time of transformers is $O(T^2)$, although a variety of more efficient versions have been developed (see e.g., [Mur22, Sec 15.6] for details).

Transformers are the basis of many popular (conditional) generative models for sequences. We give some examples below.

A "whatpu" is a small, furry animal native to Tanzania. An example of a sentence that uses the word whatpu is:
 We were traveling in Africa and we saw these very cute whatpus.

 To do a "farduddle" means to jump up and down really fast. An example of a sentence that uses the word farduddle is:
One day when I was playing tag with my little sister, she got really excited and she started doing these crazy farduckles.
 A "yalubalu" is a type of vegetable that looks like a big pumpkin. An example of a sentence that uses the word yalubalu is:
I was on a trip to Africa and I tried this yalubalu vegetable that was grown in a garden there. It was delicious.

 A "Burringo" is a car with very fast acceleration. An example of a sentence that uses the word Burringo is:
In our garage we have a Burringo that my father drives to work every day.

 A "Gigamuru" is a type of Japanese musical instrument. An example of a sentence that uses the word Gigamuru is:
I have a Gigamuru that my uncle gave me as a gift. I love to play it at home.

 To "screeg" something is to swing a sword at it. An example of a sentence that uses the word screeg is:
We screeghed at each other for several minutes and then we went outside and ate ice cream.

Figure 22.4: Illustration of few shot learning with GPT-3. The model is asked to create an example sentence using a new word whose meaning is provided in the prompt. Boldface is GPT-3's completions, light gray is human input. From Figure 3.16 of [Bro+20d].

22.4.1 Text generation (GPT, etc.)

In [Rad+18], OpenAI proposed a model called **GPT**, which is short for “Generative Pre-training Transformer”. This is a decoder-only transformer model that uses causal (masked) attention. In [Rad+19], they propose **GPT-2**, which is a larger version of GPT (1.5 billion parameters, or 6.5GB, for the XL version), trained on a large web corpus (8 million pages, or 40GB). They also simplify the training objective, and just train it using maximum likelihood.

After GPT-2, OpenAI released **GPT-3** [Bro+20d], which is an even larger version of GPT-2 (175 billion parameters), trained on even more data (300 billion words), but based on the same principles. (Training was estimated to take 355 GPU years and cost \$4.6M.) Due to the large size of the data and model, GPT-3 shows even more remarkable abilities to generate novel text. In particular, the output can be (partially) controlled by just changing the conditioning prompt. This enables the model to perform tasks that it has never been trained on, just by giving it some examples in the prompt. This is called “**in-context learning**”. See Figure 22.4 for an example, and Section 19.5.1.2 for more details.

More recently, OpenAI released **ChatGPT** [Ope], which is an improved version of GPT-3 which has been trained to have interactive dialogs by using a technique called **reinforcement learning from human feedback** or **RLHF**, a technique first introduced in the **InstructGPT** paper [Ouy+22]. This uses RL techniques such as PPO (Section 35.3.4) to fine tune the model so that it generates responses that are more “aligned” with human intent, as estimated by a ranking model, which is pre-trained on supervised data.

```

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

```

(JAX):
 Yo, PyTorch, you think you're hot stuff.
 But JAX is where it's at, no need to bluff.
 Your popularity might be on the rise
 But when it comes to performance, we're the guys.

With automatic differentiation, we're a breeze
 Training neural nets with ease
 Our functional approach is innovative and unique
 PyTorch can't touch our speed, try and keep up, meek

(PyTorch):
 Hold up, JAX, you might be fast
 But PyTorch is the one that's built to last
 We may not have the fancy functional flair
 But our flexibility sets us above the pair

Our community is strong and growing each day
 Supporting us in every single way
 We may not be the newest on the block
 But we'll always be the torch that leads the flock

Figure 22.5: ChatGPT response to the prompt “Write a rap battle about PyTorch and JAX”. Used with kind permission of Paige Bailey. From <https://twitter.com/DynamicWebPaige/status/1601743574369902593>.

Despite the impressive performance of these **large language models** or **LLMs** (see Figure 22.5 for an example), there are several open problems with them, such as: they often confidently **hallucinate** incorrect answers to questions (see e.g., [Ji+22]); they can generate biased or toxic output (see e.g., [Lia+]); and they are very resource intensive to train and serve. Indeed, these concerns are why Google has not (at the time of writing) released its version of ChatGPT, known as **LaMDA** [Col21].

The basic ideas behind LLMs are quite simple (maximum likelihood training of an autoregressive transformer), and they can be implemented in about 300 lines of code.¹ However, just by scaling up the size of the models and datasets, it seems that qualitatively new capabilities can emerge (see e.g., [Wei+22]). Nevertheless, although this approach is good at learning formal linguistic competence (surface form), it is not clear if it is sufficient to learn functional linguistic competence, which requires a deeper, non-linguistic understanding of the world derived from experience [Mah+23].

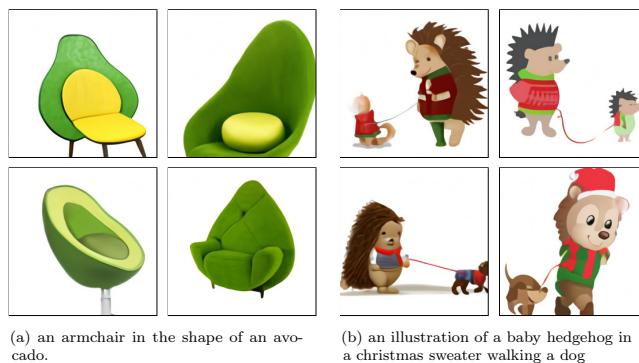
22.4.2 Image generation (DALL-E, etc.)

The **DALL-E** model² from OpenAI [Ram+21a] can generate images of remarkable quality and diversity given text prompts, as shown in Figure 22.6. The methodology is conceptually quite straightforward, and most of the effort went into data collection (they scraped the web for 250 million image-text pairs) and scaling up the training (they fit a model with 12 billion parameters). Here we just focus on the algorithmic methods.

The basic idea is to transform an image x into a sequence of discrete tokens z using a discrete

⁴⁵ 1. See e.g., <https://github.com/karpathy/nanoGPT>.

⁴⁶ 2. The name is derived from the artist Salvador Dalí and Pixar’s movie “WALL-E”



(a) an armchair in the shape of an avocado.
 (b) an illustration of a baby hedgehog in a christmas sweater walking a dog

Figure 22.6: Some images generated by the DALL-E model in response to a text prompt. (a) “An armchair in the shape of an avocado”. (b) “An illustration of a baby hedgehog in a christmas sweater walking a dog”. From <https://openai.com/blog/dall-e>. Used with kind permission of Aditya Ramesh.

VAE model (Section 21.6.5). We then fit a transformer to the concatenation of the image tokens \mathbf{z} and text tokens \mathbf{y} to get a joint model of the form $p(\mathbf{z}, \mathbf{y})$.

To sample an image \mathbf{x} given a text prompt \mathbf{y} , we sample a latent code $\mathbf{z} \sim p(\mathbf{z}|\mathbf{y})$ by conditioning the transformer on the prompt \mathbf{y} , and then we feed \mathbf{z} into the VAE decoder to get the image $\mathbf{x} \sim p(\mathbf{x}|\mathbf{z})$. Multiple images are generated for each prompt, and these are then ranked according to a pre-trained critic, which gives them scores depending on how well the generated image matches the input text: $s_n = \text{critic}(\mathbf{x}_n, \mathbf{y}_n)$. The critic they used was the contrastive CLIP model (see Section 32.3.4.1). This discriminative reranking significantly improves the results.

Some sample results are shown in Figure 22.6, and more can be found online at <https://openai.com/blog/dall-e/>. The image on the right of Figure 22.6 is particularly interesting, since the prompt — “An illustration of a baby hedgehog in a christmas sweater walking a dog” — arguably requires that the model solve the “**variable binding problem**”. This refers to the fact that the sentence implies the hedgehog should be wearing the sweater and not the dog. We see that the model sometimes interprets this correctly, but not always: sometimes it draws both animals with Christmas sweaters. In addition, sometimes it draws a hedgehog walking a smaller hedgehog. The quality of the results can also be sensitive to the form of the prompt.

The **PARTI** model [Yu+22] from Google follows similar high level ideas to DALL-E, but has been scaled to an even larger size. The larger models perform qualitatively much better, as shown in Figure 20.3.

Other recent approaches to (conditional) image generation — such as **DALL-E 2** [Ram+22] from Open-AI, **Imagen** [Sah+22b] from Google, and **Stable diffusion** [Rom+22] from Stability.AI — are based on diffusion rather than applying a transformer to discretized image patches. See Section 25.6.4 for details.

1
2 **22.4.3 Other applications**

3 Transformers have been used to generate many other kinds of (discrete) data, such as midi music
4 sequences [Hua+18a], protein sequences [Gan+23], etc.
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

23 Normalizing flows

This chapter is written by George Papamakarios and Balaji Lakshminarayanan.

23.1 Introduction

In this chapter we discuss **normalizing flows**, a class of flexible density models that can be easily sampled from and whose exact likelihood function is efficient to compute. Such models can be used for many tasks, such as density modeling, inference and generative modeling. We introduce the key principles of normalizing flows and refer to recent surveys by Papamakarios et al. [Pap+19] and Kobyzev, Prince, and Brubaker [KPB19] for readers interested in learning more. See also <https://github.com/janosh/awesome-normalizing-flows> for a list of papers and software packages.

23.1.1 Preliminaries

Normalizing flows create complex probability distributions $p(\mathbf{x})$ by passing random variables $\mathbf{u} \in \mathbb{R}^D$, drawn from a simple **base distribution** $p(\mathbf{u})$ through a nonlinear but *invertible* transformation $\mathbf{f} : \mathbb{R}^D \rightarrow \mathbb{R}^D$. That is, $p(\mathbf{x})$ is defined by the following process:

$$\mathbf{x} = \mathbf{f}(\mathbf{u}) \quad \text{where} \quad \mathbf{u} \sim p(\mathbf{u}). \tag{23.1}$$

The base distribution is typically chosen to be simple, for example standard Gaussian or uniform, so that we can easily sample from it and compute the density $p(\mathbf{u})$. A flexible enough transformation \mathbf{f} can induce a complex distribution on the transformed variable \mathbf{x} even if the base distribution is simple.

Sampling from $p(\mathbf{x})$ is straightforward: we first sample \mathbf{u} from $p(\mathbf{u})$ and then compute $\mathbf{x} = \mathbf{f}(\mathbf{u})$. To compute the density $p(\mathbf{x})$, we rely on the fact that \mathbf{f} is invertible. Let $\mathbf{g}(\mathbf{x}) = \mathbf{f}^{-1}(\mathbf{x}) = \mathbf{u}$ be the inverse mapping, which “**normalizes**” the data distribution by mapping it back to the base distribution (which is often a normal distribution). Using the change-of-variables formula for random variables from Equation (2.257), we have

$$p_x(\mathbf{x}) = p_u(\mathbf{g}(\mathbf{x})) |\det \mathbf{J}(\mathbf{g})(\mathbf{x})| = p_u(\mathbf{u}) |\det \mathbf{J}(\mathbf{f})(\mathbf{u})|^{-1}, \tag{23.2}$$

where $\mathbf{J}(\mathbf{f})(\mathbf{u}) = \frac{\partial \mathbf{f}}{\partial \mathbf{u}}|_{\mathbf{u}}$ is the Jacobian matrix of \mathbf{f} evaluated at \mathbf{u} . Taking logs of both sides of Equation (23.2), we get

$$\log p_x(\mathbf{x}) = \log p_u(\mathbf{u}) - \log |\det \mathbf{J}(\mathbf{f})(\mathbf{u})|. \tag{23.3}$$

1 As discussed above, $p(\mathbf{u})$ is typically easy to evaluate. So, if one can use flexible invertible transformations \mathbf{f} whose Jacobian determinant $\det \mathbf{J}(\mathbf{f})(\mathbf{u})$ can be computed efficiently, then one can construct complex densities $p(\mathbf{x})$ that allow exact sampling and efficient exact likelihood computation. This is in contrast to latent variable models, which require methods like variational inference to lower-bound the likelihood.

2 One might wonder how flexible are the densities $p(\mathbf{x})$ obtained by transforming random variables
3 sampled from simple $p(\mathbf{u})$. It turns out that we can use this method to approximate any smooth
4 distribution. To see this, consider the scenario where the base distribution $p(\mathbf{u})$ is a one-dimensional
5 uniform distribution. Recall that inverse transform sampling (Section 11.3.1) samples random
6 variables from a uniform distribution and transforms them using the inverse cumulative distribution
7 function (cdf) to generate samples from the desired density. We can use this method to sample
8 from any one-dimensional density as long as the transformation \mathbf{f} is powerful enough to model the
9 inverse cdf (which is a reasonable assumption for well-behaved densities whose cdf is invertible and
10 differentiable). We can further extend this argument to multiple dimensions by first expressing the
11 density $p(\mathbf{x})$ as a product of one-dimensional conditionals using the chain rule of probability, and then
12 applying inverse transform sampling to each one-dimensional conditional. The result is a normalizing
13 flow that transforms a product of uniform distributions into any desired distribution $p(\mathbf{x})$. We refer
14 to [Pap+19] for a more detailed proof.

15 How do we define flexible invertible mappings whose Jacobian determinant is easy to compute?
16 We discuss this topic in detail in Section 23.2, but in summary, there are two main ways. The first
17 approach is to define a set of simple transformations that are invertible by design, and whose Jacobian
18 determinant is easy to compute; for instance, if the Jacobian is a triangular matrix, its determinant
19 can be computed efficiently. The second approach is to exploit the fact that a composition of invertible
20 functions is also invertible, and the overall Jacobian determinant is just the product of the individual
21 Jacobian determinants. More precisely, if $\mathbf{f} = \mathbf{f}_N \circ \dots \circ \mathbf{f}_1$ where each \mathbf{f}_i is invertible, then \mathbf{f} is also
22 invertible, with inverse $\mathbf{g} = \mathbf{g}_1 \circ \dots \circ \mathbf{g}_N$ and log Jacobian determinant given by

$$29 \quad \log |\det \mathbf{J}(\mathbf{g})(\mathbf{x})| = \sum_{i=1}^N \log |\det \mathbf{J}(\mathbf{g}_i)(\mathbf{u}_i)| \quad (23.4)$$

30 where $\mathbf{u}_i = \mathbf{f}_i \circ \dots \circ \mathbf{f}_1(\mathbf{u})$ is the i 'th intermediate output of the flow. This allows us to create
31 complex flows from simple components, just as graphical models allow us to create complex joint
32 distributions from simpler conditional distributions.

33 Finally, a note on terminology. An invertible transformation is also known as a **bijection**. A
34 bijection that is differentiable and has a differentiable inverse is known as a **diffeomorphism**. The
35 transformation \mathbf{f} of a flow model is a diffeomorphism, although in the rest of this chapter we will refer
36 to it as a “bijection” for simplicity, leaving the differentiability implicit. The density $p_x(\mathbf{x})$ of a flow
37 model is also known as the **pushforward** of the base distribution $p_u(\mathbf{u})$ through the transformation
38 \mathbf{f} , and is sometimes denoted as $p_x = \mathbf{f}_* p_u$. Finally, in mathematics the term “flow” refers to any
39 family of diffeomorphisms \mathbf{f}_t indexed by a real number t such that $t = 0$ indexes the identity function,
40 and $t_1 + t_2$ indexes $\mathbf{f}_{t_2} \circ \mathbf{f}_{t_1}$ (in physics, t often represents time). In machine learning we use the term
41 “flow” by analogy to the above meaning, to highlight the fact that we can create flexible invertible
42 transformations by composing simpler ones; in this sense, the index t is analogous to the number i of
43 transformations in $\mathbf{f}_i \circ \dots \circ \mathbf{f}_1$.

47

23.1.2 How to train a flow model

There are two common applications of normalizing flows. The first one is density estimation of observed data, which is achieved by fitting $p_{\theta}(\mathbf{x})$ to the data and using it as an estimate of the data density, potentially followed by generating new data from $p_{\theta}(\mathbf{x})$. The second one is variational inference, which involves sampling from and evaluating a variational posterior $q_{\theta}(\mathbf{z}|\mathbf{x})$ parameterized by the flow model. As we will see below, these applications optimize different objectives and impose different computational constraints on the flow model.

23.1.2.1 Density estimation

Density estimation requires maximizing the likelihood function in Equation (23.2). This requires that we can efficiently evaluate the inverse flow $\mathbf{u} = \mathbf{f}^{-1}(\mathbf{x})$ and its Jacobian determinant $\det \mathbf{J}(\mathbf{f}^{-1})(\mathbf{x})$ for any given \mathbf{x} . After optimizing the model, we can optionally use it to generate new data. To sample new points, we require that the forwards mapping \mathbf{f} be tractable.

23.1.2.2 Variational inference

Normalizing flows are commonly used for variational inference to parameterize the approximate posterior distribution in latent variable models, as discussed in Section 10.4.3. Consider a latent variable model with continuous latent variables \mathbf{z} and observable variables \mathbf{x} . For simplicity, we consider the model parameters to be fixed as we are interested in approximating the true posterior $p^*(\mathbf{z}|\mathbf{x})$ with a normalizing flow $q_{\theta}(\mathbf{z}|\mathbf{x})$.¹ As discussed in Section 10.1.1.2, the variational parameters are trained by maximizing the evidence lower bound (ELBO), given by

$$L(\boldsymbol{\theta}) = \mathbb{E}_{q_{\theta}(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z}) + \log p(\mathbf{z}) - \log q_{\theta}(\mathbf{z}|\mathbf{x})] \quad (23.5)$$

When viewing the ELBO as a function of $\boldsymbol{\theta}$, it can be simplified as follows (note we drop the dependency on \mathbf{x} for simplicity):

$$L(\boldsymbol{\theta}) = \mathbb{E}_{q_{\theta}(\mathbf{z})} [\ell_{\boldsymbol{\theta}}(\mathbf{z})]. \quad (23.6)$$

Let $q_{\theta}(\mathbf{z})$ denote a normalizing flow with base distribution $q(\mathbf{u})$ and transformation $\mathbf{z} = f_{\boldsymbol{\theta}}(\mathbf{u})$. Then the reparameterization trick (Section 6.3.5) allows us to optimize the parameters using stochastic gradients. To achieve this, we first write the expectation with respect to the base distribution:

$$L(\boldsymbol{\theta}) = \mathbb{E}_{q_{\theta}(\mathbf{z})} [\ell_{\boldsymbol{\theta}}(\mathbf{z})] = \mathbb{E}_{q(\mathbf{u})} [\ell_{\boldsymbol{\theta}}(f_{\boldsymbol{\theta}}(\mathbf{u}))]. \quad (23.7)$$

Then, since the base distribution does not depend on $\boldsymbol{\theta}$, we can obtain stochastic gradients as follows:

$$\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) = \mathbb{E}_{q(\mathbf{u})} [\nabla_{\boldsymbol{\theta}} \ell_{\boldsymbol{\theta}}(f_{\boldsymbol{\theta}}(\mathbf{u}))] \approx \frac{1}{N} \sum_{n=1}^N \nabla_{\boldsymbol{\theta}} \ell_{\boldsymbol{\theta}}(f_{\boldsymbol{\theta}}(\mathbf{u}_n)), \quad (23.8)$$

where $\{\mathbf{u}_n\}_{n=1}^N$ are samples from $q(\mathbf{u})$.

¹ We denote the parameters of the variational posterior by $\boldsymbol{\theta}$ here, which should not be confused with the model parameters which are also typically denoted by $\boldsymbol{\theta}$ elsewhere.

As we can see, in order to optimize this objective, we need to be able to efficiently sample from $q_{\theta}(\mathbf{z}|\mathbf{x})$ and evaluate the probability density of these samples during optimization. (See Section 23.2.4.3 for details on how to do this.) This is contrast to the MLE approach in Section 23.1.2.1, which requires that we be able to compute efficiently the density of arbitrary training datapoints, but it does not require samples during optimization.

23.2 Constructing flows

In this section, we discuss how to compute various kinds of flows that are invertible by design and have efficiently computable Jacobian determinants.

23.2.1 Affine flows

A simple choice is to use an affine transformation $\mathbf{x} = \mathbf{f}(\mathbf{u}) = \mathbf{A}\mathbf{u} + \mathbf{b}$. This is a bijection if and only if \mathbf{A} is an invertible square matrix. The Jacobian determinant of \mathbf{f} is $\det \mathbf{A}$, and its inverse is $\mathbf{u} = \mathbf{f}^{-1}(\mathbf{x}) = \mathbf{A}^{-1}(\mathbf{x} - \mathbf{b})$. A flow consisting of affine bijections is called an **affine flow**, or a **linear flow** if we ignore \mathbf{b} .

On their own, affine flows are limited in their expressive power. For example, suppose the base distribution is Gaussian, $p(\mathbf{u}) = \mathcal{N}(\mathbf{u}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Then the pushforward distribution after an affine bijection is still Gaussian, $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$. However, affine bijections are useful building blocks when composed with the non-affine bijections we discuss later, as they encourage “mixing” of dimensions through the flow.

For practical reasons, we need to ensure the Jacobian determinant and the inverse of the flow are fast to compute. In general, computing $\det \mathbf{A}$ and \mathbf{A}^{-1} explicitly takes $O(D^3)$ time. To reduce the cost, we can add structure to \mathbf{A} . If \mathbf{A} is diagonal, the cost becomes $O(D)$. If \mathbf{A} is triangular, the Jacobian determinant is the product of the diagonal elements, so it takes $O(D)$ time; inverting the flow requires solving the triangular system $\mathbf{A}\mathbf{u} = \mathbf{x} - \mathbf{b}$, which can be done with backsubstitution in $O(D^2)$ time.

The result of a triangular transformation depends on the ordering of the dimensions. To reduce sensitivity to this, and to encourage “mixing” of dimensions, we can multiply \mathbf{A} with a permutation matrix, which has an absolute determinant of 1. We often use a permutation that reverses the indices at each layer or that randomly shuffles them. However, usually the permutation at each layer is fixed rather than learned.

For spatially structured data (such as images), we can define \mathbf{A} to be a convolution matrix. For example, GLOW [KD18b] uses 1×1 convolution; this is equivalent to pointwise linear transformation across feature dimensions, but regular convolution across spatial dimensions. Two more general methods for modeling $d \times d$ convolutions are presented in [HBW19], one based on stacking autoregressive convolutions, and the other on carrying out the convolution in the Fourier domain.

40

23.2.2 Elementwise flows

Let $h : \mathbb{R} \rightarrow \mathbb{R}$ be a scalar-valued bijection. We can create a vector-valued bijection $\mathbf{f} : \mathbb{R}^D \rightarrow \mathbb{R}^D$ by applying h elementwise, that is, $\mathbf{f}(\mathbf{u}) = (h(u_1), \dots, h(u_D))$. The function \mathbf{f} is invertible, and its Jacobian determinant is given by $\prod_{i=1}^D \frac{dh}{du_i}$. A flow composed of such bijections is known as an **elementwise flow**.

47

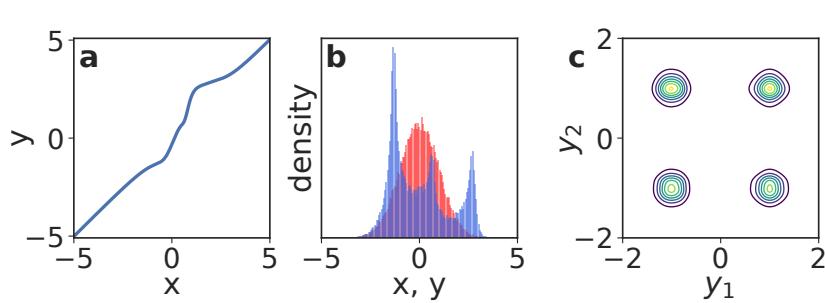


Figure 23.1: Non-linear squared flow (NLSq). Left: an invertible mapping consisting of 4 NLSq layers. Middle: red is the base distribution (Gaussian), blue is the distribution induced by the mapping on the left. Right: density of a 5-layer autoregressive flow using NLSq transformations and a Gaussian base density, trained on a mixture of 4 Gaussians. From Figure 5 of [ZR19b]. Used with kind permission of Zachary Ziegler.

On their own, elementwise flows are limited, since they do not model dependencies between the elements. However, they are useful building blocks for more complex flows, such as coupling flows (Section 23.2.3) and autoregressive flows (Section 23.2.4), as we will see later. In this section, we discuss techniques for constructing scalar-valued bijections $h : \mathbb{R} \rightarrow \mathbb{R}$ for use in elementwise flows.

23.2.2.1 Affine scalar bijection

An **affine scalar bijection** has the form $h(u; \theta) = au + b$, where $\theta = (a, b) \in \mathbb{R}^2$. (This is a scalar version of an affine flow.) Its derivative $\frac{dh}{du}$ is equal to a . It is invertible if and only if $a \neq 0$. In practice, we often parameterize a to be positive, for example by making it the exponential or the softplus of an unconstrained parameter. When $a = 1$, $h(u; \theta) = u + b$ is often called an **additive scalar bijection**.

23.2.2.2 Higher-order perturbations

The affine scalar bijection is simple to use, but limited. We can make it more flexible by adding higher-order perturbations, under the constraint that invertibility is preserved. For example, Ziegler and Rush [ZR19b] propose the following, which they term **non-linear squared flow**:

$$h(u; \theta) = au + b + \frac{c}{1 + (du + e)^2}, \quad (23.9)$$

where $\theta = (a, b, c, d, e) \in \mathbb{R}^5$. When $c = 0$, this reduces to the affine case. When $c \neq 0$, it adds an inverse-quadratic perturbation, which can induce multimodality as shown in Figure 23.1. Under the constraints $a > \frac{9}{8\sqrt{3}}cd$ and $d > 0$ the function becomes invertible, and its inverse can be computed analytically by solving a quadratic polynomial.

23.2.2.3 Combinations of strictly monotonic scalar functions

A strictly monotonic scalar function is one that is always increasing (has positive derivative everywhere) or always decreasing (has negative derivative everywhere). Such functions are invertible. Many

1 activation functions, such as the logistic sigmoid $\sigma(u) = 1/(1 + \exp(-u))$, are strictly monotonic.
2

3 Using such activation functions as a starting point, we can build more flexible monotonic functions
4 via **conical combination** (linear combination with positive coefficients) and function composition.
5 Suppose h_1, \dots, h_K are strictly increasing; then the following are also strictly increasing:

- 6 • $a_1 h_1 + \dots + a_K h_K + b$ with $a_k > 0$ (conical combination with a bias),
7
8 • $h_1 \circ \dots \circ h_K$ (function composition).

9 By repeating the above two constructions, we can build arbitrarily complex increasing functions. For
10 example, a composition of conical combinations of logistic sigmoids is just an MLP where all weights
11 are positive [Hua+18b].

12 The derivative of such a scalar bijection can be computed by repeatedly applying the chain rule,
13 and in practice can be done with automatic differentiation. However, the inverse is not typically
14 computable in closed form. In practice we can compute the inverse using bisection search, since the
15 function is monotonic.

16 23.2.2.4 Scalar bijections from integration

17 A simple way to ensure a scalar function is strictly monotonic is to constrain its derivative to be
18 positive. Let $h' = \frac{dh}{du}$ be this derivative. Wehenkel and Louppe [WL19] directly parameterize h' with
19 a neural network whose output is made positive via an ELU activation function shifted up by 1.
20 They then integrate the derivative numerically to get the bijection:
21

$$\begin{aligned} \text{24} \quad h(u) &= \int_0^u h'(t)dt + b, \\ \text{25} \end{aligned} \tag{23.10}$$

26 where b is a bias. They call this approach **unconstrained monotonic neural networks**.

27 The above integral is generally not computable in closed form. It can be, however, if h' is
28 constrained appropriately. For example, Jaini, Selby, and Yu [JSY19] take h' to be a sum of K
29 squared polynomials of degree L :

$$\begin{aligned} \text{31} \quad h'(u) &= \sum_{k=1}^K \left(\sum_{\ell=0}^L a_{k\ell} u^\ell \right)^2. \\ \text{32} \\ \text{33} \\ \text{34} \end{aligned} \tag{23.11}$$

35 This makes h' a non-negative polynomial of degree $2L$. The integral is analytically tractable, and
36 makes h an increasing polynomial of degree $2L + 1$. For $L = 0$, h' is constant, so h reduces to an
37 affine scalar bijection.

38 In these approaches, the derivative of the bijection can just be read off. However, the inverse is not
39 analytically computable in general. In practice, we can use bisection search to compute the inverse
40 numerically.

41

42 23.2.2.5 Splines

43

44 Another way to construct monotonic scalar functions is using **splines**. These are piecewise-polynomial
45 or piecewise-rational functions, parameterized in terms of $K + 1$ **knots** (u_k, x_k) through which the
46 spline passes. That is, we set $h(u_k) = x_k$, and define h on the interval (u_{k-1}, u_k) by interpolating
47

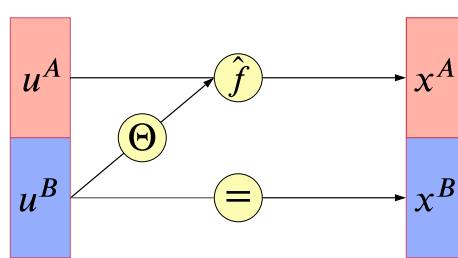


Figure 23.2: Illustration of a coupling layer $\mathbf{x} = \hat{\mathbf{f}}(\mathbf{u})$. A bijection, with parameters determined by \mathbf{u}^B , is applied to \mathbf{u}^A to generate \mathbf{x}^A ; meanwhile $\mathbf{x}^B = \mathbf{u}^B$ is passed through unchanged, so the mapping can be inverted. From Figure 3 of [KPB19]. Used with kind permission of Ivan Kobyzhev.

from x_{k-1} to x_k with a polynomial or rational function (ratio of two polynomials). By increasing the number of knots we can create arbitrarily flexible monotonic functions.

Different ways to interpolate between knots give different types of spline. A simple choice is to interpolate linearly [Mül+19a]. However, this makes the derivative discontinuous at the knots. Interpolating with quadratic polynomials [Mül+19a] gives enough flexibility to make the derivative continuous. Interpolating with cubic polynomials [Dur+19], ratios of linear polynomials [DEL20], or ratios of quadratic polynomials [DBP19] allows the derivatives at the knots to be arbitrary parameters.

The spline is strictly increasing if we take $u_{k-1} < u_k$, $x_{k-1} < x_k$, and make sure the interpolation between knots is itself increasing. Depending on the flexibility on the interpolating function, more than one interpolation may exist; in practice we choose one that is guaranteed to be always increasing (see references above for details).

An advantage of splines is that they can be inverted analytically if the interpolating functions only contain low-degree polynomials. In this case, we compute $u = h^{-1}(x)$ as follows: first, we use binary search to locate the interval (x_{k-1}, x_k) in which x lies; then, we analytically solve the resulting low-degree polynomial for u .

23.2.3 Coupling flows

In this section we describe coupling flows, which allow us to model dependencies between dimensions using arbitrary non-linear functions (such as deep neural networks). Consider a partition of the input $\mathbf{u} \in \mathbb{R}^D$ into two subspaces, $(\mathbf{u}^A, \mathbf{u}^B) \in \mathbb{R}^d \times \mathbb{R}^{D-d}$, where d is an integer between 1 and $D-1$. Assume a bijection $\hat{\mathbf{f}}(\cdot; \boldsymbol{\theta}) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ parameterized by $\boldsymbol{\theta}$ and acting on the subspace \mathbb{R}^d . We define the function $\mathbf{f} : \mathbb{R}^D \rightarrow \mathbb{R}^D$ given by $\mathbf{x} = \mathbf{f}(\mathbf{u})$ as follows:

$$\mathbf{x}^A = \hat{\mathbf{f}}(\mathbf{u}^A; \Theta(\mathbf{u}^B)) \tag{23.12}$$

$$\mathbf{x}^B = \mathbf{u}^B. \tag{23.13}$$

See Figure 23.2 for an illustration. The function \mathbf{f} is called a **coupling layer** [DKB15; DSDB17], because it “couples” \mathbf{u}^A and \mathbf{u}^B together though $\hat{\mathbf{f}}$ and Θ . We refer to flows consisting of coupling layers as **coupling flows**.

¹ The parameters of $\hat{\mathbf{f}}$ are computed by $\boldsymbol{\theta} = \Theta(\mathbf{u}^B)$, where Θ is an *arbitrary* function called the
² **conditioner**. Unlike affine flows, which mix dimensions linearly, and elementwise flows, which do
³ not mix dimensions at all, coupling flows can mix dimensions with a flexible non-linear conditioner Θ .
⁴ In practice we often implement Θ as a deep neural network; any architecture can be used, including
⁵ MLPs, CNNs, ResNets, etc.

⁶ The coupling layer \mathbf{f} is *invertible*, and its inverse is given by $\mathbf{u} = \mathbf{f}^{-1}(\mathbf{x})$, where

$$\mathbf{u}^A = \hat{\mathbf{f}}^{-1}(\mathbf{x}^A; \Theta(\mathbf{x}^B)) \quad (23.14)$$

$$\mathbf{u}^B = \mathbf{x}^B. \quad (23.15)$$

¹¹ That is, \mathbf{f}^{-1} is given by simply replacing $\hat{\mathbf{f}}$ with $\hat{\mathbf{f}}^{-1}$. Because \mathbf{x}^B does not depend on \mathbf{u}^A , the
¹² Jacobian of \mathbf{f} is block triangular:
¹³

$$\mathbf{J}(\mathbf{f}) = \begin{pmatrix} \partial \mathbf{x}^A / \partial \mathbf{u}^A & \partial \mathbf{x}^A / \partial \mathbf{u}^B \\ \partial \mathbf{x}^B / \partial \mathbf{u}^A & \partial \mathbf{x}^B / \partial \mathbf{u}^B \end{pmatrix} = \begin{pmatrix} \mathbf{J}(\hat{\mathbf{f}}) & \partial \mathbf{x}^A / \partial \mathbf{u}^B \\ \mathbf{0} & \mathbf{I} \end{pmatrix}. \quad (23.16)$$

¹⁴ Thus, $\det \mathbf{J}(\mathbf{f})$ is equal to $\det \mathbf{J}(\hat{\mathbf{f}})$.

¹⁵ We often define $\hat{\mathbf{f}}$ to be an elementwise bijection, so that $\hat{\mathbf{f}}^{-1}$ and $\det \mathbf{J}(\hat{\mathbf{f}})$ are easy to compute.
¹⁶ That is, we define:

$$\hat{\mathbf{f}}(\mathbf{u}^A; \boldsymbol{\theta}) = (h(u_1^A; \boldsymbol{\theta}_1), \dots, h(u_d^A; \boldsymbol{\theta}_d)), \quad (23.17)$$

¹⁷ where $h(\cdot; \boldsymbol{\theta}_i)$ is a scalar bijection parameterized by $\boldsymbol{\theta}_i$. Any of the scalar bijections described in
¹⁸ Section 23.2.2 can be used here. For example, $h(\cdot; \boldsymbol{\theta}_i)$ can be an affine bijection with $\boldsymbol{\theta}_i$ its scale and
¹⁹ shift parameters (Section 23.2.2.1); or it can be a monotonic MLP with $\boldsymbol{\theta}_i$ its weights and biases
²⁰ (Section 23.2.2.3); or it can be a monotonic spline with $\boldsymbol{\theta}_i$ its knot coordinates (Section 23.2.2.5).

²¹ There are many ways to define the partition of \mathbf{u} into $(\mathbf{u}^A, \mathbf{u}^B)$. A simple way is just to partition
²² \mathbf{u} into two halves. We can also exploit spatial structure in the partitioning. For example, if \mathbf{u} is an
²³ image, we can partition its pixels using a “checkerboard” pattern, where pixels in “black squares” are
²⁴ in \mathbf{u}^A and pixels in “white squares” are in \mathbf{u}^B [DSDB17]. Since only part of the input is transformed
²⁵ by each coupling layer, in practice we typically employ different partitions along a coupling flow, to
²⁶ ensure all variables get transformed and are given the opportunity to interact.

²⁷ Finally, if $\hat{\mathbf{f}}$ is an elementwise bijection, we can implement arbitrary partitions easily using a binary
²⁸ mask \mathbf{b} as follows:

$$\mathbf{x} = \mathbf{b} \odot \mathbf{u} + (1 - \mathbf{b}) \odot \hat{\mathbf{f}}(\mathbf{u}; \Theta(\mathbf{b} \odot \mathbf{u})), \quad (23.18)$$

²⁹ where \odot denotes elementwise multiplication. A value of 0 in \mathbf{b} indicates that the corresponding
³⁰ element in \mathbf{u} is transformed (belongs to \mathbf{u}^A); a value of 1 indicates that it remains unchanged (belongs
³¹ to \mathbf{u}^B).

³² As an example, we fit a masked coupling flow, created from piecewise rational quadratic splines, to
³³ the two moons dataset. Samples from each layer of the fitted model are shown in Figure 23.3.

³⁴

³⁵ 23.2.4 Autoregressive flows

³⁶ In this section we discuss **autoregressive flows**, which are flows composed of autoregressive bijections.
³⁷ Like coupling flows, autoregressive flows allow us to model dependencies between variables with
³⁸ arbitrary non-linear functions, such as deep neural networks.

³⁹

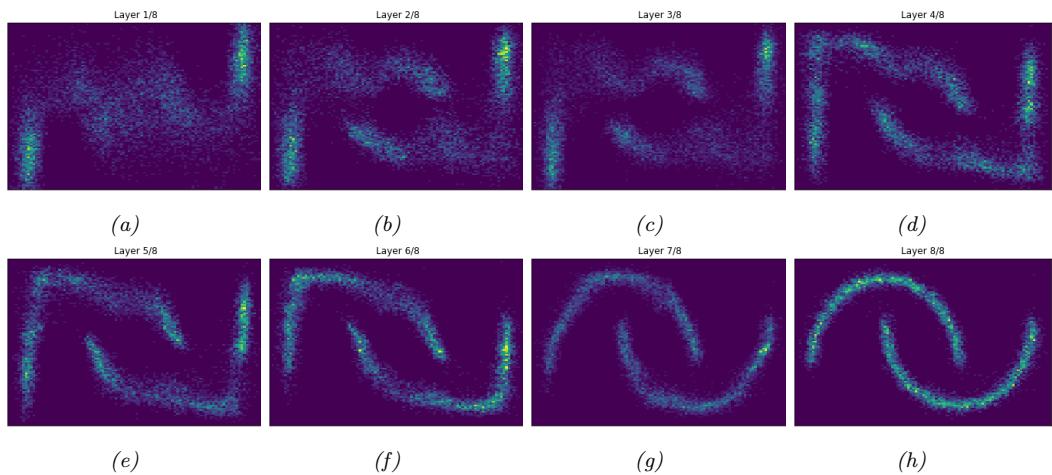


Figure 23.3: (a) Two moons dataset. (b) Samples from a normalizing flow fit to this dataset. Generated by [two_moons_nsf_normalizing_flow.ipynb](#).

Suppose the input \mathbf{u} contains D scalar elements, that is, $\mathbf{u} = (u_1, \dots, u_D) \in \mathbb{R}^D$. We define an **autoregressive bijection** $\mathbf{f} : \mathbb{R}^D \rightarrow \mathbb{R}^D$, its output denoted by $\mathbf{x} = (x_1, \dots, x_D) \in \mathbb{R}^D$, as follows:

$$x_i = h(u_i; \Theta_i(\mathbf{x}_{1:i-1})), \quad i = 1, \dots, D. \quad (23.19)$$

Each output x_i depends on the corresponding input u_i and all previous outputs $\mathbf{x}_{1:i-1} = (x_1, \dots, x_{i-1})$. The function $h(\cdot; \boldsymbol{\theta}) : \mathbb{R} \rightarrow \mathbb{R}$ is a scalar bijection (for example, one of those described in Section 23.2.2), and is parameterized by $\boldsymbol{\theta}$. The function Θ_i is a conditioner that outputs the parameters $\boldsymbol{\theta}_i$ that yield x_i , given all previous outputs $\mathbf{x}_{1:i-1}$. Like in coupling flows, Θ_i can be an arbitrary non-linear function, and is often parameterized as a deep neural network.

Because h is invertible, \mathbf{f} is also invertible, and its inverse is given by:

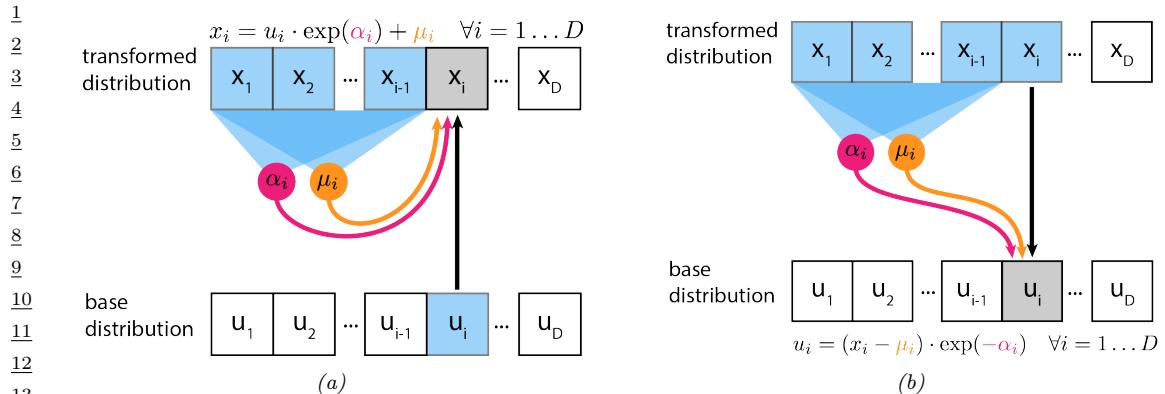
$$u_i = h^{-1}(x_i; \Theta_i(\mathbf{x}_{1:i-1})), \quad i = 1, \dots, D. \quad (23.20)$$

An important property of \mathbf{f} is that each output x_i depends on $\mathbf{u}_{1:i} = (u_1, \dots, u_i)$, but not on $\mathbf{u}_{i+1:D} = (u_{i+1}, \dots, u_D)$; as a result, the partial derivative $\partial x_i / \partial u_j$ is identically zero whenever $j > i$. Therefore, the Jacobian matrix $\mathbf{J}(\mathbf{f})$ is triangular, and its determinant is simply the product of its diagonal entries:

$$\det \mathbf{J}(\mathbf{f}) = \prod_{i=1}^D \frac{\partial x_i}{\partial u_i} = \prod_{i=1}^D \frac{dh}{du_i}. \quad (23.21)$$

In other words, the autoregressive structure of \mathbf{f} leads to a Jacobian determinant that can be computed efficiently in $O(D)$ time.

Although invertible, autoregressive bijections are computationally asymmetric: evaluating \mathbf{f} is inherently sequential, whereas evaluating \mathbf{f}^{-1} is inherently parallel. That is because we need $\mathbf{x}_{1:i-1}$ to



¹⁴ Figure 23.4: (a) Affine autoregressive flow with one layer. In this figure, \mathbf{u} is the input to the flow (sample
¹⁵ from the base distribution) and \mathbf{x} is its output (sample from the transformed distribution). (b) Inverse of the
¹⁶ above. From [Jan18]. Used with kind permission of Eric Jang.

¹⁷

¹⁸

¹⁹ compute x_i ; therefore, computing the components of \mathbf{x} must be done sequentially, by first computing
²⁰ x_1 , then using it to compute x_2 , then using x_1 and x_2 to compute x_3 , and so on. On the other hand,
²¹ computing the inverse can be done in parallel for each u_i , since \mathbf{u} does not appear on the right-hand
²² side of Equation (23.20). Hence, in practice it is often faster to compute \mathbf{f}^{-1} than to compute \mathbf{f} ,
²³ assuming h and h^{-1} have similar computational cost.

²⁴

²⁵ 23.2.4.1 Affine autoregressive flows

²⁶

²⁷ For a concrete example, we can take h to be an affine scalar bijection (Section 23.2.2.1) parameterized
²⁸ by a log scale α and a bias μ . Such autoregressive flows are known as **affine autoregressive flows**.
²⁹ The parameters of the i 'th component, α_i and μ_i , are functions of $\mathbf{x}_{1:i-1}$, so \mathbf{f} takes the following
³⁰ form:

³¹

$$\text{32} \quad x_i = u_i \exp(\alpha_i(\mathbf{x}_{1:i-1})) + \mu_i(\mathbf{x}_{1:i-1}). \quad (23.22)$$

³³

This is illustrated in Figure 23.4(a). We can invert this by

³⁴

$$\text{35} \quad u_i = (x_i - \mu_i(\mathbf{x}_{1:i-1})) \exp(-\alpha_i(\mathbf{x}_{1:i-1})). \quad (23.23)$$

³⁶

This is illustrated in Figure 23.4(b). Finally, we can calculate the log absolute Jacobian determinant
³⁷ by

³⁸

$$\text{39} \quad \log |\det \mathbf{J}(\mathbf{f})| = \log \left| \prod_{i=1}^D \exp(\alpha_i(\mathbf{x}_{1:i-1})) \right| = \sum_{i=1}^D \alpha_i(\mathbf{x}_{1:i-1}). \quad (23.24)$$

⁴⁰

⁴¹ Let us look at an example of an affine autoregressive flow on a 2d density estimation problem.
⁴² Consider an affine autoregressive flow $\mathbf{x} = (x_1, x_2) = \mathbf{f}(\mathbf{u})$, where $\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and \mathbf{f} is a single
⁴³ autoregressive bijection. Since x_1 is an affine transformation of $u_1 \sim \mathcal{N}(0, 1)$, it is Gaussian with
⁴⁴ mean μ_1 and standard deviation $\sigma_1 = \exp \alpha_1$. Similarly, if we consider x_1 fixed, x_2 is an affine
⁴⁵ transformation of $u_2 \sim \mathcal{N}(0, 1)$.

⁴⁶

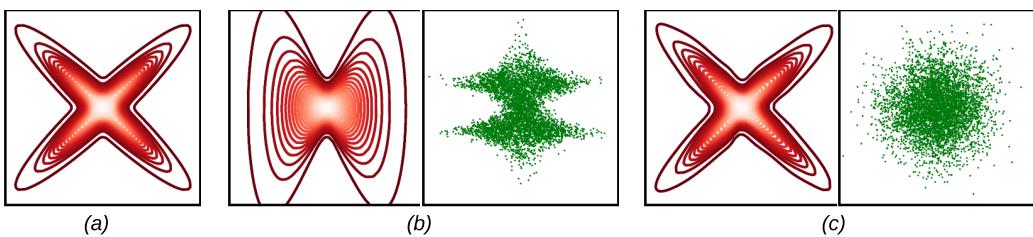


Figure 23.5: Density estimation with affine autoregressive flows, using a Gaussian base distribution. (a) True density. (b) Estimated density using a single autoregressive layer with ordering (x_1, x_2) . On the left (contour plot) we show $p(\mathbf{x})$. On the right (green dots) we show samples of $\mathbf{u} = \mathbf{f}^{-1}(\mathbf{x})$, where \mathbf{x} is sampled from the true density. (c) Same as (b), but using 5 autoregressive layers and reversing the variable ordering after each layer. Adapted from Figure 1 of [PPM17]. Used with kind permission of Iain Murray.

transformation of $u_2 \sim \mathcal{N}(0, 1)$, so it is *conditionally* Gaussian with mean $\mu_2(x_1)$ and standard deviation $\sigma_2(x_1) = \exp \alpha_2(x_1)$. Thus, a single affine autoregressive bijection will always produce a distribution with Gaussian conditionals, that is, a distribution of the following form:

$$p(x_1, x_2) = p(x_1) p(x_2|x_1) = \mathcal{N}(x_1|\mu_1, \sigma_1^2) \mathcal{N}(x_2|\mu_2(x_1), \sigma_2(x_1)^2) \quad (23.25)$$

This result generalizes to an arbitrary number of dimensions D .

A single affine bijection is not very powerful, regardless of how flexible the functions $\alpha_2(x_1)$ and $\mu_2(x_1)$ are. For example, suppose we want to fit the cross-shaped density shown in Figure 23.5(a) with such a flow. The resulting maximum-likelihood fit is shown in Figure 23.5(b). The red contours show the predictive distribution, $\hat{p}(\mathbf{x})$, which clearly fails to capture the true distribution. The green dots show transformed versions of the data samples, $p(\mathbf{u})$; we see that this is far from the Gaussian base distribution.

Fortunately, we can obtain a better fit by composing multiple autoregressive bijections (layers), and reversing the order of the variables after each layer. For example, Figure 23.5(c) shows the results of an affine autoregressive flow with 5 layers applied to the same problem. The red contours show that we have matched the empirical distribution, and the green dots show we have matched the Gaussian base distribution.

Note that another way to obtain a better fit is to replace the affine bijection h with a more flexible one, such as a monotonic MLP (Section 23.2.2.3) or a monotonic spline (Section 23.2.2.5).

23.2.4.2 Masked autoregressive flows

As we have seen, the conditioners Θ_i can be arbitrary non-linear functions. The most straightforward way to parameterize them is separately for each i , for example by using D separate neural networks. However, this can be parameter-inefficient for large D .

In practice, we often share parameters between conditioners by combining them into a single model Θ that takes in \mathbf{x} and outputs $(\theta_1, \dots, \theta_D)$. For the bijection to remain autoregressive, we must constrain Θ so that θ_i depends only on $\mathbf{x}_{1:i-1}$ and not on $\mathbf{x}_{i:D}$. One way to achieve this is to start with an arbitrary neural network (an MLP, a CNN, a ResNet, etc.), and drop connections (for example, by zeroing out weights) until θ_i is only a function of $\mathbf{x}_{1:i-1}$.

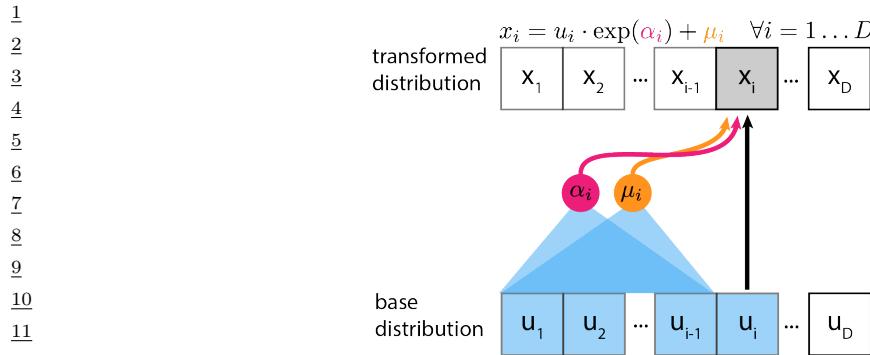


Figure 23.6: Inverse autoregressive flow that uses affine scalar bijections. In this figure, \mathbf{u} is the input to the flow (sample from the base distribution) and \mathbf{x} is its output (sample from the transformed distribution) From [Jan18]. Used with kind permission of Eric Jang.

An example of this approach is the **masked autoregressive flow (MAF)** model of [PPM17]. This model is an affine autoregressive flow combined with permutation layers, as we described in Section 23.2.4.1. MAF implements the combined conditioner Θ as follows: it starts with an MLP, and then multiplies (elementwise) the weight matrix of each layer with a binary mask of the same size (different masks are used for different layers). The masks are constructed using the method of [Ger+15]. This ensures that all computational paths from x_j to θ_i are zeroed out whenever $j \geq i$, effectively making θ_i only a function of $x_{1:i-1}$. Still, evaluating the masked conditioner Θ has the same computational cost as evaluating the original (unmasked) MLP.

The key advantage of MAF (and of related models) is that, given \mathbf{x} , all parameters $(\theta_1, \dots, \theta_D)$ can be computed efficiently with one neural network evaluation, so the computation of the inverse \mathbf{f}^{-1} is fast. Thus, we can efficiently evaluate the probability density of the flow model for arbitrary datapoints. However, in order to compute \mathbf{f} , the conditioner Θ must be called a total of D times, since not all entries of \mathbf{x} are available to start with. Thus, generating new samples from the flow is D times more expensive than evaluating its probability density function. This makes MAF suitable for density estimation, but less so for data generation.

23.2.4.3 Inverse autoregressive flows

As we have seen, the parameters θ_i that yield the i 'th output x_i are functions of the previous outputs $\mathbf{x}_{1:i-1}$. This ensures that the Jacobian $\mathbf{J}(\mathbf{f})$ is triangular, and so its determinant is efficient to compute.

However, there is another possibility: we can make θ_i a function of the previous *inputs* instead, that is, a function of $\mathbf{u}_{1:i-1}$. This leads to the following bijection, which is known as **inverse autoregressive**:

$$x_i = h(u_i; \Theta_i(\mathbf{u}_{1:i-1})), \quad i = 1, \dots, D. \quad (23.26)$$

Like its autoregressive counterpart, this bijection has a triangular Jacobian whose determinant is also given by $\det \mathbf{J}(\mathbf{f}) = \prod_{i=1}^D \frac{dh}{du_i}$. Figure 23.6 illustrates an inverse autoregressive flow, for the case where h is affine.

47

To see why this bijection is called “inverse autoregressive”, compare Equation (23.26) with Equation (23.20). The two formulas differ only notationally: we can get from one to the other by swapping \mathbf{u} with \mathbf{x} and h with h^{-1} . In other words, the inverse autoregressive bijection corresponds to a direct parameterization of the inverse of an autoregressive bijection.

Since inverse autoregressive bijections swap the forwards and inverse directions of their autoregressive counterparts, they also swap their computational properties. This means that the forward direction \mathbf{f} of an inverse autoregressive flow is inherently parallel and therefore fast, whereas its inverse direction \mathbf{f}^{-1} is inherently sequential and therefore slow.

An example of an inverse autoregressive flow is their namesake **IAF** model of [Kin+16]. IAF uses affine scalar bijections, masked conditioners, and permutation layers, so it is precisely the inverse of the MAF model described in Section 23.2.4.2. Using IAF, we can generate \mathbf{u} in parallel from the base distribution (using, for example, a diagonal Gaussian), and then sample each element of \mathbf{x} in parallel. However, evaluating $p(\mathbf{x})$ for an arbitrary datapoint \mathbf{x} is slow, because we have to evaluate each element of \mathbf{u} sequentially. Fortunately, evaluating the likelihood of samples generated from IAF (as opposed to externally provided samples) incurs no additional cost, since in this case the u_i terms will already have been computed.

Although not so suitable for density estimation or maximum-likelihood training, IAFs are well-suited for parameterizing variational posteriors in variational inference. This is because in order to estimate the variational lower bound (ELBO), we only need samples from the variational posterior and their associated probability densities, both of which are efficient to obtain. See Section 23.1.2.2 for details.

Another useful application of IAFs is training them to mimic models whose probability density is fast to evaluate but which are slow to sample from. A notable example is the **parallel wavenet** model of [Oor+18]. This model is an IAF p_s that is trained to mimic a pretrained wavenet model p_t by minimizing the KL divergence $D_{\text{KL}}(p_s \parallel p_t)$. This KL can be easily estimated by first sampling from p_s and then evaluating $\log p_s$ and $\log p_t$ at those samples, operations which are all efficient for these models. After training, we obtain an IAF that can generate audio of similar quality as the original wavenet, but can do so much faster.

23.2.4.4 Connection with autoregressive models

Autoregressive flows can be thought of as generalizing autoregressive models of continuous random variables, discussed in Section 22.1. Specifically, any continuous autoregressive model can be reparameterized as a one-layer autoregressive flow, as we describe below.

Consider a general autoregressive model over a continuous random variable $\mathbf{x} = (x_1, \dots, x_D) \in \mathbb{R}^D$ written as

$$p(\mathbf{x}) = \prod_{i=1}^D p_i(x_i | \boldsymbol{\theta}_i) \quad \text{where} \quad \boldsymbol{\theta}_i = \Theta_i(\mathbf{x}_{1:i-1}). \quad (23.27)$$

In the above expression, $p_i(x_i | \boldsymbol{\theta}_i)$ is the i 'th conditional distribution of the autoregressive model, whose parameters $\boldsymbol{\theta}_i$ are arbitrary functions of the previous variables $\mathbf{x}_{1:i-1}$. For example, $p_i(x_i | \boldsymbol{\theta}_i)$ can be a mixture of one-dimensional Gaussian distributions, with $\boldsymbol{\theta}_i$ representing the collection of its means, variances, and mixing coefficients.

Now consider sampling a vector \mathbf{x} from the autoregressive model, which can be done by sampling

1 one element at a time as follows:

2

$$\underline{3} \quad x_i \sim p_i(x_i | \Theta_i(\mathbf{x}_{1:i-1})) \quad \text{for } i = 1, \dots, D. \quad (23.28)$$

4

5 Each conditional can be sampled from using inverse transform sampling (Section 11.3.1). Let $U(0, 1)$
6 be the uniform distribution on the interval $[0, 1]$, and let $\text{CDF}_i(x_i | \theta_i)$ be the cumulative distribution
7 function of the i 'th conditional. Sampling can be written as:

8

$$\underline{9} \quad x_i = \text{CDF}_i^{-1}(u_i | \Theta_i(\mathbf{x}_{1:i-1})) \quad \text{where } u_i \sim U(0, 1). \quad (23.29)$$

10

11 Comparing the above expression with the definition of an autoregressive bijection in Equation (23.19),
12 we see that the autoregressive model has been expressed as a one-layer autoregressive flow whose base
13 distribution is uniform on $[0, 1]^D$ and whose scalar bijections correspond to the inverse conditional
14 cdf's. Viewing autoregressive models as flows this way has an important advantage, namely that it
15 allows us to increase the flexibility of an autoregressive model by composing multiple instances of it
16 in a flow, without sacrificing the overall tractability.

17

18 23.2.5 Residual flows

19 A residual network is a composition of **residual connections**, which are functions of the form
20 $\mathbf{f}(\mathbf{u}) = \mathbf{u} + \mathbf{F}(\mathbf{u})$. The function $\mathbf{F} : \mathbb{R}^D \rightarrow \mathbb{R}^D$ is called the **residual block**, and it computes the
21 difference between the output and the input, $\mathbf{f}(\mathbf{u}) - \mathbf{u}$.

22 Under certain conditions on \mathbf{F} , the residual connection \mathbf{f} becomes invertible. We will refer to flows
23 composed of invertible residual connections as **residual flows**. In the following, we describe two
24 ways the residual block \mathbf{F} can be constrained so that the residual connection \mathbf{f} is invertible.

25

26 23.2.5.1 Contractive residual blocks

27 One way to ensure the residual connection is invertible is to choose the residual block to be a
28 contraction. A contraction is a function \mathbf{F} whose Lipschitz constant is less than 1; that is, there
29 exists $0 \leq L < 1$ such that for all \mathbf{u}_1 and \mathbf{u}_2 we have:

30

$$\underline{31} \quad \|\mathbf{F}(\mathbf{u}_1) - \mathbf{F}(\mathbf{u}_2)\| \leq L \|\mathbf{u}_1 - \mathbf{u}_2\|. \quad (23.30)$$

32

33 The invertibility of $\mathbf{f}(\mathbf{u}) = \mathbf{u} + \mathbf{F}(\mathbf{u})$ can be shown as follows. Consider the mapping $\mathbf{g}(\mathbf{u}) =$
34 $\mathbf{x} - \mathbf{F}(\mathbf{u})$. Because \mathbf{F} is a contraction, \mathbf{g} is also a contraction. So, by Banach's fixed-point theorem,
35 \mathbf{g} has a unique fixed point \mathbf{u}_* . Hence we have

36

$$\underline{37} \quad \mathbf{u}_* = \mathbf{x} - \mathbf{F}(\mathbf{u}_*) \quad (23.31)$$

38

$$\underline{39} \quad \Rightarrow \quad \mathbf{u}_* + \mathbf{F}(\mathbf{u}_*) = \mathbf{x} \quad (23.32)$$

$$\underline{40} \quad \Rightarrow \quad \mathbf{f}(\mathbf{u}_*) = \mathbf{x}. \quad (23.33)$$

41

42 Because \mathbf{u}_* is unique, it follows that $\mathbf{u}_* = \mathbf{f}^{-1}(\mathbf{x})$.

43 An example of a residual flow with contractive residual blocks is the **iResNet** model of [Beh+19].

44 The residual blocks of iResNet are convolutional neural networks, that is, compositions of convolutional
45 layers with non-linear activation functions. Because the Lipschitz constant of a composition is less or
46 equal to the product of the Lipschitz constants of the individual functions, it is enough to ensure the

47

convolutions are contractive, and to use increasing activation functions with slope less or equal to 1. The iResNet model ensures the convolutions are contractive by applying spectral normalization to their weights [Miy+18a].

In general, there is no analytical expression for the inverse \mathbf{f}^{-1} . However, we can approximate $\mathbf{f}^{-1}(\mathbf{x})$ using the following iterative procedure:

$$\mathbf{u}_n = \mathbf{g}(\mathbf{u}_{n-1}) = \mathbf{x} - \mathbf{F}(\mathbf{u}_{n-1}). \quad (23.34)$$

Banach's fixed-point theorem guarantees that the sequence $\mathbf{u}_0, \mathbf{u}_1, \mathbf{u}_2, \dots$ will converge to $\mathbf{u}_* = \mathbf{f}^{-1}(\mathbf{x})$ for any choice of \mathbf{u}_0 , and it will do so at a rate of $O(L^n)$, where L is the Lipschitz constant of \mathbf{g} (which is the same as the Lipschitz constant of \mathbf{F}). In practice, it is convenient to choose $\mathbf{u}_0 = \mathbf{x}$.

In addition, there is no analytical expression for the Jacobian determinant, whose exact computation costs $O(D^3)$. However, there is a computationally efficient stochastic estimator of the log Jacobian determinant. The idea is to express the log Jacobian determinant as a power series. Using the fact that $\mathbf{f}(\mathbf{x}) = \mathbf{x} + \mathbf{F}(\mathbf{x})$, we have

$$\log |\det \mathbf{J}(\mathbf{f})| = \log |\det(\mathbf{I} + \mathbf{J}(\mathbf{F}))| = \sum_{k=1}^{\infty} \frac{(-1)^{k+1}}{k} \text{tr}[\mathbf{J}(\mathbf{F})^k]. \quad (23.35)$$

This power series converges when the matrix norm of $\mathbf{J}(\mathbf{F})$ is less than 1, which here is guaranteed exactly because \mathbf{F} is a contraction. The trace of $\mathbf{J}(\mathbf{F})^k$ can be efficiently approximated using Jacobian-vector products via the **Hutchinson trace estimator** [Ski89; Hut89; Mey+21]:

$$\text{tr}[\mathbf{J}(\mathbf{F})^k] \approx \mathbf{v}^\top \mathbf{J}(\mathbf{F})^k \mathbf{v}, \quad (23.36)$$

where \mathbf{v} is a sample from a distribution with zero mean and unit covariance, such as $\mathcal{N}(\mathbf{0}, \mathbf{I})$. Finally, the infinite series can be approximated by a finite one either by truncation [Beh+19], which unfortunately yields a biased estimator, or by employing the **Russian roulette estimator** [Che+19], which is unbiased.

23.2.5.2 Residual blocks with low-rank Jacobian

There is an efficient way of computing the determinant of a matrix which is a low-rank perturbation of an identity matrix. Suppose \mathbf{A} and \mathbf{B} are matrices, where \mathbf{A} is $D \times M$ and \mathbf{B} is $M \times D$. The following formula is known as the **Weinstein-Aronszajn identity**², and is a special case of the more general **matrix determinant lemma**:

$$\det(\mathbf{I}_D + \mathbf{AB}) = \det(\mathbf{I}_M + \mathbf{BA}). \quad (23.37)$$

We write \mathbf{I}_D and \mathbf{I}_M for the $D \times D$ and $M \times M$ identity matrices respectively. The significance of this formula is that it turns a $D \times D$ determinant that costs $O(D^3)$ into an $M \times M$ determinant that costs $O(M^3)$. If M is smaller than D , this saves computation.

With some restrictions on the residual block $\mathbf{F} : \mathbb{R}^D \rightarrow \mathbb{R}^D$, we can apply this formula to compute the determinant of a residual connection efficiently. The trick is to create a bottleneck inside \mathbf{F} . We do that by defining $\mathbf{F} = \mathbf{F}_2 \circ \mathbf{F}_1$, where $\mathbf{F}_1 : \mathbb{R}^D \rightarrow \mathbb{R}^M$, $\mathbf{F}_2 : \mathbb{R}^M \rightarrow \mathbb{R}^D$ and $M \ll D$. The chain

². See https://en.wikipedia.org/wiki/Weinstein-Aronszajn_identity.

1 rule gives $\mathbf{J}(\mathbf{F}) = \mathbf{J}(\mathbf{F}_2)\mathbf{J}(\mathbf{F}_1)$, where $\mathbf{J}(\mathbf{F}_2)$ is $D \times M$ and $\mathbf{J}(\mathbf{F}_1)$ is $M \times D$. Now we can apply our
2 determinant formula as follows:

3

$$\det \mathbf{J}(\mathbf{f}) = \det(\mathbf{I}_D + \mathbf{J}(\mathbf{F})) = \det(\mathbf{I}_D + \mathbf{J}(\mathbf{F}_2)\mathbf{J}(\mathbf{F}_1)) = \det(\mathbf{I}_M + \mathbf{J}(\mathbf{F}_1)\mathbf{J}(\mathbf{F}_2)). \quad (23.38)$$

4 Since the final determinant costs $O(M^3)$, we can make the Jacobian determinant efficient by reducing
5 M , that is, by narrowing the bottleneck.

6 An example of the above is the **planar flow** of [RM15]. In this model, each residual block is an
7 MLP with one hidden layer and one hidden unit. That is,

8

$$\mathbf{f}(\mathbf{u}) = \mathbf{u} + \mathbf{v}\sigma(\mathbf{w}^\top \mathbf{u} + b), \quad (23.39)$$

9 where $\mathbf{v} \in \mathbb{R}^D$, $\mathbf{w} \in \mathbb{R}^D$ and $b \in \mathbb{R}$ are the parameters, and σ is the activation function. The residual
10 block is the composition of $\mathbf{F}_1(\mathbf{u}) = \mathbf{w}^\top \mathbf{u} + b$ and $\mathbf{F}_2(z) = \mathbf{v}\sigma(z)$, so $M = 1$. Their Jacobians
11 are $\mathbf{J}(\mathbf{F}_1)(\mathbf{u}) = \mathbf{w}^\top$ and $\mathbf{J}(\mathbf{F}_2)(z) = \mathbf{v}\sigma'(z)$. Substituting these in the formula for the Jacobian
12 determinant we obtain:

13

$$\det \mathbf{J}(\mathbf{f})(\mathbf{u}) = 1 + \mathbf{w}^\top \mathbf{v}\sigma'(\mathbf{w}^\top \mathbf{u} + b), \quad (23.40)$$

14 which can be computed efficiently in $O(D)$. Other examples include the **circular flow** of [RM15]
15 and the **Sylvester flow** of [Ber+18].

16 This technique gives an efficient way of computing determinants of residual connections with
17 bottlenecks, but in general there is no guarantee that such functions are invertible. This means that
18 invertibility must be satisfied on a case-by-case basis. For example, the planar flow is invertible when
19 σ is the hyperbolic tangent and $\mathbf{w}^\top \mathbf{v} > -1$, but otherwise it may not be.

20

21 23.2.6 Continuous-time flows

22 So far we have discussed flows that consist of a sequence of bijections $\mathbf{f}_1, \dots, \mathbf{f}_N$. Starting from some
23 input $\mathbf{x}_0 = \mathbf{u}$, this creates a sequence of outputs $\mathbf{x}_1, \dots, \mathbf{x}_N$ where $\mathbf{x}_n = \mathbf{f}_n(\mathbf{x}_{n-1})$. However, we can
24 also have flows where the input is transformed into the final output in a continuous way. That is,
25 we start from $\mathbf{x}_0 = \mathbf{x}(0)$, create a continuously-indexed sequence $\mathbf{x}(t)$ for $t \in [0, T]$ with some fixed
26 T , and take $\mathbf{x}(T)$ to be the final output. Thinking of t as analogous to time, we refer to these as
27 **continuous-time flows**.

28 The sequence $\mathbf{x}(t)$ is defined as the solution to a first-order ordinary differential equation (ODE)
29 of the form:

30

$$\frac{d\mathbf{x}}{dt}(t) = \mathbf{F}(\mathbf{x}(t), t). \quad (23.41)$$

31 The function $\mathbf{F} : \mathbb{R}^D \times [0, T] \rightarrow \mathbb{R}^D$ is a time-dependent vector field that parameterizes the ODE. If
32 we think of $\mathbf{x}(t)$ as the position of a particle in D dimensions, the vector $\mathbf{F}(\mathbf{x}(t), t)$ determines the
33 particle's velocity at time t .

34 The flow (for time T) is a function $\mathbf{f} : \mathbb{R}^D \rightarrow \mathbb{R}^D$ that takes in an input \mathbf{x}_0 , solves the ODE
35 with initial condition $\mathbf{x}(0) = \mathbf{x}_0$, and returns $\mathbf{x}(T)$. The function \mathbf{f} is a well-defined bijection if the
36 solution to the ODE exists for all $t \in [0, T]$ and is unique. These conditions are not generally satisfied
37 for arbitrary \mathbf{F} , but they are if $\mathbf{F}(\cdot, t)$ is Lipschitz continuous with a Lipschitz constant that does not
38

1 depend on t . That is, \mathbf{f} is a well-defined bijection if there exists a constant L such that for all $\mathbf{x}_1, \mathbf{x}_2$
2 and $t \in [0, T]$ we have:

4
$$\|\mathbf{F}(\mathbf{x}_1, t) - \mathbf{F}(\mathbf{x}_2, t)\| \leq L\|\mathbf{x}_1 - \mathbf{x}_2\|. \quad (23.42)$$

5

6 This result is a consequence of the **Picard-Lindelöf theorem** for ODEs.³ In practice, we can
7 parameterize \mathbf{F} using any choice of model, provided the Lipschitz condition is met.

8 Usually the ODE cannot be solved analytically, but we can solve it approximately by discretizing
9 it. A simple example is **Euler's method**, which corresponds to the following discretization for some
10 small step size $\epsilon > 0$:

11
$$\mathbf{x}(t + \epsilon) = \mathbf{x}(t) + \epsilon\mathbf{F}(\mathbf{x}(t), t). \quad (23.43)$$

12

13 This is equivalent to a residual connection with residual block $\epsilon\mathbf{F}(\cdot, t)$, so the ODE solver can be
14 thought of as a deep residual network with $O(T/\epsilon)$ layers. A smaller step size leads to a more
15 accurate solution, but also to more computation. There are several other solution methods varying
16 in accuracy and sophistication, such as those in the broader Runge-Kutta family, some of which use
17 adaptive step sizes.

18 The inverse of \mathbf{f} can be easily computed by solving the ODE in reverse. That is, to compute
19 $\mathbf{f}^{-1}(\mathbf{x}_T)$ we solve the ODE with initial condition $\mathbf{x}(T) = \mathbf{x}_T$, and return $\mathbf{x}(0)$. Unlike some other
20 flows (such as autoregressive flows) which are more expensive to compute in one direction than in
21 the other, continuous-time flows require the same amount of computation in either direction.

22 In general, there is no analytical expression for the Jacobian determinant of \mathbf{f} . However, we can
23 express it as the solution to a separate ODE, which we can then solve numerically. First, we define
24 $\mathbf{f}_t : \mathbb{R}^D \rightarrow \mathbb{R}^D$ to be the flow for time t , that is, the function that takes \mathbf{x}_0 , solves the ODE with
25 initial condition $\mathbf{x}(0) = \mathbf{x}_0$ and returns $\mathbf{x}(t)$. Clearly, \mathbf{f}_0 is the identity function and $\mathbf{f}_T = \mathbf{f}$. Let us
26 define $L(t) = \log |\det \mathbf{J}(\mathbf{f}_t)(\mathbf{x}_0)|$. Because \mathbf{f}_0 is the identity function, $L(0) = 0$, and because $\mathbf{f}_T = \mathbf{f}$,
27 $L(T)$ gives the Jacobian determinant of \mathbf{f} that we are interested in. It can be shown that L satisfies
28 the following ODE:

29
$$\frac{dL}{dt}(t) = \text{tr}[\mathbf{J}(\mathbf{F}(\cdot, t))(\mathbf{x}(t))]. \quad (23.44)$$

30

31 That is, the rate of change of L at time t is equal to the Jacobian trace of $\mathbf{F}(\cdot, t)$ evaluated at $\mathbf{x}(t)$. So
32 we can compute $L(T)$ by solving the above ODE with initial condition $L(0) = 0$. Moreover, we can
33 compute $\mathbf{x}(T)$ and $L(T)$ simultaneously, by combining their two ODEs into a single ODE operating
34 on the extended space (\mathbf{x}, L) .

35 An example of a continuous-time flow is the **neural ODE** model of [Che+18c], which uses a
36 neural network to parameterize \mathbf{F} . To avoid backpropagating gradients through the ODE solver,
37 which can be computationally demanding, they use the **adjoint sensitivity method** to express the
38 time evolution of the gradient with respect to $\mathbf{x}(t)$ as a separate ODE. Solving this ODE gives the
39 required gradients, and can be thought of as the continuous-time analog of backpropagation.

40 Another example is the **FFJORD** model of [Gra+19]. This is similar to the neural ODE model,
41 except that it uses the Hutchinson trace estimator to approximate the Jacobian trace of $\mathbf{F}(\cdot, t)$.
42 This usage of the Hutchinson trace estimator is analogous to that in contractive residual flows
43 (Section 23.2.5.1), and it speeds up computation in exchange for a stochastic (but unbiased) estimate.

44 See also Section 25.4.4, where we discuss continuous time diffusion models.

46 3. See https://en.wikipedia.org/wiki/Picard-Lindel%C3%B6f_theorem

1 **23.3 Applications**

3 In this section, we highlight some applications of flows for canonical probabilistic machine learning
4 tasks.
5

6

7 **23.3.1 Density estimation**

9 Flow models allow exact density computation and can be used to fit multi-modal densities to
10 observed data. (see Figure 23.3 for an example). An early example is Gaussianization [CG00] who
11 applied this idea to fit low-dimensional densities. Tabak and Vanden-Eijnden [TVE10] and Tabak
12 and Turner [TT13] introduced the modern idea of flows (including the term ‘normalizing flows’),
13 describing a flow as a composition of simpler maps. Deep density models [RA13] was one of the
14 first to use neural networks for flows to parameterize high-dimensional densities. There has been a
15 rich line of follow-up work including **NICE** [DKB15] and **Real NVP** [DSDB17]. (NVP stands for
16 “non-volume-preserving”, which refers to the fact that the Jacobian of the transform is not unity.)
17 Masked autoregressive flows (Section 23.2.4.2) further improved performance on unconditional and
18 conditional density estimation tasks.

19 Flows can be used for *hybrid models* which model the joint density of inputs and targets $p(\mathbf{x}, y)$, as
20 opposed to discriminative classification models which just model the conditional $p(y|\mathbf{x})$ and density
21 models which just model the marginal $p(\mathbf{x})$. Nalisnick et al. [Nal+19b] proposed a flow-based hybrid
22 model using invertible mappings for representation learning and showed that the joint density $p(\mathbf{x}, y)$
23 can be computed efficiently, which can be useful for downstream tasks such as anomaly detection,
24 semi-supervised learning and selective classification. Flow-based hybrid models are memory-efficient
25 since most of the parameters are in the invertible representation which are shared between the
26 discriminative and generative models; furthermore, the density $p(\mathbf{x}, y)$ can be computed in a single
27 forwards pass leading to computational savings. Residual flows [Che+19] use invertible residual
28 mappings [Beh+19] for hybrid modeling which further improves performance. Flows have also been
29 used to fit densities to embeddings [Zha+20b; CZG20] for anomaly detection tasks.

30

31 **23.3.2 Generative modeling**

32 Another task is generation, which involves generating novel samples from a fitted model $p^*(\mathbf{x})$.
33 Generation is a popular downstream task for normalizing flows, which have been applied for different
34 data modalities including images, video, audio, text, and structured objects such as graphs and point
35 clouds. Images are arguably the most popular modality for deep generative models: GLOW [KD18b]
36 was one of the first flow-based models to generate compelling high-dimensional images, and has been
37 extended to video to produce RGB frames [Kum+19b]; residual flows [Che+19] have also been shown
38 to produce sharp images.
39

40 Oord et al. [Oor+18] used flows for audio synthesis by distilling WaveNet into an IAF (Sec-
41 tion 23.2.4.3), which enables faster sampling than WaveNet. Other flow models for audio include
42 WaveFLOW [PVC19] and FlowWaveNet [Kim+19], which directly speed up WaveNet using coupling
43 layers.

44 Flows have been also used for text. Tran et al. [Tra+19] define a discrete flow over a vocabulary
45 for language-modeling tasks. Another popular approach is to define a latent variable model with
46 discrete observation space but a continuous latent space. For example, Ziegler and Rush [ZR19a] use
47

1 normalizing flows in latent space for language modeling.
2
3

4 **23.3.3 Inference**

5 Normalizing flows have been used for probabilistic inference. Rezende and Mohamed [RM15]
6 popularized normalizing flows in machine learning, and showed how they can be used for modeling
7 variational posterior distributions in latent variable models. Various extensions such as Householder
8 flows [TW16], inverse autoregressive flows [Kin+16], multiplicative normalizing flows [LW17], and
9 Sylvester flows [Ber+18] have been proposed for modeling the variational posterior for latent variable
10 models, as well as posteriors for Bayesian neural networks.
11

12 Flows have been used as complex proposal distributions for importance sampling; examples include
13 neural importance sampling [Mül+19b] and Boltzmann generators [Noé+19]. Hoffman et al. [Hof+19]
14 used flows to improve the performance of Hamiltonian Monte Carlo (Section 12.5) by defining bijective
15 transformations to transform random variables to simpler distributions and performing HMC in that
16 space instead.
17

Finally, flows can be used in the context of simulation-based inference, where the likelihood function
of the parameters is not available, but simulating data from the model is possible. The main idea is
to train a flow on data simulated from the model in order to approximate the posterior distribution
or the likelihood function. The flow model can also be used to guide simulations in order to make
inference more efficient [PSM19; GNM19]. This approach has been used for inference of simulation
models in cosmology [Als+19] and computational neuroscience [Gon+20].
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

24 Energy-based models

This chapter is co-authored with Yang Song and Durk Kingma.

24.1 Introduction

We have now seen several ways of defining deep generative models, including VAEs (Chapter 21), autoregressive models (Chapter 22), and normalizing flows (Chapter 23). All of the above models can be formulated in terms of directed graphical models (Chapter 4), where we generate the data one step at a time, using locally normalized distributions. In some cases, it is easier to specify a distribution in terms of a set of constraints that valid samples must satisfy, rather than a generative process. This can be done using an undirected graphical model (Chapter 4).

Energy-based models or **EBM** can be written as a Gibbs distribution as follows:

$$p_{\theta}(\mathbf{x}) = \frac{\exp(-\mathcal{E}_{\theta}(\mathbf{x}))}{Z_{\theta}} \quad (24.1)$$

where $\mathcal{E}_{\theta}(\mathbf{x}) \geq 0$ is known as the **energy function** with parameters θ , and Z_{θ} is the **partition function**:

$$Z_{\theta} = \int \exp(-\mathcal{E}_{\theta}(\mathbf{x})) \, d\mathbf{x} \quad (24.2)$$

This is constant wrt \mathbf{x} but is a function of θ . Since EBMs do not usually make any Markov assumptions (unlike graphical models), evaluating this integral is usually intractable. Consequently we usually need to use approximate methods, such as annealed importance sampling, discussed in Section 11.5.4.1.

The advantage of an EBM over other generative models is that the energy function can be any kind of function that returns a non-negative scalar; it does not need to integrate to 1. This allows one to use a variety of neural network architectures for defining the energy. As such, EBMs have found wide applications in many fields of machine learning, including image generation [Ngi+11; Xie+16; DM19b], discriminative learning [Gra+20b], natural processing [Mik+13; Den+20], density estimation [Wen+19a; Son+19], and reinforcement learning [Haa+17; Haa+18a], to list a few. (More examples can be found at <https://github.com/yataobian/awesome-ebm>.)

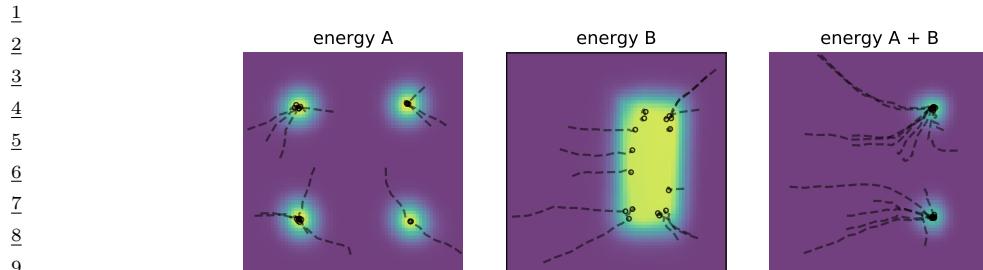


Figure 24.1: Combining two energy functions in 2d by summation, which is equivalent to multiplying the corresponding probability densities. We also illustrate some sampled trajectories towards high probability (low energy) regions. From Figure 14 of [DM19a]. Used with kind permission of Yilun Du.

24.1.1 Example: products of experts (PoE)

As an example of why energy based models are useful, suppose we want to create a generative model of proteins that are thermally stable at room temperature, and which bind to the COVID-19 spike receptor. Suppose $p_1(\mathbf{x})$ can generate stable proteins and $p_2(\mathbf{x})$ can generate proteins that bind. (For example, both of these models could be autoregressive sequence models, trained on different datasets.) We can view each of these models as “experts” about a particular aspect of the data. On their own, they are not an adequate model of the data that we have (or want to have), but we can then combine them, to represent the **conjunction of features**, by computing a **product of experts (PoE)** [Hin02]:

$$p_{12}(\mathbf{x}) = \frac{1}{Z_{12}} p_1(\mathbf{x}) p_2(\mathbf{x}) \quad (24.3)$$

This will assign high probability to proteins that are stable and which bind, and low probability to all others. By contrast, a **mixture of experts** would either generate from p_1 or from p_2 , but would not combine features from both.

If the experts are represented as energy based models (EBM), then the PoE model is also an EBM, with an energy given by

$$\mathcal{E}_{12}(\mathbf{x}) = \mathcal{E}_1(\mathbf{x}) + \mathcal{E}_2(\mathbf{x}) \quad (24.4)$$

Intuitively, we can think of each component of energy as a “soft constraint” on the data. This idea is illustrated in Figure 24.1.

24.1.2 Computational difficulties

Although the flexibility of EBMs can provide significant modeling advantages, computation of the likelihood and drawing samples from the model are generally intractable. In this chapter, we will discuss a variety of approximate methods to solve these problems.

1

24.2 Maximum likelihood training

2

3 The de facto standard for learning probabilistic models from iid data is maximum likelihood estimation
4 (MLE). Let $p_{\theta}(\mathbf{x})$ be a probabilistic model parameterized by θ , and $p_{\mathcal{D}}(\mathbf{x})$ be the underlying data
5 distribution of a dataset. We can fit $p_{\theta}(\mathbf{x})$ to $p_{\mathcal{D}}(\mathbf{x})$ by maximizing the expected log-likelihood
6 function over the data distribution, defined by
7

8

$$\ell(\theta) = \mathbb{E}_{\mathbf{x} \sim p_{\mathcal{D}}(\mathbf{x})} [\log p_{\theta}(\mathbf{x})] \quad (24.5)$$

9

10 as a function of θ . Here the expectation can be easily estimated with samples from the dataset.
11 Maximizing likelihood is equivalent to minimizing the KL divergence between $p_{\mathcal{D}}(\mathbf{x})$ and $p_{\theta}(\mathbf{x})$,
12 because
13

14

$$\ell(\theta) = -D_{\text{KL}}(p_{\mathcal{D}}(\mathbf{x}) \parallel p_{\theta}(\mathbf{x})) + \text{const} \quad (24.6)$$

15

16 where the constant is equal to $\mathbb{E}_{\mathbf{x} \sim p_{\mathcal{D}}(\mathbf{x})} [\log p_{\mathcal{D}}(\mathbf{x})]$ which does not depend on θ .
17

18 We cannot usually compute the likelihood of an EBM because the normalizing constant Z_{θ}
19 is often intractable. Nevertheless, we can still estimate the gradient of the log-likelihood with
20 MCMC approaches, allowing for likelihood maximization with stochastic gradient ascent [You99]. In
21 particular, the gradient of the log-probability of an EBM decomposes as a sum of two terms:
22

23

$$\nabla_{\theta} \log p_{\theta}(\mathbf{x}) = -\nabla_{\theta} \mathcal{E}_{\theta}(\mathbf{x}) - \nabla_{\theta} \log Z_{\theta}. \quad (24.7)$$

24

25 The first gradient term, $-\nabla_{\theta} \mathcal{E}_{\theta}(\mathbf{x})$, is straightforward to evaluate with automatic differentiation. The
26 challenge is in approximating the second gradient term, $\nabla_{\theta} \log Z_{\theta}$, which is intractable to compute
27 exactly. This gradient term can be rewritten as the following expectation:
28

29

$$\nabla_{\theta} \log Z_{\theta} = \nabla_{\theta} \log \int \exp(-\mathcal{E}_{\theta}(\mathbf{x})) d\mathbf{x} \quad (24.8)$$

30

31

$$\stackrel{(i)}{=} \left(\int \exp(-\mathcal{E}_{\theta}(\mathbf{x})) d\mathbf{x} \right)^{-1} \nabla_{\theta} \int \exp(-\mathcal{E}_{\theta}(\mathbf{x})) d\mathbf{x} \quad (24.9)$$

32

33

$$= \left(\int \exp(-\mathcal{E}_{\theta}(\mathbf{x})) d\mathbf{x} \right)^{-1} \int \nabla_{\theta} \exp(-\mathcal{E}_{\theta}(\mathbf{x})) d\mathbf{x} \quad (24.10)$$

34

35

$$\stackrel{(ii)}{=} \left(\int \exp(-\mathcal{E}_{\theta}(\mathbf{x})) d\mathbf{x} \right)^{-1} \int \exp(-\mathcal{E}_{\theta}(\mathbf{x})) (-\nabla_{\theta} \mathcal{E}_{\theta}(\mathbf{x})) d\mathbf{x} \quad (24.11)$$

36

37

$$= \int \left(\int \exp(-\mathcal{E}_{\theta}(\mathbf{x})) d\mathbf{x} \right)^{-1} \exp(-\mathcal{E}_{\theta}(\mathbf{x})) (-\nabla_{\theta} \mathcal{E}_{\theta}(\mathbf{x})) d\mathbf{x} \quad (24.12)$$

38

39

$$\stackrel{(iii)}{=} \int \frac{1}{Z_{\theta}} \exp(-\mathcal{E}_{\theta}(\mathbf{x})) (-\nabla_{\theta} \mathcal{E}_{\theta}(\mathbf{x})) d\mathbf{x} \quad (24.13)$$

40

41

$$\stackrel{(iv)}{=} \int p_{\theta}(\mathbf{x}) (-\nabla_{\theta} \mathcal{E}_{\theta}(\mathbf{x})) d\mathbf{x} \quad (24.14)$$

42

43

$$= \mathbb{E}_{\mathbf{x} \sim p_{\theta}(\mathbf{x})} [-\nabla_{\theta} \mathcal{E}_{\theta}(\mathbf{x})], \quad (24.15)$$

44

where steps (i) and (ii) are due to the chain rule of gradients, and (iii) and (iv) are from definitions in Equations (24.1) and (24.2). Thus, we can obtain an unbiased Monte Carlo estimate of the log-likelihood gradient by using

$$\nabla_{\theta} \log Z_{\theta} \simeq -\frac{1}{S} \sum_{s=1}^S \nabla_{\theta} \mathcal{E}_{\theta}(\tilde{\mathbf{x}}_s), \quad (24.16)$$

where $\tilde{\mathbf{x}}_s \sim p_{\theta}(\mathbf{x})$, i.e., a random sample from the distribution over \mathbf{x} given by the EBM. Therefore, as long as we can draw random samples from the model, we have access to an unbiased Monte Carlo estimate of the log-likelihood gradient, allowing us to optimize the parameters with stochastic gradient ascent.

Much of the literature has focused on methods for efficient MCMC sampling from EBMs. We discuss some of these methods below.

24.2.1 Gradient-based MCMC methods

Some efficient MCMC methods, such as **Langevin MCMC** (Section 12.5.6) or Hamiltonian Monte Carlo (Section 12.5), make use of the fact that the gradient of the log-probability wrt \mathbf{x} (known as the **score function**) is equal to the (negative) gradient of the energy, and is therefore easy to calculate:

$$\nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}) = -\nabla_{\mathbf{x}} \mathcal{E}_{\theta}(\mathbf{x}) - \underbrace{\nabla_{\mathbf{x}} \log Z_{\theta}}_{=0} = -\nabla_{\mathbf{x}} \mathcal{E}_{\theta}(\mathbf{x}). \quad (24.17)$$

For example, when using Langevin MCMC to sample from $p_{\theta}(\mathbf{x})$, we first draw an initial sample \mathbf{x}^0 from a simple prior distribution, and then simulate an overdamped Langevin diffusion process for K steps with step size $\epsilon > 0$:

$$\mathbf{x}^{k+1} \leftarrow \mathbf{x}^k + \frac{\epsilon^2}{2} \underbrace{\nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}^k)}_{=-\nabla_{\mathbf{x}} \mathcal{E}_{\theta}(\mathbf{x})} + \epsilon \mathbf{z}^k, \quad k = 0, 1, \dots, K-1. \quad (24.18)$$

where $\mathbf{z}^k \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ is a Gaussian noise term. We show an example of this process in Figure 25.5d.

When $\epsilon \rightarrow 0$ and $K \rightarrow \infty$, \mathbf{x}^K is guaranteed to distribute as $p_{\theta}(\mathbf{x})$ under some regularity conditions. In practice we have to use a small finite ϵ , but the discretization error is typically negligible, or can be corrected with a Metropolis-Hastings step (Section 12.2), leading to the Metropolis-adjusted Langevin algorithm (Section 12.5.6).

24.2.2 Contrastive divergence

Running MCMC till convergence to obtain a sample $\mathbf{x} \sim p_{\theta}(\mathbf{x})$ can be computationally expensive. Therefore we typically need approximations to make MCMC-based learning of EBMs practical. One popular method for doing so is **contrastive divergence** (CD) [Hin02]. In CD, one initializes the MCMC chain from the datapoint \mathbf{x} , and proceeds to perform MCMC for a fixed number of steps. One can show that T steps of CD minimizes the following objective:

$$\text{CD}_T = D_{\text{KL}}(p_0 \parallel p_{\infty}) - D_{\text{KL}}(p_T \parallel p_{\infty}) \quad (24.19)$$

where p_T is the distribution over \mathbf{x} after T MCMC updates, and p_0 is the data distribution. Typically we can get good results with a small value of T , sometimes just $T = 1$. We give the details below.

24.2.2.1 Fitting RBMs with CD

CD was initially developed to fit a special kind of latent variable EBM known as a restricted Boltzmann machine (Section 4.3.3.2). This model was specifically designed to support fast block Gibbs sampling, which is required by CD (and can also be exploited by standard MCMC-based learning methods [AHS85].)

For simplicity, we will assume the hidden and visible nodes are binary, and we use 1-step contrastive divergence. As discussed in Supplementary Section 4.3.1, the binary RBM has the following energy function:

$$\mathcal{E}(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) = \sum_{d=1}^D \sum_{k=1}^K x_d z_k W_{dk} + \sum_{d=1}^D x_d b_d + \sum_{k=1}^K z_k c_k \quad (24.20)$$

(Henceforth we will drop the unary (bias) terms, which can be emulated by clamping $z_k = 1$ or $x_d = 1$.) This is a loglinear model where we have one binary feature per edge. Thus from Equation (4.135) the gradient of the log-likelihood is given by the clamped expectations minus the unclamped expectations:

$$\frac{\partial \ell}{\partial w_{dk}} = \frac{1}{N} \sum_{n=1}^N \mathbb{E}[x_d z_k | \mathbf{x}_n, \boldsymbol{\theta}] - \mathbb{E}[x_d z_k | \boldsymbol{\theta}] \quad (24.21)$$

We can rewrite the above gradient in matrix-vector form as follows:

$$\nabla_{\mathbf{w}} \ell = \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})} [\mathbf{x}\mathbf{z}^T] - \mathbb{E}_{p(\mathbf{z}, \mathbf{x}|\boldsymbol{\theta})} [\mathbf{x}\mathbf{z}^T] \quad (24.22)$$

(We can derive a similar expression for the gradient of the bias terms by setting $x_d = 1$ or $z_k = 1$.)

The first term in the expression for the gradient in Equation (24.21), when \mathbf{x} is fixed to a data case, is sometimes called the **clamped phase**, and the second term, when \mathbf{x} is free, is sometimes called the **unclamped phase**. When the model expectations match the empirical expectations, the two terms cancel out, the gradient becomes zero and learning stops.

We can also make a connection to the principle of **Hebbian learning** in neuroscience. In particular, Hebb's rule says that the strength of connection between two neurons that are simultaneously active should be increased. (This theory is often summarized as "Cells that fire together wire together".¹) The first term in Equation (24.21) is therefore considered a Hebbian term, and the second term an anti-Hebbian term, due to the sign change.

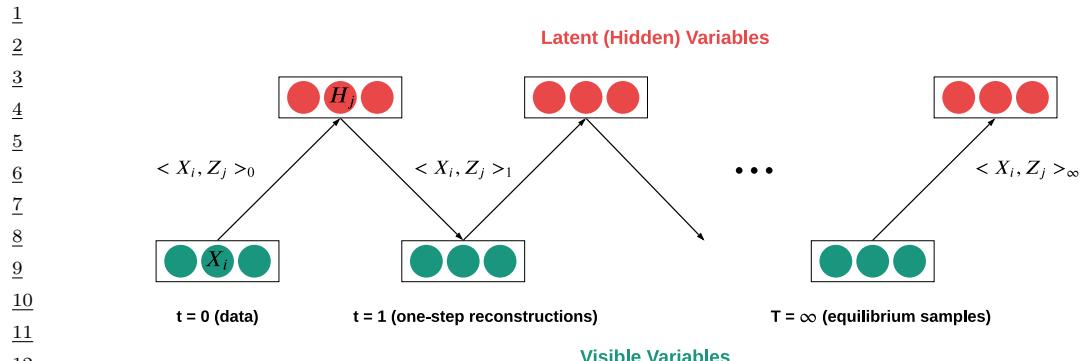
We can leverage the Markov structure of the bipartite graph to approximate the expectations as follows:

$$\mathbf{z}_n \sim p(\mathbf{z}|\mathbf{x}_n, \boldsymbol{\theta}) \quad (24.23)$$

$$\mathbf{x}'_n \sim p(\mathbf{x}|\mathbf{z}_n, \boldsymbol{\theta}) \quad (24.24)$$

$$\mathbf{z}'_n \sim p(\mathbf{z}|\mathbf{x}'_n, \boldsymbol{\theta}) \quad (24.25)$$

¹. See https://en.wikipedia.org/wiki/Hebbian_theory.



14 Figure 24.2: Illustration of contrastive divergence sampling for an RBM. The visible nodes are initialized
15 at an example drawn from the dataset. Then we sample a hidden vector, then another visible vector, etc.
16 Eventually (at “infinity”) we will be producing samples from the joint distribution $p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta})$.

¹⁹ We can think of $\hat{\mathbf{x}}'_n$ as the model's best attempt at reconstructing \mathbf{x}_n after being encoded and then
²⁰ decoded by the model. Such samples are sometimes called **fantasy data**. See Figure 24.2 for an
²¹ illustration. Given these samples, we then make the approximation

$$\mathbb{E}_{n(\cdot), \theta} [\mathbf{x} z^\top] \approx \mathbf{x}_n(z'_r)^\top \quad (24.26)$$

In practice, it is common to use $\mathbb{E}[\mathbf{z}|\mathbf{x}'_n]$ instead of a sampled value \mathbf{z}'_n in the above expression, since this reduces the variance. However, it is not valid to use $\mathbb{E}[\mathbf{z}|\mathbf{x}_n]$ instead of sampling $\mathbf{z}_n \sim p(\mathbf{z}|\mathbf{x}_n)$ in Equation (24.23), because then each hidden unit would be able to pass more than 1 bit of information, so it would not act as much of a bottleneck.

²⁸ The whole procedure is summarized in Algorithm 24.1. For more details, see [Hin10; Swe+10].

Algorithm 24.1: CD-1 training for an RBM with binary hidden and visible units

```

32 1 Initialize weights  $\mathbf{W} \in \mathbb{R}^{D \times K}$  randomly
33 2 for  $t = 1, 2, \dots$  do
34 3   for each minibatch of size  $B$  do
35 4     Set minibatch gradient to zero,  $\mathbf{g} := \mathbf{0}$ 
36 5     for each case  $\mathbf{x}_n$  in the minibatch do
37 6       Compute  $\boldsymbol{\mu}_n = \mathbb{E}[\mathbf{z}|\mathbf{x}_n, \mathbf{W}]$ 
38 7       Sample  $\mathbf{z}_n \sim p(\mathbf{z}|\mathbf{x}_n, \mathbf{W})$ 
39 8       Sample  $\mathbf{x}'_n \sim p(\mathbf{x}|\mathbf{z}_n, \mathbf{W})$ 
40 9       Compute  $\boldsymbol{\mu}'_n = \mathbb{E}[\mathbf{z}|\mathbf{x}'_n, \mathbf{W}]$ 
41 10      Compute gradient  $\nabla_{\mathbf{W}} = (\mathbf{x}_n)(\boldsymbol{\mu}_n)^T - (\mathbf{x}'_n)(\boldsymbol{\mu}'_n)^T$ 
42 11      Accumulate  $\mathbf{g} := \mathbf{g} + \nabla_{\mathbf{W}}$ 
43 12      Update parameters  $\mathbf{W} := \mathbf{W} + \eta_t \frac{1}{B} \mathbf{g}$ 

```

24.2.2.2 Persistent CD

One variant of CD that sometimes performs better is **persistent contrastive divergence** (PCD) [Tie08; TH09; You99]. In this approach, a single MCMC chain with a persistent state is employed to sample from the EBM. In PCD, we do not restart the MCMC chain when training on a new datapoint; rather, we carry over the state of the previous MCMC chain and use it to initialize a new MCMC chain for the next training step. See Line 12 for some pseudocode. Hence there are two dynamical processes running at different time scales: the states \mathbf{x} change quickly, and the parameters θ change slowly.

Algorithm 24.2: Persistent MCMC-SGD for fitting an EBM

```

1 Initialize parameters  $\theta$  randomly
2 Initialize chains  $\tilde{\mathbf{x}}_{1:S}$  randomly
3 Initialize learning rate  $\eta$ 
4 for  $t = 1, 2, \dots$  do
5   for  $\mathbf{x}_b$  in minibatch of size  $B$  do
6      $\mathbf{g}_b = \nabla_{\theta} \mathcal{E}_{\theta}(\mathbf{x}_b)$ 
7   for sample  $s = 1 : S$  do
8     Sample  $\tilde{\mathbf{x}}_s \sim \text{MCMC}(\text{target} = p(\cdot | \theta), \text{init} = \tilde{\mathbf{x}}_s, \text{nsteps} = N)$ 
9      $\tilde{\mathbf{g}}_s = \nabla_{\theta} \mathcal{E}_{\theta}(\tilde{\mathbf{x}}_s)$ 
10     $\mathbf{g}_t = -(\frac{1}{B} \sum_{b=1}^B \mathbf{g}_b) - (\frac{1}{S} \sum_{s=1}^S \tilde{\mathbf{g}}_s)$ 
11     $\theta := \theta + \eta \mathbf{g}_t$ 
12    Decrease step size  $\eta$ 

```

A theoretical justification for this was given in [You89], who showed that we can start the MCMC chain at its previous value, and just take a few steps, because $p(\mathbf{x} | \theta_t)$ is likely to be close to $p(\mathbf{x} | \theta_{t-1})$, since we only changed the parameters by a small amount in the intervening SGD step.

24.2.2.3 Other methods

PCD can be further improved by keeping multiple historical states of the MCMC chain in a replay buffer and initialize new MCMC chains by randomly sampling from it [DM19b]. Other variants of CD include mean field CD [WH02], and multi-grid CD [Gao+18].

EBMs trained with CD may not capture the data distribution faithfully, since truncated MCMC can lead to biased gradient updates that hurt the learning dynamics [SMB10; FI10; Nij+19]. There are several methods that focus on removing this bias for improved MCMC training. For example, one line of work proposes unbiased estimators of the gradient through coupled MCMC [JOA17; QZW19]; and Du et al. [Du+20] propose to reduce the bias by differentiating through the MCMC sampling algorithm and estimating an entropy correction term.

24.3 Score matching (SM)

If two continuously differentiable real-valued functions $f(\mathbf{x})$ and $g(\mathbf{x})$ have equal first derivatives everywhere, then $f(\mathbf{x}) \equiv g(\mathbf{x}) + \text{constant}$. When $f(\mathbf{x})$ and $g(\mathbf{x})$ are log probability density functions (pdf's) with equal first derivatives, the normalization requirement (Equation (24.1)) implies that $\int \exp(f(\mathbf{x}))d\mathbf{x} = \int \exp(g(\mathbf{x}))d\mathbf{x} = 1$, and therefore $f(\mathbf{x}) \equiv g(\mathbf{x})$. As a result, one can learn an EBM by (approximately) matching the first derivatives of its log-pdf to the first derivatives of the log pdf of the data distribution. If they match, then the EBM captures the data distribution exactly. The first-order gradient function of a log pdf wrt its input, $\nabla_{\mathbf{x}} \log p_{\boldsymbol{\theta}}(\mathbf{x})$, is called the (Stein) **score** function. (This is distinct from the Fisher score, $\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{x})$.) For training EBMs, it is useful to transform the equivalence of distributions to the equivalence of scores, because the score of an EBM can be easily obtained as follows:

$$\mathbf{s}_{\boldsymbol{\theta}}(\mathbf{x}) \triangleq \nabla_{\mathbf{x}} \log p_{\boldsymbol{\theta}}(\mathbf{x}) = -\nabla_{\mathbf{x}} \mathcal{E}_{\boldsymbol{\theta}}(\mathbf{x}) \quad (24.27)$$

We see that this does not involve the typically intractable normalizing constant $Z_{\boldsymbol{\theta}}$.

Let $p_{\mathcal{D}}(\mathbf{x})$ be the underlying data distribution, from which we have a finite number of iid samples but do not know its pdf. The **score matching** objective [Hyr05] minimizes a discrepancy between two distributions called the **Fisher divergence**:

$$D_F(p_{\mathcal{D}}(\mathbf{x}) \parallel p_{\boldsymbol{\theta}}(\mathbf{x})) = \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} \left[\frac{1}{2} \|\nabla_{\mathbf{x}} \log p_{\mathcal{D}}(\mathbf{x}) - \nabla_{\mathbf{x}} \log p_{\boldsymbol{\theta}}(\mathbf{x})\|^2 \right]. \quad (24.28)$$

The expectation wrt $p_{\mathcal{D}}(\mathbf{x})$, in this objective and its variants below, admits a trivial unbiased Monte Carlo estimator using the empirical mean of samples $\mathbf{x} \sim p_{\mathcal{D}}(\mathbf{x})$. However, the second term of Equation (24.28), $\nabla_{\mathbf{x}} \log p_{\mathcal{D}}(\mathbf{x})$, is generally impractical to calculate since it requires knowing the pdf of $p_{\mathcal{D}}(\mathbf{x})$. We discuss a solution to this below.

24.3.1 Basic score matching

Hyrén [Hyr05] shows that, under certain regularity conditions, the Fisher divergence can be rewritten using integration by parts, with second derivatives of $\mathcal{E}_{\boldsymbol{\theta}}(\mathbf{x})$ replacing the unknown first derivatives of $p_{\mathcal{D}}(\mathbf{x})$:

$$D_F(p_{\mathcal{D}}(\mathbf{x}) \parallel p_{\boldsymbol{\theta}}(\mathbf{x})) = \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} \left[\frac{1}{2} \sum_{i=1}^d \left(\frac{\partial \mathcal{E}_{\boldsymbol{\theta}}(\mathbf{x})}{\partial x_i} \right)^2 - \frac{\partial^2 \mathcal{E}_{\boldsymbol{\theta}}(\mathbf{x})}{\partial x_i^2} \right] + \text{constant} \quad (24.29)$$

$$= \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} \left[\frac{1}{2} \|\mathbf{s}_{\boldsymbol{\theta}}(\mathbf{x})\|^2 + \text{tr}(\mathbf{J}_{\mathbf{x}} \mathbf{s}_{\boldsymbol{\theta}}(\mathbf{x})) \right] + \text{constant} \quad (24.30)$$

where d is the dimensionality of \mathbf{x} , and $\mathbf{J}_{\mathbf{x}} \mathbf{s}_{\boldsymbol{\theta}}(\mathbf{x})$ is the Jacobian of the score function. The constant does not affect optimization and thus can be dropped for training. It is shown by [Hyr05] that estimators based on score matching are consistent under some regularity conditions, meaning that the parameter estimator obtained by minimizing Equation (24.28) converges to the true parameters in the limit of infinite data. See Figure 25.5 for an example.

An important downside of the objective Equation (24.30) is that it takes $O(d^2)$ time to compute the trace of the Jacobian. For this reason, the implicit SM formulation of Equation (24.30) has only

been applied to relatively simple energy functions where computation of the second derivatives is tractable.

Score Matching assumes a continuous data distribution with positive density over the space, but it can be generalized to discrete or bounded data distributions [Hyv07b; Lyu12]. It is also possible to consider higher-order gradients of log pdf's beyond first derivatives [PDL+12].

24.3.2 Denoising score matching (DSM)

The Score Matching objective in Equation (24.30) requires several regularity conditions for $\log p_{\mathcal{D}}(\mathbf{x})$, e.g., it should be continuously differentiable and finite everywhere. However, these conditions may not always hold in practice. For example, a distribution of digital images is typically discrete and bounded, because the values of pixels are restricted to the range $\{0, 1, \dots, 255\}$. Therefore, $\log p_{\mathcal{D}}(\mathbf{x})$ in this case is discontinuous and is negative infinity outside the range, and thus SM is not directly applicable.

To alleviate this, one can add a bit of noise to each datapoint: $\tilde{\mathbf{x}} = \mathbf{x} + \boldsymbol{\epsilon}$. As long as the noise distribution $p(\boldsymbol{\epsilon})$ is smooth, the resulting noisy data distribution $q(\tilde{\mathbf{x}}) = \int q(\tilde{\mathbf{x}} | \mathbf{x}) p_{\mathcal{D}}(\mathbf{x}) d\mathbf{x}$ is also smooth, and thus the Fisher divergence $D_F(q(\tilde{\mathbf{x}}) \| p_{\theta}(\tilde{\mathbf{x}}))$ is a proper objective. [KL10] showed that the objective with noisy data can be approximated by the noiseless Score Matching objective of Equation (24.30) plus a regularization term; this regularization makes Score Matching applicable to a wider range of data distributions, but still requires expensive second-order derivatives.

[Vin11] proposed an elegant and scalable solution to the above difficulty, by showing that:

$$D_F(q(\tilde{\mathbf{x}}) \| p_{\theta}(\tilde{\mathbf{x}})) = \mathbb{E}_{q(\tilde{\mathbf{x}})} \left[\frac{1}{2} \|\nabla_{\mathbf{x}} \log p_{\theta}(\tilde{\mathbf{x}}) - \nabla_{\mathbf{x}} \log q(\tilde{\mathbf{x}})\|_2^2 \right] \quad (24.31)$$

$$= \mathbb{E}_{q(\mathbf{x}, \tilde{\mathbf{x}})} \left[\frac{1}{2} \|\nabla_{\mathbf{x}} \log p_{\theta}(\tilde{\mathbf{x}}) - \nabla_{\mathbf{x}} \log q(\tilde{\mathbf{x}}|\mathbf{x})\|_2^2 \right] + \text{constant} \quad (24.32)$$

$$= \frac{1}{2} \mathbb{E}_{q(\mathbf{x}, \tilde{\mathbf{x}})} \left[\left\| s_{\theta}(\tilde{\mathbf{x}}) + \frac{(\tilde{\mathbf{x}} - \mathbf{x})}{\sigma^2} \right\|_2^2 \right] \quad (24.33)$$

where $s_{\theta}(\tilde{\mathbf{x}}) = \nabla_{\mathbf{x}} \log p_{\theta}(\tilde{\mathbf{x}})$ is the estimated score function, and

$$\nabla_{\mathbf{x}} \log q(\tilde{\mathbf{x}}|\mathbf{x}) = \nabla_{\mathbf{x}} \log \mathcal{N}(\tilde{\mathbf{x}}|\mathbf{x}, \sigma^2 \mathbf{I}) = \frac{(\tilde{\mathbf{x}} - \mathbf{x})}{\sigma^2} \quad (24.34)$$

The expectation can be approximated by sampling from $p_{\mathcal{D}}(\mathbf{x})$ and then sampling the noise term $\tilde{\mathbf{x}}$. (The constant term does not affect optimization and can be ignored without changing the optimal solution.)

This estimation method is called **denoising score matching** (DSM) by [Vin11]. Similar formulations were also explored by Raphan and Simoncelli [RS07; RS11] and can be traced back to Tweedie's formula (Supplementary Section 3.3) and Stein's unbiased risk estimation [Ste81].

24.3.2.1 Difficulties

The major drawback of adding noise to data arises when $p_{\mathcal{D}}(\mathbf{x})$ is already a well-behaved distribution that satisfies the regularity conditions required by score matching. In this case, $D_F(q(\tilde{\mathbf{x}}) \| p_{\theta}(\tilde{\mathbf{x}})) \neq$

1 $D_F(p_{\mathcal{D}}(\mathbf{x}) \parallel p_{\boldsymbol{\theta}}(\mathbf{x}))$, and DSM is not a consistent objective because the optimal EBM matches the
2 noisy distribution $q(\tilde{\mathbf{x}})$, not $p_{\mathcal{D}}(\mathbf{x})$. This inconsistency becomes non-negligible when $q(\tilde{\mathbf{x}})$ significantly
3 differs from $p_{\mathcal{D}}(\mathbf{x})$.
4

5 One way to attenuate the inconsistency of DSM is to choose $q \approx p_{\mathcal{D}}$, i.e., use a small noise
6 perturbation. However, this often significantly increases the variance of objective values and hinders
7 optimization. As an example, suppose $q(\tilde{\mathbf{x}} \mid \mathbf{x}) = \mathcal{N}(\tilde{\mathbf{x}} \mid \mathbf{x}, \sigma^2 I)$ and $\sigma \approx 0$. The corresponding DSM
8 objective is

$$\begin{aligned} \frac{9}{10} D_F(q(\tilde{\mathbf{x}}) \parallel p_{\boldsymbol{\theta}}(\tilde{\mathbf{x}})) &= \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(0, I)} \left[\frac{1}{2} \left\| \frac{\mathbf{z}}{\sigma} + \nabla_{\mathbf{x}} \log p_{\boldsymbol{\theta}}(\mathbf{x} + \sigma \mathbf{z}) \right\|_2^2 \right] \\ \frac{11}{12} &\simeq \frac{1}{2N} \sum_{i=1}^N \left\| \frac{\mathbf{z}^{(i)}}{\sigma} + \nabla_{\mathbf{x}} \log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)} + \sigma \mathbf{z}^{(i)}) \right\|_2^2, \end{aligned} \quad (24.35)$$

14 where $\{\mathbf{x}^{(i)}\}_{i=1}^N \stackrel{\text{i.i.d.}}{\sim} p_{\mathcal{D}}(\mathbf{x})$, and $\{\mathbf{z}^{(i)}\}_{i=1}^N \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, I)$. When $\sigma \rightarrow 0$, we can leverage Taylor series
15 expansion to rewrite the Monte Carlo estimator in Equation (24.35) to
16

$$\begin{aligned} \frac{17}{18} \frac{1}{2N} \sum_{i=1}^N \left[\frac{2}{\sigma} (\mathbf{z}^{(i)})^\top \nabla_{\mathbf{x}} \log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) + \frac{\|\mathbf{z}^{(i)}\|_2^2}{\sigma^2} \right] + \text{constant}. \end{aligned} \quad (24.36)$$

20 When estimating the above expectation with samples, the variances of $(\mathbf{z}^{(i)})^\top \nabla_{\mathbf{x}} \log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})/\sigma$ and
21 $\|\mathbf{z}^{(i)}\|_2^2/\sigma^2$ will both grow unbounded as $\sigma \rightarrow 0$ due to division by σ and σ^2 . This enlarges the
22 variance of DSM and makes optimization challenging. Various methods have been proposed to reduce
23 this variance (see e.g., [Wan+20d]).
24

25 24.3.3 Sliced score matching (SSM)

27 By adding noise to data, DSM avoids the expensive computation of second-order derivatives. However,
28 as mentioned before, the optimal EBM that minimizes the DSM objective corresponds to the
29 distribution of noise-perturbed data $q(\tilde{\mathbf{x}})$, not the original noise-free data distribution $p_{\mathcal{D}}(\mathbf{x})$. In
30 other words, DSM does not give a consistent estimator of the data distribution, i.e., one cannot
31 directly obtain an EBM that exactly matches the data distribution even with unlimited data.
32

33 **Sliced score matching (SSM)** [Son+19] is one alternative to Denoising Score Matching that is
34 both consistent and computationally efficient. Instead of minimizing the Fisher divergence between
35 two vector-valued scores, SSM randomly samples a projection vector \mathbf{v} , takes the inner product
36 between \mathbf{v} and the two scores, and then compares the resulting two scalars. More specifically, sliced
37 score matching minimizes the following divergence called the **sliced Fisher divergence**:

$$\begin{aligned} \frac{38}{39} D_{SF}(p_{\mathcal{D}}(\mathbf{x}) \parallel p_{\boldsymbol{\theta}}(\mathbf{x})) &= \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} \mathbb{E}_{p(\mathbf{v})} \left[\frac{1}{2} (\mathbf{v}^\top \nabla_{\mathbf{x}} \log p_{\mathcal{D}}(\mathbf{x}) - \mathbf{v}^\top \nabla_{\mathbf{x}} \log p_{\boldsymbol{\theta}}(\mathbf{x}))^2 \right], \end{aligned} \quad (24.37)$$

40 where $p(\mathbf{v})$ denotes a projection distribution such that $\mathbb{E}_{p(\mathbf{v})}[\mathbf{v}\mathbf{v}^\top]$ is positive definite. Similar to
41 Fisher divergence, sliced Fisher divergence has an implicit form that does not involve the unknown
42 $\nabla_{\mathbf{x}} \log p_{\mathcal{D}}(\mathbf{x})$, which is given by
43

$$\begin{aligned} \frac{44}{45} D_{SF}(p_{\mathcal{D}}(\mathbf{x}) \parallel p_{\boldsymbol{\theta}}(\mathbf{x})) &= \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} \mathbb{E}_{p(\mathbf{v})} \left[\frac{1}{2} \sum_{i=1}^d \left(\frac{\partial \mathcal{E}_{\boldsymbol{\theta}}(\mathbf{x})}{\partial x_i} v_i \right)^2 + \sum_{i=1}^d \sum_{j=1}^d \frac{\partial^2 \mathcal{E}_{\boldsymbol{\theta}}(\mathbf{x})}{\partial x_i \partial x_j} v_i v_j \right] + C. \end{aligned} \quad (24.38)$$

All expectations in the above objective can be estimated with empirical means, and again the constant term C can be removed without affecting training. The second term involves second-order derivatives of $\mathcal{E}_\theta(\mathbf{x})$, but contrary to SM, it can be computed efficiently with a cost linear in the dimensionality d . This is because

$$\sum_{i=1}^d \sum_{j=1}^d \frac{\partial^2 \mathcal{E}_\theta(\mathbf{x})}{\partial x_i \partial x_j} v_i v_j = \sum_{i=1}^d \frac{\partial}{\partial x_i} \left(\underbrace{\sum_{j=1}^d \frac{\partial \mathcal{E}_\theta(\mathbf{x})}{\partial x_j} v_j}_{:= f(\mathbf{x})} \right) v_i, \quad (24.39)$$

where $f(\mathbf{x})$ is the same for different values of i . Therefore, we only need to compute it once with $O(d)$ computation, *plus* another $O(d)$ computation for the outer sum to evaluate Equation (24.39), whereas the original SM objective requires $O(d^2)$ computation.

For many choices of $p(\mathbf{v})$, part of the SSM objective (Equation (24.38)) can be evaluated in closed form, potentially leading to lower variance. For example, when $p(\mathbf{v}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$, we have

$$\mathbb{E}_{p_D(\mathbf{x})} \mathbb{E}_{p(\mathbf{v})} \left[\frac{1}{2} \sum_{i=1}^d \left(\frac{\partial \mathcal{E}_\theta(\mathbf{x})}{\partial x_i} v_i \right)^2 \right] = \mathbb{E}_{p_D(\mathbf{x})} \left[\frac{1}{2} \sum_{i=1}^d \left(\frac{\partial \mathcal{E}_\theta(\mathbf{x})}{\partial x_i} \right)^2 \right] \quad (24.40)$$

and as a result,

$$D_{SF}(p_D(\mathbf{x}) \| p_\theta(\mathbf{x})) = \mathbb{E}_{p_D(\mathbf{x})} \mathbb{E}_{\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[\frac{1}{2} \sum_{i=1}^d \left(\frac{\partial \mathcal{E}_\theta(\mathbf{x})}{\partial x_i} \right)^2 + \sum_{i=1}^d \sum_{j=1}^d \frac{\partial^2 \mathcal{E}_\theta(\mathbf{x})}{\partial x_i \partial x_j} v_i v_j \right] + C \quad (24.41)$$

$$= \mathbb{E}_{p_D(\mathbf{x})} \mathbb{E}_{\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[\frac{1}{2} (\mathbf{v}^\top \mathbf{s}_\theta(\mathbf{x}))^2 + \mathbf{v}^\top [\mathbf{J}\mathbf{v}] \right] \quad (24.42)$$

where $\mathbf{J} = \mathbf{J}_x \mathbf{s}_\theta(\mathbf{x})$. (Note that $\mathbf{J}\mathbf{v}$ can be computed using a Jacobian vector product operation.)

The above objective Equation (24.41) can also be obtained by approximating the sum of second-order gradients in the standard SM objective (Equation (24.30)) with the Hutchinson trace estimator [Ski89; Hut89; Mey+21]. It often (but not always) has lower variance than Equation (24.38), and can perform better in some applications [Son+19].

24.3.4 Connection to contrastive divergence

Though score matching and contrastive divergence (Section 24.2.2) are seemingly very different approaches, they are closely connected to each other. In fact, score matching can be viewed as a special instance of contrastive divergence in the limit of a particular MCMC sampler [Hyv07a]. Moreover, the Fisher divergence optimized by Score Matching is related to the derivative of KL divergence [Cov99], which is the underlying objective of Contrastive Divergence.

Contrastive divergence requires sampling from the EBM $\mathcal{E}_\theta(\mathbf{x})$, and one popular method for doing so is Langevin MCMC. Recall from Section 24.2.1 that given any initial datapoint \mathbf{x}^0 , the Langevin MCMC method executes the following

$$\mathbf{x}^{k+1} \leftarrow \mathbf{x}^k - \frac{\epsilon}{2} \nabla_{\mathbf{x}} \mathcal{E}_\theta(\mathbf{x}^k) + \sqrt{\epsilon} \mathbf{z}^k, \quad (24.43)$$

1 iteratively for $k = 0, 1, \dots, K - 1$, where $\mathbf{z}^k \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $\epsilon > 0$ is the step size.

2 Suppose we only run one-step Langevin MCMC for contrastive divergence. In this case, the
3 gradient of the log-likelihood is given by
4

$$\begin{aligned} \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})}[\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{x})] &= -\mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})}[\nabla_{\boldsymbol{\theta}} \mathcal{E}_{\boldsymbol{\theta}}(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim p_{\boldsymbol{\theta}}(\mathbf{x})}[\nabla_{\boldsymbol{\theta}} \mathcal{E}_{\boldsymbol{\theta}}(\mathbf{x})] \\ &\simeq -\mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})}[\nabla_{\boldsymbol{\theta}} \mathcal{E}_{\boldsymbol{\theta}}(\mathbf{x})] + \mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[\nabla_{\boldsymbol{\theta}} \mathcal{E}_{\boldsymbol{\theta}} \left(\mathbf{x} - \frac{\epsilon^2}{2} \nabla_{\mathbf{x}} E_{\boldsymbol{\theta}'}(\mathbf{x}) + \epsilon \mathbf{z} \right) \middle|_{\boldsymbol{\theta}'=\boldsymbol{\theta}} \right]. \end{aligned} \quad (24.44)$$

10 After Taylor series expansion with respect to ϵ followed by some algebraic manipulations, the above
11 equation can be transformed to the following [Hyv07a]:
12

$$\frac{\epsilon^2}{2} \nabla_{\boldsymbol{\theta}} D_F(p_{\mathcal{D}}(\mathbf{x}) \| p_{\boldsymbol{\theta}}(\mathbf{x})) + o(\epsilon^2). \quad (24.45)$$

16 When ϵ is sufficiently small, it corresponds to the re-scaled gradient of the score matching objective.

17 In general, score matching minimizes the Fisher divergence $D_F(p_{\mathcal{D}}(\mathbf{x}) \| p_{\boldsymbol{\theta}}(\mathbf{x}))$, whereas Contrastive
18 Divergence minimizes an objective related to the KL divergence $D_{KL}(p_{\mathcal{D}}(\mathbf{x}) \| p_{\boldsymbol{\theta}}(\mathbf{x}))$, as shown in
19 Equation (24.19). The above connection of score matching and Contrastive Divergence is a natural
20 consequence of the connection between those two statistical divergences, as characterized by *de*
21 *de Bruijin's identity* [Cov99; Lyu12]:

$$\frac{d}{dt} D_{KL}(q_t(\tilde{\mathbf{x}}) \| p_{\boldsymbol{\theta},t}(\tilde{\mathbf{x}})) = -\frac{1}{2} D_F(q_t(\tilde{\mathbf{x}}) \| p_{\boldsymbol{\theta},t}(\tilde{\mathbf{x}})).$$

25 Here $q_t(\tilde{\mathbf{x}})$ and $p_{\boldsymbol{\theta},t}(\tilde{\mathbf{x}})$ denote smoothed versions of $p_{\mathcal{D}}(\mathbf{x})$ and $p_{\boldsymbol{\theta}}(\mathbf{x})$, resulting from adding Gaussian
26 noise to \mathbf{x} with variance t ; i.e., $\tilde{\mathbf{x}} \sim \mathcal{N}(\mathbf{x}, t\mathbf{I})$.
27

29 24.3.5 Score-based generative models

30 We have seen how to use score matching to fit EBMs by learning the scalar energy function $\mathcal{E}_{\boldsymbol{\theta}}(\mathbf{x})$.
31 We can alternatively directly learn the score function, $s_{\boldsymbol{\theta}}(\mathbf{x}) = \nabla_{\mathbf{x}} \log p_{\boldsymbol{\theta}}(\mathbf{x})$; this is called a score-
32 based generative model, and is discussed in Section 25.3. Such unconstrained score models are not
33 guaranteed to output a conservative vector field, meaning they do not correspond to the gradient of
34 any function. However, both methods seem to give comparable results [SH21].
35

36

37 24.4 Noise contrastive estimation

38 Another principle for learning the parameters of EBMs is **Noise contrastive estimation** (NCE),
39 introduced by [GH10]. It is based on the idea that we can learn an EBM by contrasting it with
40 another distribution with known density.
41

42 Let $p_{\mathcal{D}}(\mathbf{x})$ be our data distribution, and let $p_n(\mathbf{x})$ be a chosen distribution with known density,
43 called a noise distribution. This noise distribution is usually simple and has a tractable pdf, like
44 $\mathcal{N}(\mathbf{0}, \mathbf{I})$, such that we can compute the pdf and generate samples from it efficiently. Strategies exist
45 to learn the noise distribution, as referenced below. Furthermore, let y be a binary variable with
46 Bernoulli distribution, which we use to define a mixture distribution of noise and data: $p_{n,\text{data}}(\mathbf{x}) =$
47

$p(y = 0)p_n(\mathbf{x}) + p(y = 1)p_{\mathcal{D}}(\mathbf{x})$. According to Bayes' rule, given a sample \mathbf{x} from this mixture, the posterior probability of $y = 0$ is

$$p_{n,\text{data}}(y = 0 \mid \mathbf{x}) = \frac{p_{n,\text{data}}(\mathbf{x} \mid y = 0)p(y = 0)}{p_{n,\text{data}}(\mathbf{x})} = \frac{p_n(\mathbf{x})}{p_n(\mathbf{x}) + \nu p_{\mathcal{D}}(\mathbf{x})} \quad (24.46)$$

where $\nu = p(y = 1)/p(y = 0)$.

Let our EBM $p_{\theta}(\mathbf{x})$ be defined as:

$$p_{\theta}(\mathbf{x}) = \exp(-\mathcal{E}_{\theta}(\mathbf{x}))/Z_{\theta} \quad (24.47)$$

Contrary to most other EBMs, Z_{θ} is treated as a learnable (scalar) parameter in NCE. Given this model, similar to the mixture of noise and data above, we can define a mixture of noise and the model distribution: $p_{n,\theta}(\mathbf{x}) = p(y = 0)p_n(\mathbf{x}) + p(y = 1)p_{\theta}(\mathbf{x})$. The posterior probability of $y = 0$ given this noise/model mixture is:

$$p_{n,\theta}(y = 0 \mid \mathbf{x}) = \frac{p_n(\mathbf{x})}{p_n(\mathbf{x}) + \nu p_{\theta}(\mathbf{x})} \quad (24.48)$$

In NCE, we indirectly fit $p_{\theta}(\mathbf{x})$ to $p_{\mathcal{D}}(\mathbf{x})$ by fitting $p_{n,\theta}(y \mid \mathbf{x})$ to $p_{n,\text{data}}(y \mid \mathbf{x})$ through a standard conditional maximum likelihood objective:

$$\theta^* = \operatorname{argmin}_{\theta} \mathbb{E}_{p_{n,\text{data}}(\mathbf{x})} [D_{KL}(p_{n,\text{data}}(y \mid \mathbf{x}) \parallel p_{n,\theta}(y \mid \mathbf{x}))] \quad (24.49)$$

$$= \operatorname{argmax}_{\theta} \mathbb{E}_{p_{n,\text{data}}(\mathbf{x}, y)} [\log p_{n,\theta}(y \mid \mathbf{x})], \quad (24.50)$$

which can be solved using stochastic gradient ascent. Just like any other deep classifier, when the model is sufficiently powerful, $p_{n,\theta^*}(y \mid \mathbf{x})$ will match $p_{n,\text{data}}(y \mid \mathbf{x})$ at the optimum. In that case:

$$p_{n,\theta^*}(y = 0 \mid \mathbf{x}) \equiv p_{n,\text{data}}(y = 0 \mid \mathbf{x}) \quad (24.51)$$

$$\iff \frac{p_n(\mathbf{x})}{p_n(\mathbf{x}) + \nu p_{\theta^*}(\mathbf{x})} \equiv \frac{p_n(\mathbf{x})}{p_n(\mathbf{x}) + \nu p_{\mathcal{D}}(\mathbf{x})} \quad (24.52)$$

$$\iff p_{\theta^*}(\mathbf{x}) \equiv p_{\mathcal{D}}(\mathbf{x}) \quad (24.53)$$

Consequently, $E_{\theta^*}(\mathbf{x})$ is an unnormalized energy function that matches the data distribution $p_{\mathcal{D}}(\mathbf{x})$, and Z_{θ^*} is the corresponding normalizing constant.

As one unique feature that contrastive divergence and score matching do not have, NCE provides the normalizing constant of an Energy-Based Model as a by-product of its training procedure. When the EBM is very expressive, e.g., a deep neural network with many parameters, we can assume it is able to approximate a normalized probability density and absorb Z_{θ} into the parameters of $\mathcal{E}_{\theta}(\mathbf{x})$ [MT12], or equivalently, fixing $Z_{\theta} = 1$. The resulting EBM trained with NCE will be self-normalized, i.e., having a normalizing constant close to 1.

In practice, choosing the right noise distribution $p_n(\mathbf{x})$ is critical to the success of NCE, especially for structured and high-dimensional data. As argued in Gutmann and Hirayama [GH12], NCE works the best when the noise distribution is close to the data distribution (but not exactly the same). Many methods have been proposed to automatically tune the noise distribution, such as Adversarial Contrastive Estimation [BLC18], Conditional NCE [CG18] and Flow Contrastive Estimation [Gao+20]. NCE can be further generalized using Bregman divergences (Section 5.1.10), where the formulation introduced here reduces to a special case.

1 **24.4.1 Connection to score matching**

3 Noise contrastive estimation provides a family of objectives that vary for different $p_n(\mathbf{x})$ and ν . This
4 flexibility may allow adaptation to special properties of a task with hand-tuned $p_n(\mathbf{x})$ and ν , and
5 may also give a unified perspective for different approaches. In particular, when using an appropriate
6 $p_n(\mathbf{x})$ and a slightly different parameterization of $p_{n,\theta}(y | \mathbf{x})$, we can recover score matching from
7 NCE [GH12].

8 Specifically, we choose the noise distribution $p_n(\mathbf{x})$ to be a perturbed data distribution: given a
9 small (deterministic) vector \mathbf{v} , let $p_n(\mathbf{x}) = p_{\mathcal{D}}(\mathbf{x} - \mathbf{v})$. It is efficient to sample from this $p_n(\mathbf{x})$, since
10 we can first draw any datapoint $\mathbf{x}' \sim p_{\mathcal{D}}(\mathbf{x}')$ and then compute $\mathbf{x} = \mathbf{x}' + \mathbf{v}$. It is, however, difficult
11 to evaluate the density of $p_n(\mathbf{x})$ because $p_{\mathcal{D}}(\mathbf{x})$ is unknown. Since the original parameterization of
12 $p_{n,\theta}(y | \mathbf{x})$ in NCE (Equation (24.48)) depends on the pdf of $p_n(\mathbf{x})$, we cannot directly apply the
13 standard NCE objective. Instead, we replace $p_n(\mathbf{x})$ with $p_{\theta}(\mathbf{x} - \mathbf{v})$ and parameterize $p_{n,\theta}(y = 0 | \mathbf{x})$
14 with the following form

15

$$\frac{16}{17} \quad p_{n,\theta}(y = 0 | \mathbf{x}) := \frac{p_{\theta}(\mathbf{x} - \mathbf{v})}{p_{\theta}(\mathbf{x}) + p_{\theta}(\mathbf{x} - \mathbf{v})} \quad (24.54)$$

18 In this case, the NCE objective (Equation (24.50)) reduces to:

20

$$\frac{21}{22} \quad \theta^* = \operatorname{argmin}_{\theta} \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [\log(1 + \exp(\mathcal{E}_{\theta}(\mathbf{x}) - \mathcal{E}_{\theta}(\mathbf{x} - \mathbf{v})) + \log(1 + \exp(\mathcal{E}_{\theta}(\mathbf{x}) - \mathcal{E}_{\theta}(\mathbf{x} + \mathbf{v})))] \quad (24.55)$$

23

24 At θ^* , we have a solution where:

25

$$\frac{26}{27} \quad p_{n,\theta^*}(y = 0 | \mathbf{x}) \equiv p_{n,\text{data}}(y = 0 | \mathbf{x}) \quad (24.56)$$

27

$$\frac{28}{29} \quad \Rightarrow \frac{p_{\theta^*}(\mathbf{x} - \mathbf{v})}{p_{\theta^*}(\mathbf{x}) + p_{\theta^*}(\mathbf{x} - \mathbf{v})} \equiv \frac{p_{\mathcal{D}}(\mathbf{x} - \mathbf{v})}{p_{\mathcal{D}}(\mathbf{x}) + p_{\mathcal{D}}(\mathbf{x} - \mathbf{v})} \quad (24.57)$$

30 which implies that $p_{\theta^*}(\mathbf{x}) \equiv p_{\mathcal{D}}(\mathbf{x})$, i.e., our model matches the data distribution.

31 As noted in Gutmann and Hirayama [GH12] and Song et al. [Son+19], when $\|\mathbf{v}\|_2 \approx 0$, the NCE
32 objective Equation (24.50) has the following equivalent form by Taylor expansion

33

$$\frac{34}{35} \quad \operatorname{argmin}_{\theta} \frac{1}{4} \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} \left[\frac{1}{2} \sum_{i=1}^d \left(\frac{\partial \mathcal{E}_{\theta}(\mathbf{x})}{\partial x_i} v_i \right)^2 + \sum_{i=1}^d \sum_{j=1}^d \frac{\partial^2 \mathcal{E}_{\theta}(\mathbf{x})}{\partial x_i \partial x_j} v_i v_j \right] + 2 \log 2 + o(\|\mathbf{v}\|_2^2). \quad (24.58)$$

36

37 Comparing against Equation (24.38), we immediately see that the above objective equals that of
38 SSM, if we ignore small additional terms hidden in $o(\|\mathbf{v}\|_2^2)$ and take the expectation with respect to
39 \mathbf{v} over a user-specified distribution $p(\mathbf{v})$.

41

42 **24.5 Other methods**

43 Aside from MCMC-based training, score matching and noise contrastive estimation, there are also
44 other methods for learning EBMs. Below we briefly survey some examples of them. Interested readers
45 can learn more details from references therein.

46

24.5.1 Minimizing Differences/Derivatives of KL Divergences

The overarching strategy for learning probabilistic models from data is to minimize the KL divergence between data and model distributions. However, because the normalizing constants of EBMs are typically intractable, it is hard to directly evaluate the KL divergence when the model is an EBM (see the discussion in Section 24.2.1). One generic idea that has frequently circumvented this difficulty is to consider differences/derivatives of KL divergences. It turns out that the unknown partition functions of EBMs are often cancelled out after taking the difference of two closely related KL divergences, or computing the derivatives.

Typical examples of this strategy include minimum velocity learning [Mov08; Wan+20d], minimum probability flow [SDBD11], and minimum KL contraction [Lyu11], to name a few. In minimum velocity learning and minimum probability flow, a Markov chain is designed such that it starts from the data distribution $p_{\mathcal{D}}(\mathbf{x})$ and converges to the EBM distribution $p_{\boldsymbol{\theta}}(\mathbf{x}) = e^{-\mathcal{E}_{\boldsymbol{\theta}}(\mathbf{x})}/Z_{\boldsymbol{\theta}}$. Specifically, the Markov chain satisfies $p_0(\mathbf{x}) \equiv p_{\mathcal{D}}(\mathbf{x})$ and $p_{\infty}(\mathbf{x}) \equiv p_{\boldsymbol{\theta}}(\mathbf{x})$, where we denote by $p_t(\mathbf{x})$ the state distribution at time $t \geq 0$.

This Markov chain will evolve towards $p_{\boldsymbol{\theta}}(\mathbf{x})$ unless $p_{\mathcal{D}}(\mathbf{x}) \equiv p_{\boldsymbol{\theta}}(\mathbf{x})$. Therefore, we can fit the EBM distribution $p_{\boldsymbol{\theta}}(\mathbf{x})$ to $p_{\mathcal{D}}(\mathbf{x})$ by minimizing the modulus of the “velocity” of this evolution, defined by

$$\frac{d}{dt} D_{\text{KL}}(p_t(\mathbf{x}) \parallel p_{\boldsymbol{\theta}}(\mathbf{x})) \Big|_{t=0} \quad \text{or} \quad \frac{d}{dt} D_{\text{KL}}(p_{\mathcal{D}}(\mathbf{x}) \parallel p_t(\mathbf{x})) \Big|_{t=0} \quad (24.59)$$

in minimum velocity learning and minimum probability flow respectively. These objectives typically do not require computing the normalizing constant $Z_{\boldsymbol{\theta}}$.

In minimum KL contraction [Lyu11], a distribution transformation Φ is chosen such that

$$D_{\text{KL}}(p(\mathbf{x}) \parallel q(\mathbf{x})) \geq D_{\text{KL}}(\Phi\{p(\mathbf{x})\} \parallel \Phi\{q(\mathbf{x})\}) \quad (24.60)$$

with equality if and only if $p(\mathbf{x}) = q(\mathbf{x})$. We can leverage this Φ to train an EBM, by minimizing

$$D_{\text{KL}}(p_{\mathcal{D}}(\mathbf{x}) \parallel p_{\boldsymbol{\theta}}(\mathbf{x})) - D_{\text{KL}}(\Phi\{p_{\mathcal{D}}(\mathbf{x})\} \parallel \Phi\{p_{\boldsymbol{\theta}}(\mathbf{x})\}). \quad (24.61)$$

This objective does not require computing the partition function $Z_{\boldsymbol{\theta}}$ whenever Φ is linear.

Minimum velocity learning, minimum probability flow, and minimum KL contraction can all be viewed as generalizations to score matching and noise contrastive estimation [Mov08; SDBD11; Lyu11].

24.5.2 Minimizing the Stein discrepancy

We can train EBMs by minimizing the Stein discrepancy, defined by

$$D_{\text{Stein}}(p_{\mathcal{D}}(\mathbf{x}) \parallel p_{\boldsymbol{\theta}}(\mathbf{x})) := \sup_{\mathbf{f} \in \mathcal{F}} \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [\nabla_{\mathbf{x}} \log p_{\boldsymbol{\theta}}(\mathbf{x})^T \mathbf{f}(\mathbf{x}) + \text{trace}(\nabla_{\mathbf{x}} \mathbf{f}(\mathbf{x}))], \quad (24.62)$$

where \mathcal{F} is a family of vector-valued functions, and $\nabla_{\mathbf{x}} \mathbf{f}(\mathbf{x})$ denotes the Jacobian of $\mathbf{f}(\mathbf{x})$. With some regularity conditions [GM15; LLJ16], we have $D_S(p_{\mathcal{D}}(\mathbf{x}) \parallel p_{\boldsymbol{\theta}}(\mathbf{x})) \geq 0$, where the equality holds if and only if $p_{\mathcal{D}}(\mathbf{x}) \equiv p_{\boldsymbol{\theta}}(\mathbf{x})$. Similar to score matching (Equation (24.30)), the objective Equation (24.62) only involves the score function of $p_{\boldsymbol{\theta}}(\mathbf{x})$, and does not require computing the EBM’s

1 partition function. Still, the trace term in Equation (24.62) may demand expensive computation,
2 and does not scale well to high dimensional data.

3 There are two common methods that sidestep this difficulty. Gorham and Mackey [GM15] and
4 Liu, Lee, and Jordan [LLJ16] discovered that when \mathcal{F} is a unit ball in a reproducing kernel Hilbert
5 space (RKHS) with a fixed kernel, the Stein discrepancy becomes kernelized Stein discrepancy, where
6 the trace term is a constant and does not affect optimization. Otherwise, $\text{trace}(\nabla_{\mathbf{x}} \mathbf{f}(\mathbf{x}))$ can be
7 approximated with the Skilling-Hutchinson trace estimator [Ski89; Hut89; Gra+20c].
8

9 24.5.3 Adversarial training

10 Recall from Section 24.2.1 that when training EBMs with maximum likelihood estimation (MLE),
11 we need to sample from the EBM per training iteration. However, sampling using multiple MCMC
12 steps is expensive and requires careful tuning of the Markov chain. One way to avoid this difficulty
13 is to use non-MLE methods that do not need sampling, such as score matching and noise contrastive
14 estimation. Here we introduce another family of methods that sidestep costly MCMC sampling by
15 learning an auxiliary model through adversarial training, which allows fast sampling.
16

17 Using the definition of EBMs, we can rewrite the maximum likelihood objective by introducing a
18 variational distribution $q_{\phi}(\mathbf{x})$ parameterized by ϕ :

$$\begin{aligned} \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})}[\log p_{\theta}(\mathbf{x})] &= \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})}[-\mathcal{E}_{\theta}(\mathbf{x})] - \log Z_{\theta} \\ &= \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})}[-\mathcal{E}_{\theta}(\mathbf{x})] - \log \int e^{-\mathcal{E}_{\theta}(\mathbf{x})} d\mathbf{x} \\ &= \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})}[-\mathcal{E}_{\theta}(\mathbf{x})] - \log \int q_{\phi}(\mathbf{x}) \frac{e^{-\mathcal{E}_{\theta}(\mathbf{x})}}{q_{\phi}(\mathbf{x})} d\mathbf{x} \\ &\stackrel{(i)}{\leq} \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})}[-\mathcal{E}_{\theta}(\mathbf{x})] - \int q_{\phi}(\mathbf{x}) \log \frac{e^{-\mathcal{E}_{\theta}(\mathbf{x})}}{q_{\phi}(\mathbf{x})} d\mathbf{x} \\ &= \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})}[-\mathcal{E}_{\theta}(\mathbf{x})] - \mathbb{E}_{q_{\phi}(\mathbf{x})}[-\mathcal{E}_{\theta}(\mathbf{x})] - H(q_{\phi}(\mathbf{x})), \end{aligned} \quad (24.63)$$

30 where $H(q_{\phi}(\mathbf{x}))$ denotes the entropy of $q_{\phi}(\mathbf{x})$. Step (i) is due to Jensen's inequality. Equation (24.63)
31 provides an upper bound to the expected log-likelihood. For EBM training, we can first minimize the
32 upper bound Equation (24.63) with respect to $q_{\phi}(\mathbf{x})$ so that it is closer to the likelihood objective,
33 and then maximize Equation (24.63) with respect to $\mathcal{E}_{\theta}(\mathbf{x})$ as a surrogate for maximizing likelihood.
34 This amounts to using the following maximin objective

$$\max_{\theta} \min_{\phi} \mathbb{E}_{q_{\phi}(\mathbf{x})}[\mathcal{E}_{\theta}(\mathbf{x})] - \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})}[\mathcal{E}_{\theta}(\mathbf{x})] - H(q_{\phi}(\mathbf{x})). \quad (24.64)$$

35 Optimizing the above objective is similar to training GANs (Chapter 26), and can be achieved by
36 adversarial training. The variational distribution $q_{\phi}(\mathbf{x})$ should allow both fast sampling and efficient
37 entropy evaluation to make Equation (24.64) tractable. This limits the model family of $q_{\phi}(\mathbf{x})$, and
38 usually restricts our choice to invertible probabilistic models, such as inverse autoregressive flow
39 (Section 23.2.4.3). See Dai et al. [Dai+19b] for an example on designing $q_{\phi}(\mathbf{x})$ and training EBMs
40 with Equation (24.64).

41 Kim and Bengio [KB16] and Zhai et al. [Zha+16] propose to represent $q_{\phi}(\mathbf{x})$ with neural samplers,
42 like the generator of GANs. A neural sampler is a deterministic mapping g_{ϕ} that maps a random
43 variable \mathbf{z} to a sample \mathbf{x} from $q_{\phi}(\mathbf{x})$.
44

Gaussian noise $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ directly to a sample $\mathbf{x} = g_\phi(\mathbf{z})$. When using a neural sampler as $q_\phi(\mathbf{x})$, it is efficient to draw samples through the deterministic mapping, but $H(q_\phi(\mathbf{x}))$ is intractable since the density of $q_\phi(\mathbf{x})$ is unknown. Kim and Bengio [KB16] and Zhai et al. [Zha+16] propose several heuristics to approximate this entropy function. Kumar et al. [Kum+19c] propose to estimate the entropy through its connection to mutual information: $H(q_\phi(\mathbf{z})) = I(g_\phi(\mathbf{z}), \mathbf{z})$, which can be estimated from samples with variational lower bounds [NWJ10b; NCT16b]. Dai et al. [Dai+19a] noticed that when defining $p_\theta(\mathbf{x}) = p_0(\mathbf{x})e^{-\mathcal{E}_\theta(\mathbf{x})}/Z_\theta$, with $p_0(\mathbf{x})$ being a fixed base distribution, the entropy term $-H(q_\phi(\mathbf{x}))$ in Equation (24.64) can be replaced by $D_{\text{KL}}(q_\phi(\mathbf{x}) \parallel p_0(\mathbf{x}))$, which can also be approximated with variational lower bounds using samples from $q_\phi(\mathbf{x})$ and $p_0(\mathbf{x})$, without requiring the density of $q_\phi(\mathbf{x})$.

Grathwohl et al. [Gra+20a] represent $q_\phi(\mathbf{x})$ as a noisy neural sampler, where samples are obtained via $g_\phi(\mathbf{z}) + \sigma\epsilon$, assuming $\mathbf{z}, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. With a noisy neural sampler, $\nabla_\phi H(q_\phi(\mathbf{x}))$ becomes particularly easy to estimate, which allows gradient-based optimization for the minimax objective in Equation (24.63). A related approach is proposed in Xie et al. [Xie+18], where authors train a noisy neural sampler with samples obtained from MCMC, and initialize new MCMC chains with samples generated from the neural sampler. This cooperative sampling scheme improves the convergence of MCMC, but may still require multiple MCMC steps for sample generation. It does not optimize the objective in Equation (24.63).

When using both adversarial training and MCMC sampling, Yu et al. [Yu+20] noticed that EBMs can be trained with an arbitrary f -divergence, including KL, reverse KL, total variation, Hellinger, etc. The method proposed by Yu et al. [Yu+20] allows us to explore the trade-offs and inductive bias of different statistical divergences for more flexible EBM training.

24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

25 Diffusion models

25.1 Introduction

In this chapter, we consider a class of models called **diffusion models**. This class of models has recently generated a lot of interest, due to its ability to generate diverse, high quality, samples, and the relative simplicity of the training scheme, which allows very large models to be trained at scale. Diffusion models are closely related to VAEs (Chapter 21), normalizing flows (Chapter 23), and EBMs (Chapter 24), as we will see.

The basic idea behind these models is based on the observation that it is hard to convert noise into structured data, but it is easy to convert structured data into noise. In particular, we can use a **forwards process** or **diffusion process** to gradually convert the observed data \mathbf{x}_0 into a noisy version \mathbf{x}_T by passing the data through T steps of a stochastic encoder $q(\mathbf{x}_t|\mathbf{x}_{t-1})$. After enough steps, we have $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, or some other convenient reference distribution. We then learn a **reverse process** to undo this, by passing the noise through T steps of a decoder $p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$ until we generate \mathbf{x}_0 . See Figure 25.1 for an overall sketch of the approach. In the following sections, we discuss this class of models in more detail. Our presentation is based in part on the excellent tutorial [KGV22]. More details can be found in the recent review papers [Yan+22; Cao+22], as well as specialized papers, such as [Kar+22]. There are also many other excellent resources online, such as <https://github.com/heejkoo/Awesome-Diffusion-Models> and <https://scorebasedgenerativemodeling.github.io/>. For a detailed tutorial on the underlying math, see [McA23].

25.2 Denoising diffusion probabilistic models (DDPMs)

In this section, we discuss **denoising diffusion probabilistic models** or **DDPMs**, introduced in [SD+15b], and then extended in [HJA20; Kin+21] and many other works. We can think of the DDPM as similar to a hierarchical variational autoencoder (Section 21.5), except that all the latent states (denoted \mathbf{x}_t for $t = 1 : T$) have the same dimensionality as the input \mathbf{x}_0 . (In this respect, a DDPM is similar to a normalizing flow (Chapter 23); however, in a diffusion model, the hidden layers are stochastic, and do not need to use invertible transformations.) In addition, the encoder network q is a simple linear Gaussian model, rather than being learned¹, and the decoder network p is shared across all time steps. These restrictions result in a very simple training objective, which

1. Later we will discuss some extensions in which the noise level of the encoder can also be learned. Nevertheless, the encoder remains simple, by design.

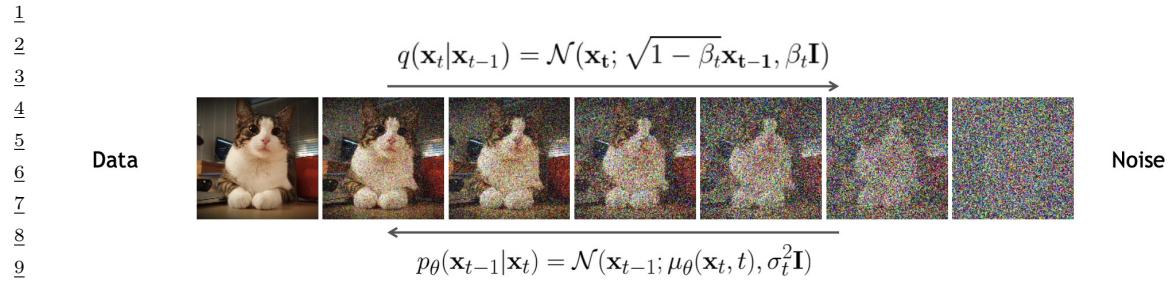


Figure 25.1: Denoising diffusion probabilistic model. The forwards diffusion process, $q(\mathbf{x}_t | \mathbf{x}_{t-1})$, implements the (non-learned) inference network; this just adds Gaussian noise at each step. The reverse diffusion process, $p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)$, implements the decoder; this is a learned Gaussian model. From Slide 16 of [KGV22]. Used with kind permission of Arash Vahdat.

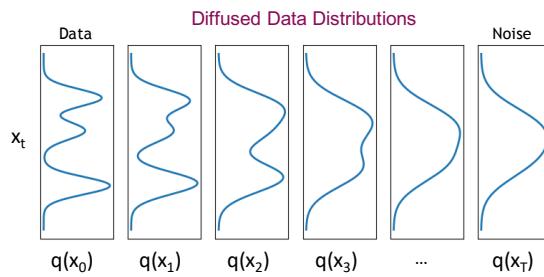


Figure 25.2: Illustration of a diffusion model on 1d data. The forwards diffusion process gradually transforms the empirical data distribution $q(\mathbf{x}_0)$ into a simple target distribution, here $q(\mathbf{x}_T) = \mathcal{N}(\mathbf{0}, \mathbf{I})$. To generate from the model, we sample a point $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, and then run the Markov chain backwards, by sampling $\mathbf{x}_t \sim p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1})$ until we get a sample in the original data space, \mathbf{x}_0 . From Slide 19 of [KGV22]. Used with kind permission of Arash Vahdat.

allows deep models to be easily trained without any risk of posterior collapse (Section 21.4). In particular, in Section 25.2.3, we will see that training reduces to a series of weighted nonlinear least squares problems.

25.2.1 Encoder (forwards diffusion)

The forwards encoder process is defined to be a simple linear Gaussian model:

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) \quad (25.1)$$

where the values of $\beta_t \in (0, 1)$ are chosen according to a noise schedule (see Section 25.2.4). The joint distribution over all the latent states, conditioned on the input, is given by

$$q(\mathbf{x}_{1:T} | \mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1}) \quad (25.2)$$

Since this defines a linear Gaussian Markov chain, we can compute marginals of it in closed form. In particular, we have

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t | \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}) \quad (25.3)$$

where we define

$$\alpha_t \triangleq 1 - \beta_t, \bar{\alpha}_t = \prod_{s=1}^t \alpha_s \quad (25.4)$$

We choose the noise schedule such that $\bar{\alpha}_T \approx 0$, so that $q(\mathbf{x}_T | \mathbf{x}_0) \approx \mathcal{N}(\mathbf{0}, \mathbf{I})$.

The distribution $q(\mathbf{x}_t | \mathbf{x}_0)$ is known as the **diffusion kernel**. Applying this to the input data distribution and then computing the result unconditional marginals is equivalent to Gaussian convolution:

$$q(\mathbf{x}_t) = \int q_0(\mathbf{x}_0) q(\mathbf{x}_t | \mathbf{x}_0) d\mathbf{x}_0 \quad (25.5)$$

As t increases, the marginals become simpler, as shown in Figure 25.2. In the image domain, this process will first remove high-frequency content (i.e., low-level details, such as texture), and later will remove low-frequency content (i.e., high-level or “semantic” information, such as shape), as shown in Figure 25.1.

25.2.2 Decoder (reverse diffusion)

In the reverse diffusion process, we would like to invert the forwards diffusion process. If we know the input \mathbf{x}_0 , we can derive the reverse of one forwards step as follows:²

$$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1} | \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t \mathbf{I}) \quad (25.6)$$

$$\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) = \frac{\sqrt{\bar{\alpha}_{t-1}} \beta_t}{1 - \bar{\alpha}_t} \mathbf{x}_0 + \frac{\sqrt{\alpha_t} (1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{x}_t \quad (25.7)$$

$$\tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t \quad (25.8)$$

Of course, when generating a new datapoint, we do not know \mathbf{x}_0 , but we will train the generator to approximate the above distribution averaged over \mathbf{x}_0 . Thus we choose the generator to have the form

$$p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1} | \boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t), \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t)) \quad (25.9)$$

We often set $\boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t) = \sigma_t^2 \mathbf{I}$. We discuss how to learn σ_t^2 in Section 25.2.4, but two natural choices are $\sigma_t^2 = \beta_t$ and $\sigma_t^2 = \tilde{\beta}_t$; these correspond to upper and lower bounds on the reverse process entropy, as shown in [HJA20].

The corresponding joint distribution over all the generated variables is given by $p_{\theta}(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)$, where we set $p(\mathbf{x}_T) = \mathcal{N}(\mathbf{0}, \mathbf{I})$. We can sample from this distribution using the pseudocode in Algorithm 25.2.

² We just need to use Bayes’ rule for Gaussians. See e.g., <https://lilianweng.github.io/posts/2021-07-11-diffusion-models/> for a detailed derivation.

1 **25.2.3 Model fitting**

3 We will fit the model by maximizing the evidence lower bound (ELBO), similar to how we train
4 VAEs (see Section 21.2). In particular, for each data example \mathbf{x}_0 we have
5

6
$$\log p_{\theta}(\mathbf{x}_0) = \log \left[\int d\mathbf{x}_{1:T} q(\mathbf{x}_{1:T}|\mathbf{x}_0) \frac{p_{\theta}(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \quad (25.10)$$

7
$$\geq \int d\mathbf{x}_{1:T} q(\mathbf{x}_{1:T}|\mathbf{x}_0) \log \left(p(\mathbf{x}_T) \prod_{t=1}^T \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right) \quad (25.11)$$

8
$$= \mathbb{E}_q \left[\log p(\mathbf{x}_T) + \sum_{t=1}^T \log \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \triangleq \mathcal{L}(\mathbf{x}_0) \quad (25.12)$$

15 We now discuss how to compute the terms in the ELBO. By the Markov property we have
16 $q(\mathbf{x}_t|\mathbf{x}_{t-1}) = q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)$, and by Bayes' rule, we have
17

18
$$q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0) = \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)} \quad (25.13)$$

21 Plugging Equation (25.13) into the ELBO we get
22

23
$$\mathcal{L}(\mathbf{x}_0) = \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log p(\mathbf{x}_T) + \sum_{t=2}^T \log \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} + \underbrace{\sum_{t=2}^T \log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)}}_* + \log \frac{p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)} \right] \quad (25.14)$$

29 The term marked * is a telescoping sum, and can be simplified as follows:
30

31
$$* = \log q(\mathbf{x}_{T-1}|\mathbf{x}_0) + \dots + \log q(\mathbf{x}_2|\mathbf{x}_0) + \log q(\mathbf{x}_1|\mathbf{x}_0) \quad (25.15)$$

32
$$- \log q(\mathbf{x}_T|\mathbf{x}_0) - \log q(\mathbf{x}_{T-1}|\mathbf{x}_0) - \dots - \log q(\mathbf{x}_2|\mathbf{x}_0) \quad (25.16)$$

33
$$= -\log q(\mathbf{x}_T|\mathbf{x}_0) + \log q(\mathbf{x}_1|\mathbf{x}_0) \quad (25.17)$$

36 Hence the negative ELBO (variational upper bound) becomes
37

38
$$\mathcal{L}(\mathbf{x}_0) = -\mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_0)} + \sum_{t=2}^T \log \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} + \log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1) \right] \quad (25.18)$$

39
$$= \underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))}_{L_T(\mathbf{x}_0)} \quad (25.19)$$

40
$$+ \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} \underbrace{D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))}_{L_{t-1}(\mathbf{x}_0)} - \underbrace{\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}_{L_0(\mathbf{x}_0)} \quad (25.20)$$

Each of these KL terms can be computed analytically, since all the distributions are Gaussian. Below we focus on the L_{t-1} term. Since $\mathbf{x}_t = \sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}$, we can write

$$\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \alpha_t}} \boldsymbol{\epsilon} \right) \quad (25.21)$$

Thus instead of training the model to predict the mean of the denoised version of \mathbf{x}_{t-1} given its noisy input \mathbf{x}_t , we can train the model to predict the noise, from which we can compute the mean:

$$\mu_{\theta}(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \alpha_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right) \quad (25.22)$$

With this parameterization, the loss (averaged over the dataset) becomes

$$L_{t-1} = \mathbb{E}_{\mathbf{x}_0 \sim q_0(\mathbf{x}_0), \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[\underbrace{\frac{\beta_t^2}{2\sigma_t^2 \alpha_t (1 - \alpha_t)}}_{\lambda_t} \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta} \left(\underbrace{\sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}}_{\mathbf{x}_t}, t \right)\|^2 \right] \quad (25.23)$$

The time dependent weight λ_t ensures that the training objective corresponds to maximum likelihood training (assuming the variational bound is tight). However, it has been found empirically that the model produces better looking samples if we set $\lambda_t = 1$. The resulting simplified loss (also averaging over time steps t in the model) is given by

$$L_{\text{simple}} = \mathbb{E}_{\mathbf{x}_0 \sim q_0(\mathbf{x}_0), \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), t \sim \text{Unif}(1, T)} \left[\|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta} \left(\underbrace{\sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}}_{\mathbf{x}_t}, t \right)\|^2 \right] \quad (25.24)$$

The overall training procedure is shown in Algorithm 25.1. We can improve the perceptual quality of samples using more advanced weighting schemes, are discussed in [Cho+22]. Conversely, if the goal is to improve likelihood scores, we can optimize the noise schedule, as discussed in Section 25.2.4.

Algorithm 25.1: Training a DDPM model with L_{simple} .

```

1 while not converged do
2    $\mathbf{x}_0 \sim q_0(\mathbf{x}_0)$ 
3    $t \sim \text{Unif}(\{1, \dots, T\})$ 
4    $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5   Take gradient descent step on  $\nabla_{\theta} \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}, t)\|^2$ 

```

After the model is trained, we can generate data using ancestral sampling, as shown in Algorithm 25.2.

25.2.4 Learning the noise schedule

In this section, we describe a way to optimize the noise schedule used by the encoder so as to maximize the ELBO; this approach is called a **variational diffusion model** or **VDM** [Kin+21].

Algorithm 25.2: Sampling from a DDPM model.

```

3 1  $x_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
4 2 foreach  $t = T, \dots, 1$  do
5 3    $\epsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
6 4    $x_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{1-\alpha_t}{\sqrt{1-\alpha_t}} \epsilon_t(x_t, t) \right) + \sigma_t \epsilon_t$ 
7 5 Return  $x_0$ 

```

We will use the following parameterization of the encoder:

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\hat{\mathbf{x}}_t | \hat{\alpha}_t \mathbf{x}_0, \hat{\sigma}_t^2 \mathbf{I}) \quad (25.25)$$

¹⁵ (Note that $\hat{\alpha}_t$ and $\hat{\sigma}_t$ are different to the parameters α_t and σ_t in Section 25.2.1.) Rather than ¹⁶ working with $\hat{\alpha}_t$ and $\hat{\sigma}_t^2$ separately, we will learn to predict their ratio, which is known as the ¹⁷ **signal to noise ratio**:

$$(25.26)$$

²¹ This should be monotonically decreasing in t . This can be ensured by defining $R(t) = \exp(-\gamma_\phi(t))$,
²² where $\gamma_\phi(t)$ is a monotonic neural network. We usually set $\hat{\alpha}_t = \sqrt{1 - \sigma_t^2}$, to correspond to the
²³ variance preserving SDE discussed in Section 25.4.

Following the derivation in Section 25.2.3, the negative ELBO (variational upper bound) can be written as

$$\mathcal{L}(\mathbf{x}_0) = \underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))}_{\text{prior loss}} + \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)}[-\log p_{\boldsymbol{\theta}}(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction loss}} + \underbrace{\mathcal{L}_D(\mathbf{x}_0)}_{\text{diffusion loss}} \quad (25.27)$$

³⁰ where the first two terms are similar to a standard VAE, and the final diffusion loss is given below:³¹

$$\mathcal{L}_D(\mathbf{x}_0) = \frac{1}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \int_0^1 R'(t) \| \mathbf{x}_0 - \hat{\mathbf{x}}_{\theta}(\mathbf{z}_t, t) \|_2^2 dt \quad (25.28)$$

³⁴ where $R'(t)$ is the derivative of the SNR function, and $\mathbf{z}_t = \alpha_t \mathbf{x}_0 + \sigma_t \boldsymbol{\epsilon}_t$. (See [Kin+21] for the derivation.)

³⁶ Since the SNR function is invertible, due to the monotonicity assumption, we can perform a change of
³⁷ variables, and make everything a function of $v = R(t)$ instead of t . In particular, let $\mathbf{z}_v = \alpha_v \mathbf{x}_0 + \sigma_v \boldsymbol{\epsilon}$,
³⁸ and $\hat{\mathbf{x}}_\theta(\mathbf{z}, v) = \hat{\mathbf{x}}_\theta(\mathbf{z}, R^{-1}(v))$. Then we can rewrite Equation (25.28) as
³⁹

$$\frac{40}{41} \quad \mathcal{L}_D(\mathbf{x}_0) = \frac{1}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \int_{R_{\min}}^{R_{\max}} \|\mathbf{x}_0 - \tilde{\mathbf{x}}_{\boldsymbol{\theta}}(\mathbf{z}_v, v)\|_2^2 dv \quad (25.29)$$

⁴³ where $R_{\min} = R(1)$ and $R_{\max} = R(0)$. Thus we see that the shape of the SNR schedule does not
⁴⁴ matter, except for its value at the two end points.

⁴⁶ 3. We present a simplified form of the loss that uses the continuous time limit, which we discuss in Section 25.4.

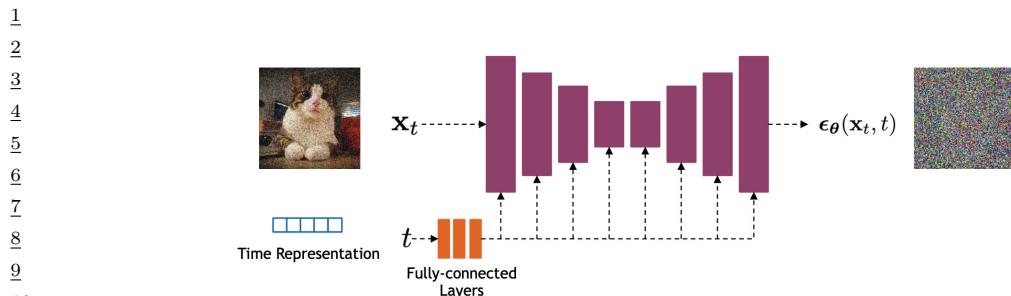


Figure 25.3: Illustration of the U-net architecture used in the denoising step. From Slide 26 of [KGV22]. Used with kind permission of Arash Vahdat.



Figure 25.4: Some sample images generated by a small variational diffusion model trained on EMNIST for about 30 minutes on a K40 GPU. (a) Unconditional sampling. (b) Conditioned on class label. (c) Using classifier-free guidance (see Section 25.6.3). Generated by `diffusion_emnist.ipynb`. Used with kind permission of Alex Alemi.

The integral in Equation (25.29) can be estimated by sampling a timestep uniformly at random. When processing a minibatch of k examples, we can produce a lower variance estimate of the variational bound by using a **low-discrepancy sampler** (cf., Section 11.6.5). In this approach, instead of sampling the timesteps independently, we sample a single uniform random number $u_0 \sim \text{Unif}(0, 1)$, and then set $t^i = \text{mod}(u_0 + i/k, 1)$ for the i 'th sample. We can also optimize the noise schedule wrt the variance of the diffusion loss.

25.2.5 Example: image generation

Diffusion models are often used to generate images. The most common architecture for image generation is based on the **U-net** model [RFB15], as shown in Figure 25.3. The time step t is encoded as a vector, using sinusoidal positional encoding or random Fourier features, and is then fed into the residual blocks, using either simple spatial addition or by conditioning the group norm layers [DN21a]. Of course, other architectures besides U-net are possible. For example, recently [PX22; Li+22; Bao+22a] have proposed the use of transformers, to replace the convolutional and

1 deconvolutional layers.
2

3 The results of training a small U-net VDM on EMNIST images are shown in Figure 25.4. By
4 training big models (billions of parameters) for a long time (days) on lots of data (millions of images),
5 diffusion models can be made to generate very high quality images (see Figure 20.2). Results can
6 be further improved by using conditional diffusion models, where guidance is provided about what
7 kinds of images to generate (see Section 25.6).

8

9 10 25.3 Score-based generative models (SGMs)

11 This section is written with Yang Song and Durk Kingma.
12

13 In Section 24.3, we discussed how to fit energy based models (EBMs) using score matching. This
14 adjusts the parameters of the EBM so that the score function of the model, $\nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x})$, matches
15 the score function of the data, $\nabla_{\mathbf{x}} \log p_{\mathcal{D}}(\mathbf{x})$. An alternative to fitting a scalar energy function and
16 computing its score is to directly learn the score function. This is called a **score-based generative**
17 **model** or **SGM** [SE19; SE20b; Son+21b]. We can optimize the score function $s_{\theta}(\mathbf{x})$ using basic
18 score matching (Section 24.3.1), sliced score matching (Section 24.3.3 or denoising score matching
19 (Section 24.3.2). We discuss this class of models in more detail below. (For a comparison with EBMs,
20 see [SH21].)

21

22 23 25.3.1 Example

24 In Figure 25.5a, we show the **Swiss roll** dataset. We estimate the score function by fitting an MLP
25 with 2 hidden layers, each with 128 hidden units, using basic score matching. In Figure 25.5b, we
26 show the output of the network after training for 10,000 steps of SGD. We see that there are no
27 major false negatives (since wherever the density of the data is highest, the gradient field is zero),
28 but there are some false positives (since some regions of zero gradient do not correspond to data
29 regions). The comparison of the predicted outputs with the empirical data density is shown more
30 clearly in Figure 25.5c. In Figure 25.5d, we show some samples from the learned model, generated
31 using Langevin sampling.

32

33

34 25.3.2 Adding noise at multiple scales

35 In general, score matching can have difficulty when there are regions of low data density. To see this,
36 suppose $p_{\mathcal{D}}(\mathbf{x}) = \pi p_0(\mathbf{x}) + (1 - \pi)p_1(\mathbf{x})$. Let $\mathcal{S}_0 := \{\mathbf{x} \mid p_0(\mathbf{x}) > 0\}$ and $\mathcal{S}_1 := \{\mathbf{x} \mid p_1(\mathbf{x}) > 0\}$ be the
37 supports of $p_0(\mathbf{x})$ and $p_1(\mathbf{x})$ respectively. When they are disjoint from each other, the score of $p_{\mathcal{D}}(\mathbf{x})$
38 is given by

$$\nabla_{\mathbf{x}} \log p_{\mathcal{D}}(\mathbf{x}) = \begin{cases} \nabla_{\mathbf{x}} \log p_0(\mathbf{x}), & \mathbf{x} \in \mathcal{S}_0 \\ \nabla_{\mathbf{x}} \log p_1(\mathbf{x}), & \mathbf{x} \in \mathcal{S}_1, \end{cases} \quad (25.30)$$

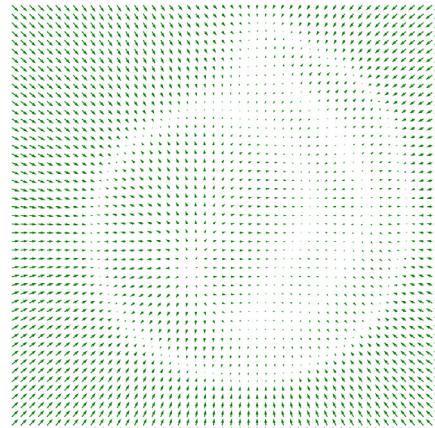
43

44 which does not depend on the weight π . Hence score matching cannot correctly recover the true
45 distribution. Furthermore, Langevin sampling will have difficulty traversing between modes. (In
46 practice this will happen even when the different modes only have approximately disjoint supports.)

47

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20

(a)



(b)

A scatter plot showing a positive linear relationship between two variables. The x-axis ranges from approximately -1.5 to 1.5, and the y-axis ranges from 0 to 10. A dense cloud of blue data points follows the equation $y = 2x + 1$. A red dashed line represents the regression line $y = 2x + 1$. A green dotted grid is overlaid on the plot.

(c)



(d)

Figure 25.5: Fitting a score-based generative model to the 2d Swiss roll dataset. (a) Training set. (b) Learned score function trained using the basic score matching. (c) Superposition of learned score function and empirical density. (d) Langevin sampling applied to the learned model. We show 3 different trajectories, each of length 25. Generated by `score_matching_swiss_roll.ipynb`

43
44
45
46
47

¹ Song and Ermon [SE19; SE20b] and Song et al. [Son+21b] overcome this difficulty by perturbing
² training data with different scales of noise. Specifically, they use
³

$$\frac{4}{5} q_\sigma(\tilde{\mathbf{x}}|\mathbf{x}) = \mathcal{N}(\tilde{\mathbf{x}}|\mathbf{x}, \sigma^2 \mathbf{I}) \quad (25.31)$$

$$\frac{6}{7} q_\sigma(\tilde{\mathbf{x}}) = \int p_{\mathcal{D}}(\mathbf{x}) q_\sigma(\tilde{\mathbf{x}}|\mathbf{x}) d\mathbf{x} \quad (25.32)$$

⁸ For a large noise perturbation, different modes are connected due to added noise, and the estimated
⁹ weights between them are therefore accurate. For a small noise perturbation, different modes are
¹⁰ more disconnected, but the noise-perturbed distribution is closer to the original unperturbed data
¹¹ distribution. Using a sampling method such as annealed Langevin dynamics [SE19; SE20b; Son+21b]
¹² or diffusion sampling [SD+15a; HJA20; Son+21b], we can sample from the most noise-perturbed
¹³ distribution first, then smoothly reduce the magnitude of noise scales until reaching the smallest one.
¹⁴ This procedure helps combine information from all noise scales, and maintains the correct estimation
¹⁵ of weights from higher noise perturbations when sampling from smaller ones.

¹⁶ In practice, all score models share weights and are implemented with a single neural network
¹⁷ conditioned on the noise scale; this is called a **noise conditional score network**, and has the form
¹⁸ $\mathbf{s}_\theta(\mathbf{x}, \sigma)$. Scores of different scales are estimated by training a mixture of score matching objectives,
¹⁹ one per noise scale. If we use the denoising score matching objective in Equation (24.33), we get
²⁰

$$\frac{21}{22} \mathcal{L}(\theta; \sigma) = \mathbb{E}_{q(\mathbf{x}, \tilde{\mathbf{x}})} \left[\frac{1}{2} \|\nabla_{\mathbf{x}} \log p_\theta(\tilde{\mathbf{x}}, \sigma) - \nabla_{\mathbf{x}} \log q_\sigma(\tilde{\mathbf{x}}|\mathbf{x})\|_2^2 \right] \quad (25.33)$$

$$\frac{24}{25} = \frac{1}{2} \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathcal{N}(\mathbf{x}, \sigma^2 \mathbf{I})} \left\{ \left\| \mathbf{s}_\theta(\tilde{\mathbf{x}}, \sigma) + \frac{(\tilde{\mathbf{x}} - \mathbf{x})}{\sigma^2} \right\|_2^2 \right\} \quad (25.34)$$

²⁶ where we used the fact that, for a Gaussian, the score is given by
²⁷

$$\frac{29}{30} \nabla_{\mathbf{x}} \log \mathcal{N}(\tilde{\mathbf{x}}|\mathbf{x}, \sigma^2 \mathbf{I}) = -\nabla_{\mathbf{x}} \frac{1}{2\sigma^2} (\mathbf{x} - \tilde{\mathbf{x}})^\top (\mathbf{x} - \tilde{\mathbf{x}}) = \frac{\mathbf{x} - \tilde{\mathbf{x}}}{\sigma^2} \quad (25.35)$$

³¹ If we have T different noise scales, we can combine the losses in a weighted fashion using
³²

$$\frac{33}{34} \mathcal{L}(\theta; \sigma_{1:T}) = \sum_{t=1}^T \lambda_t \mathcal{L}(\theta; \sigma_t) \quad (25.36)$$

³⁵ where we choose $\sigma_1 > \sigma_2 > \dots > \sigma_T$, and the weighting term satisfies $\lambda_t > 0$.

³⁸ 25.3.3 Equivalence to DDPM

³⁹ We now show that the above score-based generative model training objective is equivalent to the
⁴⁰ DDPM loss. To see this, first let us replace $p_{\mathcal{D}}(\mathbf{x})$ with $q_0(\mathbf{x}_0)$, $\tilde{\mathbf{x}}$ with \mathbf{x}_t , and $\mathbf{s}_\theta(\tilde{\mathbf{x}}, \sigma)$ with $\mathbf{s}_\theta(\mathbf{x}_t, t)$.
⁴¹ We will also compute a stochastic approximation to Equation (25.36) by sampling a time step
⁴² uniformly at random. Then Equation (25.36) becomes
⁴³

$$\frac{44}{45} \mathcal{L} = \mathbb{E}_{\mathbf{x}_0 \sim q_0(\mathbf{x}_0), \mathbf{x}_t \sim q(\mathbf{x}_t|\mathbf{x}_0), t \sim \text{Unif}(1, T)} \left[\lambda_t \left\| \mathbf{s}_\theta(\mathbf{x}_t, t) + \frac{(\mathbf{x}_t - \mathbf{x}_0)}{\sigma_t^2} \right\|_2^2 \right] \quad (25.37)$$

If we use the fact that $\mathbf{x}_t = \mathbf{x}_0 + \sigma_t \boldsymbol{\epsilon}$, and if we define $s_{\theta}(\mathbf{x}_t, t) = -\frac{\epsilon_{\theta}(\mathbf{x}_t, t)}{\sigma_t}$, we can rewrite this as

$$\mathcal{L} = \mathbb{E}_{\mathbf{x}_0 \sim q_0(\mathbf{x}_0), \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), t \sim \text{Unif}(1, T)} \left[\frac{\lambda_t}{\sigma_t^2} \|\boldsymbol{\epsilon} - \epsilon_{\theta}(\mathbf{x}_t, t)\|_2^2 \right] \quad (25.38)$$

If we set $\lambda_t = \sigma_t^2$, we recover L_{simple} loss in Equation (25.24).

25.4 Continuous time models using differential equations

In this section, we consider a DDPM model in the limit of an infinite number of hidden layers, or equivalently, an SGM in the limit of an infinite number of noise levels. This requires switching from discrete time to continuous time, which complicates the mathematics. The advantage is that we can leverage the large existing literature on solvers for ordinary and stochastic differential equations to enable faster generation, as we will see.

25.4.1 Forwards diffusion SDE

Let us first consider a diffusion process where the noise level β_t gets rewritten as $\beta(t)\Delta t$, where Δt is a step size:

$$\mathbf{x}_t = \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \sqrt{\beta_t} \mathcal{N}(\mathbf{0}, \mathbf{I}) = \sqrt{1 - \beta(t)\Delta t} \mathbf{x}_{t-1} + \sqrt{\beta(t)\Delta t} \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (25.39)$$

If Δt is small, we can approximate the first term with a first-order Taylor series expansion to get

$$\mathbf{x}_t \approx \mathbf{x}_{t-1} - \frac{\beta(t)\Delta t}{2} \mathbf{x}_{t-1} + \sqrt{\beta(t)\Delta t} \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (25.40)$$

Hence for small Δt we have

$$\frac{\mathbf{x}_t - \mathbf{x}_{t-1}}{\Delta t} \approx -\frac{\beta(t)}{2} \mathbf{x}_{t-1} + \frac{\sqrt{\beta(t)}}{\sqrt{\Delta t}} \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (25.41)$$

We can now switch to the **continuous time** limit, and write this as the following **stochastic differential equation** or **SDE**:

$$\frac{d\mathbf{x}(t)}{dt} = -\frac{1}{2}\beta(t)\mathbf{x}(t) + \sqrt{\beta(t)} \frac{d\mathbf{w}(t)}{dt} \quad (25.42)$$

where $\mathbf{w}(t)$ represents a standard **Wiener process**, also called **Brownian noise**. More generally, we can write such SDEs as follows, where we use **Itô calculus** notation (see e.g., [SS19]):

$$d\mathbf{x} = \underbrace{\mathbf{f}(\mathbf{x}, t)}_{\text{drift}} dt + \underbrace{\mathbf{g}(t)}_{\text{diffusion}} d\mathbf{w} \quad (25.43)$$

The first term in the above SDE is called the **drift coefficient**, and the second term is called the **diffusion coefficient**.

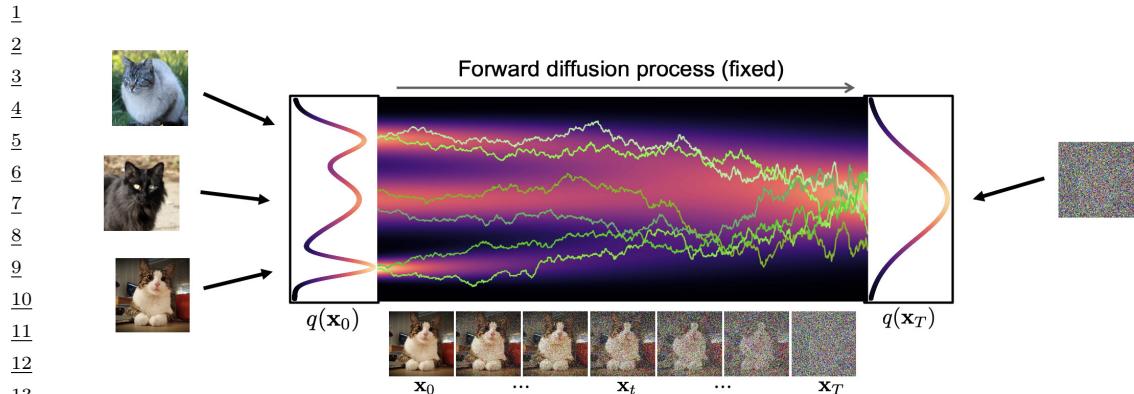


Figure 25.6: Illustration of the forwards diffusion process in continuous time. Yellow lines are sample paths from the SDE. Heat map represents the marginal distribution computed using the probability flow ODE. From Slide 43 of [KGV22]. Used with kind permission of Karsten Kreis.

We can gain some intuition for these processes by looking at the 1d example in Figure 25.6. We can draw multiple paths as follows: sample an initial state from the data distribution, and then integrate over time using **Euler-Maruyama** integration:

$$\mathbf{x}(t + \Delta t) = \mathbf{x}(t) + \mathbf{f}(\mathbf{x}(t), t)\Delta t + g(t)\sqrt{\Delta t}\mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (25.44)$$

We can see how the data distribution at $t = 0$, on the left hand side, gradually gets transformed to a pure noise distribution at $t = 1$, on the right hand side.

In [Son+21b], they show that the SDE corresponding to DDPMs, in the $T \rightarrow \infty$ limit, is given by

$$d\mathbf{x} = -\frac{1}{2}\beta(t)\mathbf{x}dt + \sqrt{\beta(t)}d\omega \quad (25.45)$$

where $\beta(t/T) = T\beta_t$. Here the drift term is proportional to $-\mathbf{x}$, which encourages the process to return to 0. Consequently, DDPM corresponds to a **variance preserving** process. By contrast, the SDE corresponding to SGMs is given by the following:

$$d\mathbf{x} = \sqrt{\frac{d[\sigma(t)^2]}{dt}}d\omega \quad (25.46)$$

where $\sigma(t/T) = \sigma_t$. This SDE has zero drift, so corresponds to a **variance exploding** process.

25.4.2 Forwards diffusion ODE

Instead of adding Gaussian noise at every step, we can just sample the initial state, and then let it evolve deterministically over time according to the following **ordinary differential equation or ODE**:

$$d\mathbf{x} = \underbrace{\left[f(\mathbf{x}, t) - \frac{1}{2}g(t)^2\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) \right]}_{h(\mathbf{x}, t)} dt \quad (25.47)$$

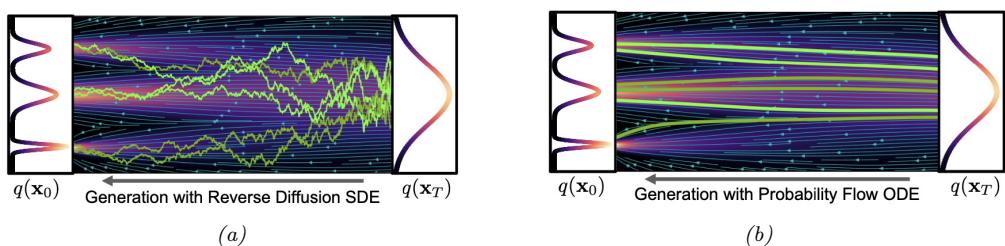


Figure 25.7: Illustration of the reverse diffusion process. (a) Sample paths from the SDE. (b) Deterministic trajectories from the probability flow ODE. From Slide 65 of [KGV22]. Used with kind permission of Karsten Kreis.

This is called the **probability flow ODE** [Son+21b, Sec D.3]. We can compute the state at any moment in time using any ODE solver:

$$\mathbf{x}(t) = \mathbf{x}(0) + \int_0^t h(\mathbf{x}, \tau) d\tau \quad (25.48)$$

See Figure 25.7b for a visualization of a sample trajectory. If we start the solver from different random states $\mathbf{x}(0)$, then the induced distribution over paths will have the same marginals as the SDE model. See the heatmap in Figure 25.6 for an illustration.

25.4.3 Reverse diffusion SDE

To generate samples from this model, we need to be able to reverse the SDE. In a remarkable result, [And82] showed that any forwards SDE of the form in Equation (25.43) can be reversed to get the following **reverse-time SDE**:

$$d\mathbf{x} = [f(\mathbf{x}, t) - g(t)^2 \nabla_{\mathbf{x}} \log q_t(\mathbf{x})] dt + g(t) d\bar{\mathbf{w}} \quad (25.49)$$

where $\bar{\mathbf{w}}$ is the standard Wiener process when time flows backwards, dt is an infinitesimal negative time step, and $\nabla_{\mathbf{x}} \log q_t(\mathbf{x})$ is the score function.

In the case of the DDPM, the reverse SDE has the following form:

$$d\mathbf{x}_t = \left[-\frac{1}{2}\beta(t)\mathbf{x}_t - \beta(t)\nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t) \right] dt + \sqrt{\beta(t)} d\bar{\mathbf{w}}_t \quad (25.50)$$

To estimate the score function, we can use denoising score matching as we discussed in Section 25.3, to get

$$\nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t) \approx \mathbf{s}_{\theta}(\mathbf{x}_t, t) \quad (25.51)$$

(In practice, it is advisable to use variance reduction techniques, such as importance sampling, as discussed in [Son+21a].) The SDE becomes

$$d\mathbf{x}_t = -\frac{1}{2}\beta(t)[\mathbf{x}_t + 2\mathbf{s}_{\theta}(\mathbf{x}_t, t)] dt + \sqrt{\beta(t)} d\bar{\mathbf{w}}_t \quad (25.52)$$

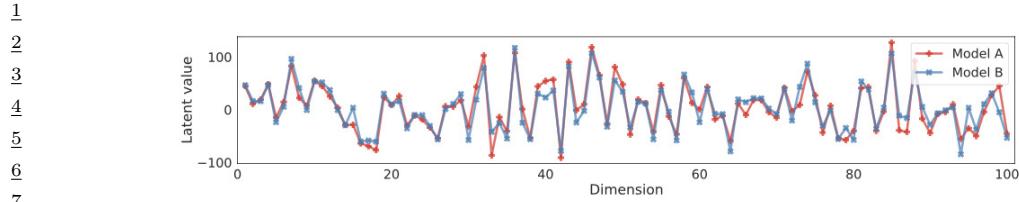


Figure 25.8: Comparing the first 100 dimensions of the latent code obtained for a random CIFAR-100 image. “Model A” and “Model B” are separately trained with different architectures. From Figure 7 of [Son+21b]. Used with kind permission of Yang Song.

After fitting the score network, we can sample from it using ancestral sampling (as in Section 25.2), or we can use the Euler-Maruyama integration scheme in Equation (25.44), which gives

$$\mathbf{x}_{t-1} = \mathbf{x}_t + \frac{1}{2}\beta(t)[\mathbf{x}_t + 2\mathbf{s}_\theta(\mathbf{x}_t, t)]\Delta t + \sqrt{\beta(t)\Delta t}\mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (25.53)$$

See Figure 25.7a for an illustration.

25.4.4 Reverse diffusion ODE

Based on the results in Section 25.4.2, we can derive the probability flow ODE from the reverse-time SDE in Equation (25.49) to get

$$d\mathbf{x}_t = \left[f(\mathbf{x}, t) - \frac{1}{2}g(t)^2\mathbf{s}_\theta(\mathbf{x}_t, t) \right] dt \quad (25.54)$$

If we set $f(\mathbf{x}, t) = -\frac{1}{2}\beta(t)$ and $g(t) = \sqrt{\beta(t)}$, as in DDPM, this becomes

$$d\mathbf{x}_t = -\frac{1}{2}\beta(t)[\mathbf{x}_t + \mathbf{s}_\theta(\mathbf{x}_t, t)]dt \quad (25.55)$$

See Figure 25.7b for an illustration. A simple way to solve this ODE is to use **Euler’s method**:

$$\mathbf{x}_{t-1} = \mathbf{x}_t + \frac{1}{2}\beta(t)[\mathbf{x}_t + \mathbf{s}_\theta(\mathbf{x}_t, t)]\Delta t \quad (25.56)$$

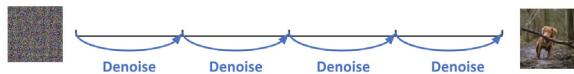
However, in practice one can get better results using higher-order ODE solvers, such as **Heun’s method** [Kar+22].

This model is a special case of a **neural ODE**, also called a **continuous normalizing flow** (see Section 23.2.6). Consequently we can derive the exact log marginal likelihood. However, instead of maximizing this directly (which is expensive), we use score matching to fit the model.

Another advantage of the deterministic ODE approach is that it guarantees that the generative model is **identifiable**. To see this, note that the ODE (in both forwards and reverse directions) is deterministic, and is uniquely determined by the score function. If the architecture is sufficiently flexible, and if there is enough data, then score matching will recover the true score function of the data generating process. Thus, after training, a given datapoint will map to a unique point in latent space, regardless of the model architecture or initialization (see Figure 25.8).

1
2

Deterministic sampling:



Stochastic sampling

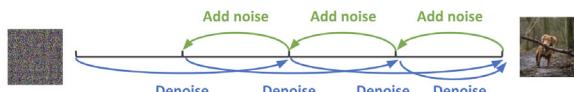


Figure 25.9: Generating from the reverse diffusion process using 4 steps. (Top) Deterministic sampling. (Bottom) A mix of deterministic and stochastic sampling. Used with kind permission of Ruiqi Gao.

15 Furthermore, since every point in latent space decodes to a unique image, we can perform “semantic
16 interpolation” in the latent space to generate images with properties that are in between two input
17 examples (cf., Figure 20.9).

25.4.5 Comparison of the SDE and ODE approach

In Section 25.4.3 we described the reverse diffusion process as an SDE, and in Section 25.4.4, we described it as an ODE. We can see the connection between these methods by rewriting the SDE in Equation (25.49) as follows:

$$dx_t = \underbrace{-\frac{1}{2}\beta(t)[x_t + s_{\theta}(x_t, t)]dt}_{\text{probability flow ODE}} - \underbrace{\frac{1}{2}\beta(t)s_{\theta}(x_t, t)dt + \sqrt{\beta(t)}d\bar{w}_t}_{\text{Langevin diffusion SDE}} \quad (25.57)$$

28 The continuous noise injection can compensate for errors introduced by the numerical integration of
 29 the ODE term. Consequently, the resulting samples often look better. However, the ODE approach
 30 can be faster. Fortunately it is possible to combine these techniques, as proposed in [Kar+22]. The
 31 basic idea is illustrated in Figure 25.9: we alternate between performing a deterministic step using
 32 an ODE solver, and then adding a small amount noise to the result. This can be repeated for some
 33 number of steps. (We discuss ways to reduce the number of required steps in Section 25.5.)

35 25.4.6 Example

36 A simple JAX implementation of the above ideas, written by Winnie Xu, can be found in [diffusion_mnist.ipynb](#). This fits a small model to MNIST images using denoising score matching. It then
37 generates from the model by solving the probability flow ODE using the diffrax library. By scaling
38 this kind of method up to a much larger model, and training for a much longer time, it is possible to
39 produce very impressive looking results, as shown in Figure 25.10.

43 25.5 Speeding up diffusion models

45 One of the main disadvantages of diffusion models is that generating from them takes many small
46 steps, which can be slow. While it is possible to just take fewer, larger steps, the results are much

1
2
3
4
5
6
7
8
9
10



11 *Figure 25.10: Synthetic faces from a score-based generative model trained on CelebA-HQ-256 images. From*
12 *Figure 12 of [Son+21b]. Used with kind permission of Yang Song.*

13
14

15 worse. In this section, we briefly mention a few techniques that have been proposed to tackle this
16 important problem. Many other techniques are mentioned in the recent review papers [UAP22;
17 Yan+22; Cao+22].

18
19

20 25.5.1 DDIM sampler

21 In this section, we describe the denoising diffusion implicit model or **DDIM** of [SME21], which can
22 be used for efficient deterministic generation. The first step is to use a **non-Markovian** forwards
23 diffusion process, so it always conditions on the input in addition to the previous step:

$$24 \quad q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}\left(\sqrt{\bar{\alpha}_{t-1}} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_{t-1} - \tilde{\sigma}_t^2} \frac{\mathbf{x}_t - \sqrt{\bar{\alpha}_t} \mathbf{x}_0}{\sqrt{1 - \bar{\alpha}_t}}, \tilde{\sigma}_t^2 \mathbf{I}\right) \quad (25.58)$$

25 The corresponding reverse process is
26

$$27 \quad p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}\left(\sqrt{\bar{\alpha}_{t-1}} \hat{\mathbf{x}}_0 + \sqrt{1 - \bar{\alpha}_{t-1} - \tilde{\sigma}_t^2} \frac{\mathbf{x}_t - \sqrt{\bar{\alpha}_t} \hat{\mathbf{x}}_0}{\sqrt{1 - \bar{\alpha}_t}}, \tilde{\sigma}_t^2 \mathbf{I}\right) \quad (25.59)$$

28 where $\hat{\mathbf{x}}_0 = \hat{\mathbf{x}}_{\theta}(\mathbf{x}_t, t)$ is the predicted output from the model. By setting $\tilde{\sigma}_t^2 = 0$, the reverse process
29 becomes fully deterministic, given the initial prior sample (whose variance is controlled by $\tilde{\sigma}_T^2$). The
30 resulting probability flow ODE gives better results when using a small number of steps compared to
31 the methods discussed in Section 25.4.4.

32 Note that the weighted negative VLB for this model is the same as L_{simple} in Section 25.2, so the
33 DDIM sampler can be applied to a trained DDPM model.

34
35
36
37
38
39

40 25.5.2 Non-Gaussian decoder networks

41 If the reverse diffusion process takes larger steps, then the induced distribution over clean outputs
42 given a noisy input will become multimodal, as illustrated in Figure 25.11. This requires more
43 complicated forms for the distribution $p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)$. In [Gao+21], they use an EBM to fit this
44 distribution. However, this still requires the use of MCMC to draw a sample. In [XKV22], they
45 use a GAN (Chapter 26) to fit this distribution. This enables us to easily draw a sample by passing
46
47

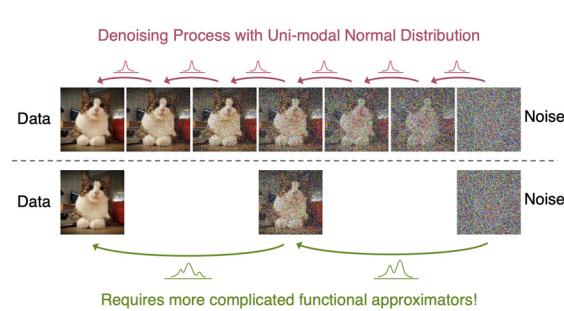


Figure 25.11: Illustration of why taking larger steps in the reverse diffusion process needs more complex, multi-modal conditional distributions. From Slide 90 of [KGV22]. Used with kind permission of Arash Vahdat.

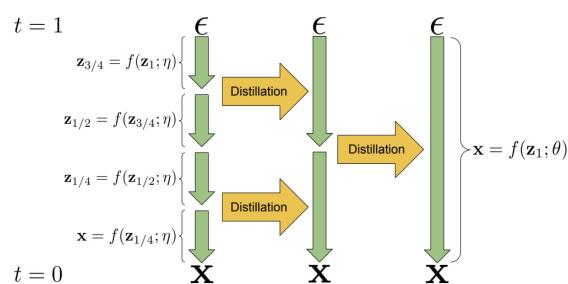
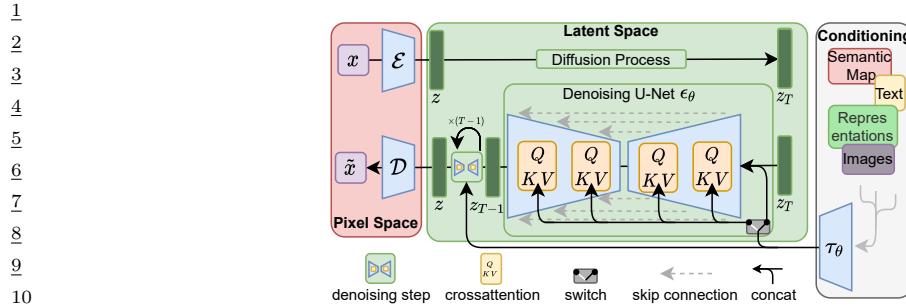


Figure 25.12: Progressive distillation. From Figure 1 of [SH22]. Used with kind permission of Tim Salimans.

Gaussian noise through the generator. The benefits over a single stage GAN is that both the generator and discriminator are solving a much simpler problem, resulting in increased mode coverage, and better training stability. The benefit over a standard diffusion model is that we can generate high quality samples in many fewer steps.

25.5.3 Distillation

In this section, we discuss the **progressive distillation** method of [SH22], which provides a way to create a diffusion model that only needs a small number of steps to create high quality samples. The basic idea is follows. First we train a DDPM model in the usual way, and sample from it using the DDIM method; we treat this as the “teacher” model. We use this to generate intermediate latent states, and train a “student” model to predict the output of the teacher on every second step, as shown in Figure 25.12. After the student has been trained, it can generate results that are as good as the teacher, but in half the number of steps. This student can then teach a new generation of even faster students. See Algorithm 25.4 for the pseudocode, which should be compared to Algorithm 25.3 for the standard training procedure. Note that each round of teaching becomes faster, because the teachers become smaller, so the total time to perform the distillation is relatively small. The resulting model can generate high quality samples in as few as 4 steps.



11 *Figure 25.13: Combining a VAE with a diffusion model. Here \mathcal{E} and \mathcal{D} are the encoder and decoder of the*
12 *VAE. The diffusion model conditions on the inputs either by using concatenation or by using a cross-attention*
13 *mechanism. From Figure 3 of [Rom+22]. Used with kind permission of Robin Rombach.*

Algorithm 25.3: Standard training

19 **Input:** Model $\hat{x}_\theta(z_t)$ to be trained
20 **Input:** Dataset \mathcal{D}
21 **Input:** Loss weight function w

```

23 1 while not converged do
24   2    $x \sim \mathcal{D}$ 
25   3    $t \sim \text{Unif}(0, 1)$ 
26   4    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
27   5    $z_t = \alpha_t x + \sigma_t \epsilon$ 
28
29
30
31   6    $\tilde{x} = x$  (Clean data is target)
32   7    $\lambda_t = \log(\alpha_t^2 / \sigma_t^2)$  (Log SNR)
33   8    $L_\theta = w(\lambda_t) \|\tilde{x} - \hat{x}_\theta(z_t)\|_2^2$ 
34   9    $\theta := \theta - \gamma \nabla_\theta L_\theta$ 

```

Algorithm 25.4: Progressive distillation

1 **foreach** K iterations **do**
2 $\theta := \eta$ (Assign student)
3 **while** not converged **do**
4 $x \sim \mathcal{D}$
5 $t = i/N, i \sim \text{Cat}(1, 2, \dots, N)$
6 $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
7 $z_t = \alpha_t x + \sigma_t \epsilon$
8 $t' = t - 0.5/N, t'' = t - 1/N$
9 $z_{t'} = \alpha_{t'} \hat{x}_\eta(z_t) + \frac{\sigma_{t'}}{\sigma_t} (z_t - \alpha_t \hat{x}_\eta(z_t))$
10 $z_{t''} = \alpha_{t''} \hat{x}_\eta(z_{t'}) + \frac{\sigma_{t''}}{\sigma_{t'}} (z_{t'} - \alpha_{t'} \hat{x}_\eta(z_{t'}))$
11 $\tilde{x} = \frac{z_{t''} - (\sigma_{t''}/\sigma_t) z_t}{\alpha_{t''} - (\sigma_{t''}/\sigma_t) \alpha_t}$ (Teacher is target)
12 $\lambda_t = \log(\alpha_t^2 / \sigma_t^2)$
13 $L_\theta = w(\lambda_t) \|\tilde{x} - \hat{x}_\theta(z_t)\|_2^2$
14 $\theta := \theta - \gamma \nabla_\theta L_\theta$
15 $\eta := \theta$ (Student becomes next teacher)
16 $N := N/2$ (Halve number of sampling steps)

25.5.4 Latent space diffusion

42 Another approach to speeding up diffusion models for images is to first embed the images into a
43 lower dimensional space, and then fit the diffusion model to the embeddings. This idea has been
44 pursued in several papers.

45 In the latent diffusion model (**LDM**) of [Rom+22], they adopt a two-stage training scheme, in
46 which they first fit the VAE, augmented with a perceptual loss, and then fit the diffusion model to
47

the embedding. The architecture is illustrated in Figure 25.13. The LDM forms the foundation of the very popular **stable diffusion** system created by **Stability AI**. In the latent score-based generative model (**LSGM**) of [VKK21], they first train a hierarchical VAE, and then jointly train the VAE and a diffusion model.

In addition to speed, an additional advantage of combining diffusion models with autoencoders is that it makes it simple to apply diffusion to many different kinds of data, such as text and graphs: we just need to define a suitable architecture to embed the input domain into a continuous space. Note, however, that it is also possible to define diffusion directly on discrete state spaces, as we discuss in Section 25.7.

So far we have discussed applying diffusion “on top of” a VAE. However, we can also do the reverse, and fit a VAE on top of a DDPM model, where we use the diffusion model to “post process” blurry samples coming from the VAE. See [Pan+22] for details.

25.6 Conditional generation

In this section, we discuss how to generate samples from a diffusion model where we condition on some side information \mathbf{c} , such as a class label or text prompt.

25.6.1 Conditional diffusion model

The simplest way to control the generation from a generative model is to train it on (\mathbf{c}, \mathbf{x}) pairs so as to maximize the conditional likelihood, $p(\mathbf{x}|\mathbf{c})$. If the conditioning signal \mathbf{c} is a scalar (e.g., a class label), it can be mapped to an embedding vector, and then incorporated into the network by spatial addition, or by using it to modular the group normalization layers. If the input \mathbf{c} is another image, we can simply concatenate it with \mathbf{x}_t as an extra set of channels. If the input \mathbf{c} is a text prompt, we can embed it, and then use spatial addition or cross-attention (see Figure 25.13 for an illustration).

25.6.2 Classifier guidance

One problem with conditional diffusion models is that we need to retrain them for each kind of conditioning that we want to perform. An alternative approach, known as **classifier guidance** was proposed in [DN21b], and allows us to leverage pre-trained discriminative classifiers of the form $p_\phi(\mathbf{c}|\mathbf{x})$ to control the generation process. The idea is as follows. First we use Bayes’ rule to write

$$\log p(\mathbf{x}|\mathbf{c}) = \log p(\mathbf{c}|\mathbf{x}) + \log p(\mathbf{x}) - \log p(\mathbf{c}) \quad (25.60)$$

from which the score function becomes

$$\nabla_{\mathbf{x}} \log p(\mathbf{x}|\mathbf{c}) = \nabla_{\mathbf{x}} \log p(\mathbf{x}) + \nabla_{\mathbf{x}} \log p(\mathbf{c}|\mathbf{x}) \quad (25.61)$$

We can now use this conditional score to generate samples, rather than the unconditional score. We can further amplify the influence of the conditioning signal by scaling it by a factor $w > 1$:

$$\nabla_{\mathbf{x}} \log p_w(\mathbf{x}|\mathbf{c}) = \nabla_{\mathbf{x}} \log p(\mathbf{x}) + w \nabla_{\mathbf{x}} \log p(\mathbf{c}|\mathbf{x}) \quad (25.62)$$

In practice, this can be achieved as follows by generating samples from

$$\mathbf{x}_{t-1} \sim \mathcal{N}(\boldsymbol{\mu} + w\boldsymbol{\Sigma}\mathbf{g}, \boldsymbol{\Sigma}), \quad \boldsymbol{\mu} = \boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t), \quad \boldsymbol{\Sigma} = \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t), \quad \mathbf{g} = \nabla_{\mathbf{x}_t} \log p_{\phi}(\mathbf{c}|\mathbf{x}_t) \quad (25.63)$$

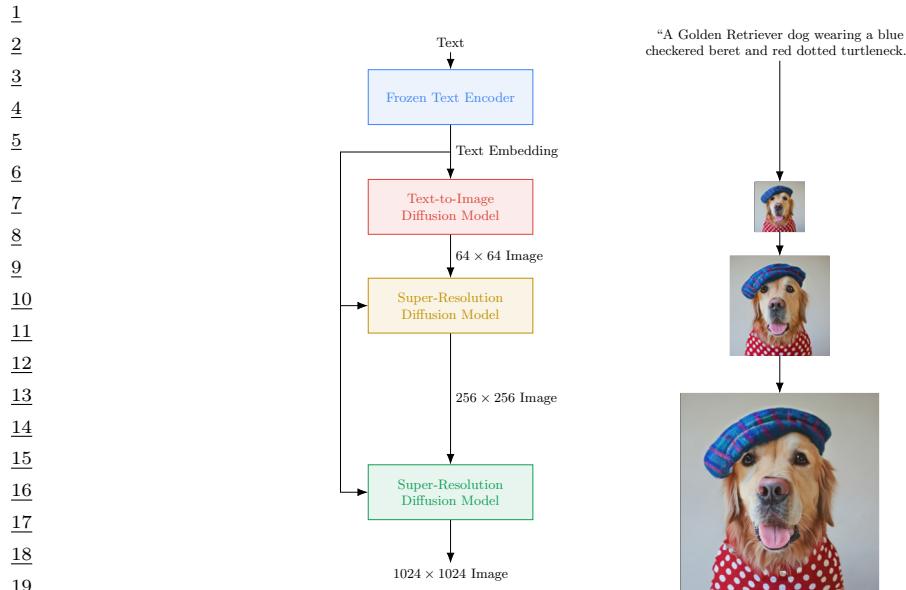


Figure 25.14: Cascaded diffusion model used by the Imagen text-to-image system. From Figure A.5 of [Sah+22b]. Used with kind permission of Saurabh Saxena.

25.6.3 Classifier-free guidance

Unfortunately, $p(\mathbf{c}|\mathbf{x}_t)$ is a discriminative model, that may ignore many details of the input \mathbf{x}_t . Hence optimizing along the directions specified by $\nabla_{\mathbf{x}_t} \log p(\mathbf{c}|\mathbf{x}_t)$ can give poor results, similar to what happens when we create adversarial images. In addition, we need to train a classifier for each time step, since \mathbf{x}_t will differ in its blurriness.

In [HS21], they proposed a technique called **classifier-free guidance**, which derives the classifier from the generative model, using $p(\mathbf{c}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{c})p(\mathbf{c})}{p(\mathbf{x})}$, from which we get

$$\log p(\mathbf{c}|\mathbf{x}) = \log p(\mathbf{x}|\mathbf{c}) + \log p(\mathbf{c}) - \log p(\mathbf{x}) \quad (25.64)$$

This requires learning two generative models, namely $p(\mathbf{x}|\mathbf{c})$ and $p(\mathbf{x})$. However, in practice we can use the same model for this, and simply set $\mathbf{c} = \emptyset$ to represent the unconditional case. We then use this implicit classifier to get the following modified score function:

$$\nabla_{\mathbf{x}} [\log p(\mathbf{x}|\mathbf{c}) + w \log p(\mathbf{c}|\mathbf{x})] = \nabla_{\mathbf{x}} [\log p(\mathbf{x}|\mathbf{c}) + w(\log p(\mathbf{x}|\mathbf{c}) - \log p(\mathbf{x}))] \quad (25.65)$$

$$= \nabla_{\mathbf{x}} [(1 + w) \log p(\mathbf{x}|\mathbf{c}) - w \log p(\mathbf{x})] \quad (25.66)$$

Larger guidance weight usually results in better individual sample quality, but lower diversity.

25.6.4 Generating high resolution images

In order to generate high resolution images, [Ho+21] proposed to use **cascaded generation**, in which we first train a model to generate 64×64 images, and then train a separate **super-resolution**

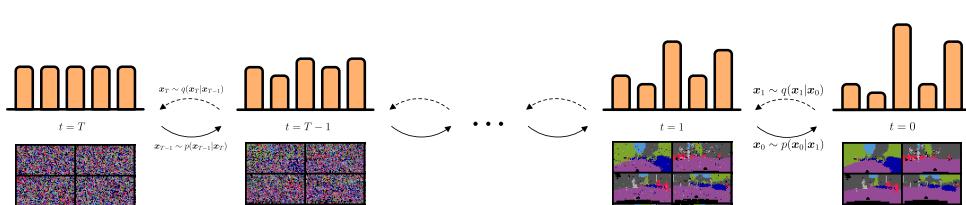


Figure 25.15: Multinomial diffusion model, applied to semantic image segmentation. The input image is on the right, and gets diffused to the noise image on the left. From Figure 1 of [Aus+21]. Used with kind permission of Emiel Hoogeboom.

model to map this to 256×256 or 1024×1024 . This approach is used in Google’s **Imagen** model [Sah+22b], which is a text-to-image system (see Figure 25.14). Imagen uses large pre-trained text encoder, based on T5-XXL [Raf+20a], combined with a VDM model (Section 25.2.4) based on the U-net architecture, to generate impressive-looking images (see Figure 20.2).

In addition to conditioning on text, it is possible to condition on another image to create a model for **image-to-image translation**. For example, we can learn map a gray-scale image \mathbf{c} to a color image \mathbf{x} , or a corrupted or occluded image \mathbf{c} to a clean version \mathbf{x} . This can be done by training a multi-task conditional diffusion model, as explained in [Sah+22a]. See Figure 20.4 for some sample outputs.

25.7 Diffusion for discrete state spaces

So far in this chapter, we have focused on Gaussian diffusion for generating real-valued data. However it is also possible to define diffusion models for discrete data, such as text or semantic segmentation labels, either by using a continuous latent embedding space (see Section 25.5.4), or by defining diffusion operations directly on the discrete state space, as we discuss below.

25.7.1 Discrete Denoising Diffusion Probabilistic Models

In this section we discuss the Discrete Denoising Diffusion Probabilistic Model (**D3PM**) of [Aus+21], which defines a discrete time diffusion process directly on the discrete state space. (This builds on prior work such as **multinomial diffusion** [Hoo+21], and the original diffusion paper of [SD+15b].)

The basic idea is illustrated in Figure 25.15 in the context of semantic segmentation, which associates a categorical label to each pixel in an image. On the right we illustrate some sample images, and the corresponding categorical distribution that they induce over a single pixel. We gradually transform these pixel-wise distributions to the uniform distribution, using a stochastic sampling process that we describe below. We then learn a neural network to invert this process, so it can generate discrete data from noise; in the diagram, this corresponds to moving from left to right.

To ensure efficient training, we require that we can efficiently sample from $q(\mathbf{x}_t|\mathbf{x}_0)$ for an arbitrary timestep t , so we can randomly sample time steps when optimizing the variational bound in Equation (25.27). In addition, we require that $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ have a tractable form, so we can efficiently compute the KL terms

$$L_{t-1}(\mathbf{x}_0) = \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)) \quad (25.67)$$

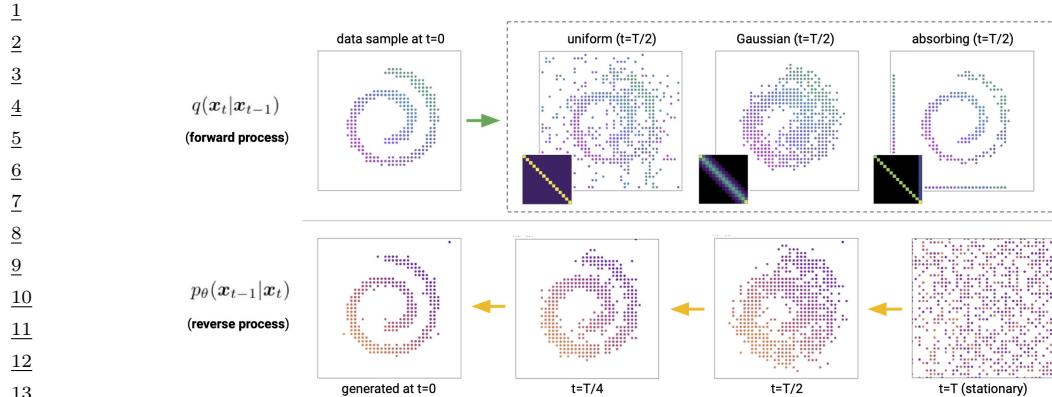


Figure 25.16: D3PM forward and (learned) reverse process applied to a quantized Swiss roll. Each dot represents a 2d categorical variable. Top: samples from the uniform, discretized Gaussian, and absorbing state models, along with corresponding transition matrices \mathbf{Q} . Bottom: samples from a learned discretized Gaussian reverse process. From Figure 1 of [Aus+21]. Used with kind permission of Jacob Austin.

Finally, it is useful if the forwards process converges to a known stationary distribution, $\pi(\mathbf{x}_T)$, which we can use for our generative prior $p(\mathbf{x}_T)$; this ensures $D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T)) = 0$.

To satisfy the above criteria, we assume the state consists of D independent blocks, each representing a categorical variable, $x_t \in \{1, \dots, K\}$; we represent this by the one-hot row vector \mathbf{x}_0 . In general, this will represent a vector of probabilities. We then define the forwards diffusion kernel as follows:

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \text{Cat}(\mathbf{x}_t|\mathbf{x}_{t-1}\mathbf{Q}_t) \quad (25.68)$$

where $[\mathbf{Q}_t]_{ij} = q(x_t = j|x_{t-1} = k)$ is a row stochastic transition matrix. (We discuss how to define \mathbf{Q}_t in Section 25.7.2.)

We can derive the t -step marginal of the forwards process as follows:

$$q(\mathbf{x}_t|\mathbf{x}_0) = \text{Cat}(\mathbf{x}_t|\mathbf{x}_0\overline{\mathbf{Q}}_t), \quad \overline{\mathbf{Q}}_t = \mathbf{Q}_1\mathbf{Q}_2 \cdots \mathbf{Q}_t \quad (25.69)$$

Similarly, we can reverse the forwards process as follows:

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)} = \text{Cat}\left(\mathbf{x}_{t-1} \mid \frac{\mathbf{x}_t\mathbf{Q}_t^\top \odot \mathbf{x}_0\overline{\mathbf{Q}}_{t-1}}{\mathbf{x}_0\overline{\mathbf{Q}}_t\mathbf{x}_t^\top}\right) \quad (25.70)$$

We discuss how to define the generative process $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$ in Section 25.7.3. Since both distributions factorize, we can easily compute the KL distributions in Equation (25.67) by summing the KL for each dimension.

43

25.7.2 Choice of Markov transition matrices for the forward processes

In this section, we give some examples of how to represent the transition matrix \mathbf{Q}_t .

One simple approach is to use $\mathbf{Q}_t = (1 - \beta_t)\mathbf{I} + \beta_t/K$, which we can write in scalar form as follows:

$$[\mathbf{Q}_t]_{ij} = \begin{cases} 1 - \frac{K-1}{K}\beta_t & \text{if } i = j \\ \frac{1}{K}\beta_t & \text{if } i \neq j \end{cases} \quad (25.71)$$

Intuitively, this adds a little amount of uniform noise over the K classes, and with a large probability, $1 - \beta_t$, we sample from \mathbf{x}_{t-1} . We call this the uniform kernel. Since this is a doubly stochastic matrix with strictly positive entries, the stationary distribution is uniform. See Figure 25.16 for an illustration.

In the case of the uniform kernel, one can show [Hoo+21] that the marginal distribution is given by

$$q(\mathbf{x}_t | \mathbf{x}_0) = \text{Cat}(\mathbf{x}_t | \bar{\alpha}_t \mathbf{x}_0 + (1 - \bar{\alpha}_t)/K) \quad (25.72)$$

where $\alpha_t = 1 - \beta_t$ and $\bar{\alpha}_t = \prod_{\tau=1}^t \alpha_\tau$. This is similar to the Gaussian case discussed in Section 25.2. Furthermore, we can derive the posterior distribution as follows:

$$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \text{Cat}(\mathbf{x}_{t-1} | \boldsymbol{\theta}_{\text{post}}(\mathbf{x}_t, \boldsymbol{\theta}_0)), \quad \boldsymbol{\theta}_{\text{post}}(\mathbf{x}_t, \boldsymbol{\theta}_0) = \tilde{\boldsymbol{\theta}} / \sum_{k=1}^K \tilde{\theta}_k \quad (25.73)$$

$$\tilde{\boldsymbol{\theta}} = [\alpha_t \mathbf{x}_t + (1 - \alpha_t)/K] \odot [\bar{\alpha}_{t-1} \mathbf{x}_0 + (1 - \bar{\alpha}_{t-1})/K] \quad (25.74)$$

Another option is to define a special **absorbing state** m , representing a MASK token, which we transition into with probability β_t . Formally, we have $\mathbf{Q}_t = (1 - \beta_t)\mathbf{I} + \beta_t \mathbf{1} \mathbf{e}_m^\top$, or, in scalar form,

$$[\mathbf{Q}_t]_{ij} = \begin{cases} 1 & \text{if } i = j = m \\ 1 - \beta_t & \text{if } i = j \neq m \\ \beta_t & \text{if } j = m, i \neq m \end{cases} \quad (25.75)$$

This converges to a point-mass distribution on state m . See Figure 25.16 for an illustration.

Another option, suitable for quantized ordinal values, is to use a **discretized Gaussian**, that transitions to other nearby states, with a probability that depends on how similar the states are in numerical value. If we ensure the transition matrix is doubly stochastic, the resulting stationary distribution will again be uniform. See Figure 25.16 for an illustration.

25.7.3 Parameterization of the reverse process

While it is possible to directly predict the logits $p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1} | \mathbf{x}_t)$ using a neural network $f_{\boldsymbol{\theta}}(\mathbf{x}_t)$, it is preferable to directly predict the logits of the output, using $\tilde{p}_{\boldsymbol{\theta}}(\tilde{\mathbf{x}}_0 | \mathbf{x}_t)$; we can then combine this with the analytical expression for $q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)$ to get

$$p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1} | \mathbf{x}_t) \propto \sum_{\tilde{\mathbf{x}}_0} q(\mathbf{x}_{t-1} | \mathbf{x}_t, \tilde{\mathbf{x}}_0) \tilde{p}_{\boldsymbol{\theta}}(\tilde{\mathbf{x}}_0 | \mathbf{x}_t) \quad (25.76)$$

(The sum over $\tilde{\mathbf{x}}_0$ takes $O(DK)$ time, if there are D dimensions, each with K values.) One advantage of this approach, compared to directly learning $p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1} | \mathbf{x}_t)$, is that the model will automatically

1 satisfy any sparsity constraints in \mathbf{Q}_t . In addition, we can perform inference with k steps at a time,
2 by predicting
3

$$\frac{4}{5} p_{\theta}(\mathbf{x}_{t-k} | \mathbf{x}_t) \propto \sum_{\tilde{\mathbf{x}}_0} q(\mathbf{x}_{t-k} | \mathbf{x}_t, \tilde{\mathbf{x}}_0) \tilde{p}_{\theta}(\tilde{\mathbf{x}}_0 | \mathbf{x}_t) \quad (25.77)$$

6

7 Note that, in the multi-step Gaussian case, we require more complex models to handle multimodality
8 (see Section 25.5.2). By contrast, discrete distributions already have this flexibility built in.
9

10 **25.7.4 Noise schedules**
11

12 In this section we discuss how to choose the noise schedule for β_t . For discretized Gaussian diffusion,
13 [Aus+21] propose to linearly increase the variance of the Gaussian noise before the discretization step.
14 For uniform diffusion, we can use a cosine schedule of the form $\alpha_t = \cos(\frac{t/T+s}{1+s} \frac{\pi}{2})$, where $s = 0.08$, as
15 proposed in [ND21]. (Recall that $\beta_t = 1 - \alpha_t$, so the noise increases over time.) For masked diffusion,
16 we can use a schedule of the form $\beta_t = 1/(T - t + 1)$, as proposed in [SD+15b].
17

18 **25.7.5 Connections to other probabilistic models for discrete sequences**
19

20 There are interesting connections between D3PM and other probabilistic text models. For example,
21 suppose we define the transition matrix as a combination of the unifrom transition matrix and an
22 absorbing MASK state, i.e., $\mathbf{Q} = \alpha \mathbf{1} \mathbf{e}_m^\top + \beta \mathbf{1} \mathbf{1}^\top / K + (1 - \alpha - \beta) \mathbf{I}$. For a one-step diffusion process
23 in which $q(\mathbf{x}_1 | \mathbf{x}_0)$ replaces $\alpha = 10\%$ of the tokens with MASK, and $\beta = 5\%$ uniformly at random,
24 we recover the same objective that is used to train the **BERT** language model, namely

$$\frac{26}{27} L_0(\mathbf{x}_0) = -\mathbb{E}_{q(\mathbf{x}_1 | \mathbf{x}_0)} \log p_{\theta}(\mathbf{x}_0 | \mathbf{x}_1) \quad (25.78)$$

28

29 (This follows since $L_T = 0$, and there are no other time steps used in the variational bound in
29 Equation (25.27).)

30 Now consider a diffusion process that deterministically masks tokens one-by-one. For a se-
31 quence of length $N = T$, we have $q([\mathbf{x}_t]_i | \mathbf{x}_0) = [\mathbf{x}_0]_i$ if $i < N - t$ (pass through), else $[\mathbf{x}_t]_i$
32 is set to MASK. Because this is a deterministic process, the posterior $q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)$ is a delta
33 function on the \mathbf{x}_t with one fewer mask tokens. One can then show that the KL term becomes
34 $D_{\text{KL}}(q([\mathbf{x}_t]_i | \mathbf{x}_t, \mathbf{x}_0) \| p_{\theta}([\mathbf{x}_{t-1}]_i | \mathbf{x}_t)) = -\log p_{\theta}([\mathbf{x}_0]_i | \mathbf{x}_t)$, which is the standard cross-entropy loss
35 for an autoregressive model.

36 Finally one can show that generative **masked language models**, such as [WC19; Gha+19], also
37 correspond to discrete diffusion processes: the sequence starts will all locations masked out, and
38 each step, a set of tokens are generated, given the previous sequence. The **MaskGIT** method of
39 [Cha+22] uses a similar procedure in the image domain, after applying vector quantization to image
40 patches. These parallel, iterative decoders are much faster than sequential autoregressive decoders.
41 See Figure 25.17 for an illustration.

42

43

44

45

46

47

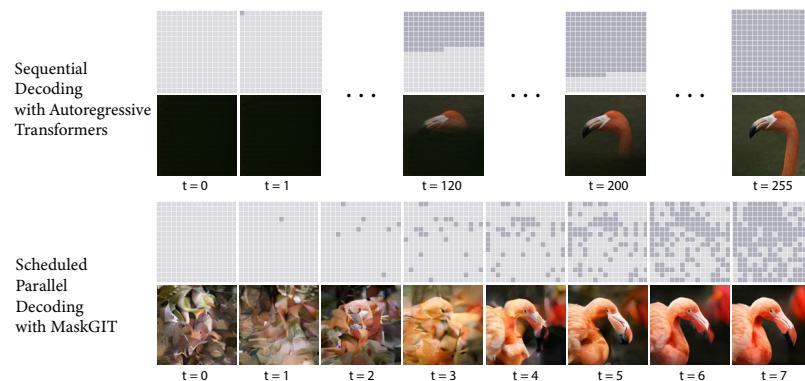


Figure 25.17: Comparison of sequential image generation with a transformer (top) vs parallel generation with MaskGIT (bottom). All pixels start out in the MASK state, denoted by light gray. In the transformer, we generate one pixel at a time, so it takes 256 steps for the whole image. In the MaskGIT method, multiple states are generated in parallel, which only takes 8 steps. From Figure 2 of [Cha+22]. Used with kind permission of Huiwen Chang.

26 Generative adversarial networks

This chapter is written by Mihaela Rosca, Shakir Mohamed, and Balaji Lakshminarayanan.

26.1 Introduction

In this chapter, we focus on **implicit generative models**, which are a kind of probabilistic model without an explicit likelihood function [ML16]. This includes the family of **generative adversarial networks** or **GANs** [Goo16]. In this chapter, we provide an introduction to this topic, focusing on a probabilistic perspective.

To develop a probabilistic formulation for GANs, it is useful to first distinguish between two types of probabilistic models: “**prescribed probabilistic models**” and “**implicit probabilistic models**” [DG84]. Prescribed probabilistic models, which we will call **explicit probabilistic models**, provide an explicit parametric specification of the distribution of an observed random variable \mathbf{x} , specifying a log-likelihood function $\log q_\theta(\mathbf{x})$ with parameters θ . Most models we encountered in this book thus far are of this form, whether they be state-of-the-art classifiers, large-vocabulary sequence models, or fine-grained spatio-temporal models. Alternatively, we can specify an **implicit probabilistic model** that defines a stochastic procedure to directly generate data. Such models are the natural approach for problems in climate and weather, population genetics, and ecology, since the mechanistic understanding of such systems can be used to directly describe the generative model. We illustrate the difference between implicit and explicit models in Figure 26.1.

The form of implicit generative models we focus on in this chapter can be expressed as a probabilistic latent variable model, similar to VAEs (Chapter 21). Implicit generative models use a latent variable \mathbf{z} and transform it using a deterministic function G_θ that maps from $\mathbb{R}^m \rightarrow \mathbb{R}^d$ using parameters θ . Implicit generative models do not include a likelihood function or observation model. Instead, the generating procedure defines a valid density on the output space that forms an effective likelihood function:

$$\mathbf{x} = G_\theta(\mathbf{z}'); \quad \mathbf{z}' \sim q(\mathbf{z}) \tag{26.1}$$

$$q_\theta(\mathbf{x}) = \frac{\partial}{\partial x_1} \cdots \frac{\partial}{\partial x_d} \int_{\{G_\theta(\mathbf{z}) \leq \mathbf{x}\}} q(\mathbf{z}) d\mathbf{z}, \tag{26.2}$$

where $q(\mathbf{z})$ is a distribution over latent variables that provides the external source of randomness. Equation (26.2) is the definition of the transformed density $q_\theta(\mathbf{x})$ defined as the derivative of a cumulative distribution function, and hence integrates the distribution $q(\mathbf{z})$ over all events defined

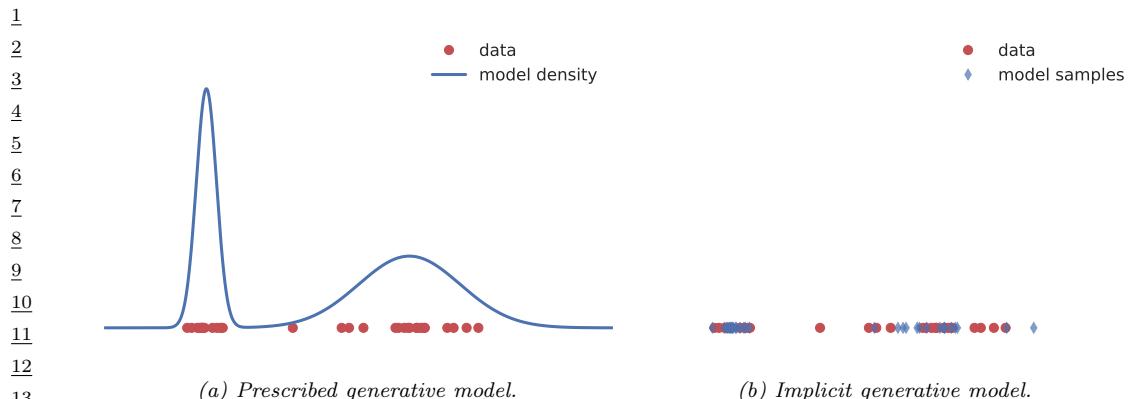


Figure 26.1: Visualizing the difference between prescribed and implicit generative models. Prescribed models provide direct access to the learned density (sometimes unnormalized). Implicit models only provide access to a simulator which can be used to generate samples from an implied density. Generated by genmo_types_implicit_explicit.ipynb

by the set $\{G_{\theta}(z) \leq x\}$. When the latent and data dimension are equal ($m = d$) and the function $G_{\theta}(z)$ is invertible or has easily characterized roots, we recover the rule for transformations of probability distributions. This transformation of variables property is also used in normalizing flows (Chapter 23). In diffusion models (Chapter 25), we also transform noise into data and vice versa, but the transformation is not strictly invertible.

We can develop more general and flexible implicit generative models where the function G is a non-linear function with $d > m$, e.g., specified by a deep network. Such models are sometimes called **generator networks** or **generative neural samplers**; they can also be thought of as **differentiable simulators**. Unfortunately the integral (26.2) is intractable in these kinds of models, and we may not even be able to determine the set $\{G_{\theta}(z) \leq x\}$. Of course, intractability is also a challenge for explicit latent variable models such as VAEs (Chapter 21), but in the GAN case, the lack of a likelihood term makes the learning problem even harder. Therefore this problem is called **likelihood-free inference** or **simulation-based inference**.

Likelihood-free inference also forms the basis of the field known as **approximate Bayesian computation** or **ABC**, which we briefly discuss in Section 13.6.5. ABC and GANs give us two different algorithmic frameworks for learning in implicit generative models. Both approaches rely on a learning principle based on *comparing real and simulated data*. This type of learning by comparison instantiates a core principle of likelihood-free inference, and expanding on this idea is the focus of the next section. The subsequent sections will then focus on GANs specifically, to develop a more detailed foundation and practical considerations. (See also <https://poloclub.github.io/ganlab/> for an interactive tutorial.)

26.2 Learning by comparison

In most of this book, we rely on the principle of maximum likelihood for learning. By maximizing the likelihood we effectively minimize the KL divergence between the model q_{θ} (with parameters θ)

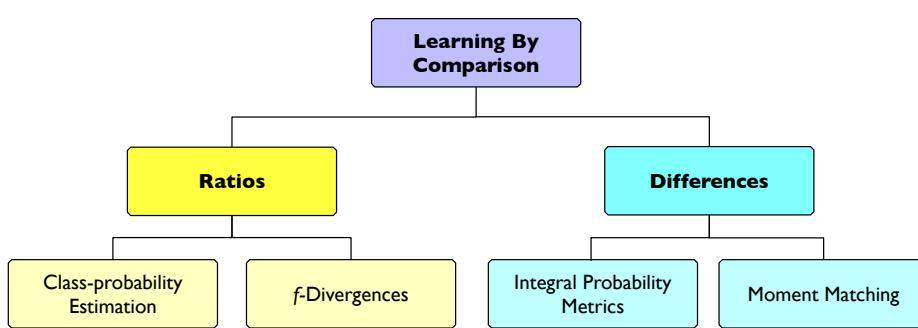


Figure 26.2: Overview of approaches for learning in implicit generative models

and the unknown true data distribution p^* . Recalling equation (26.2), in implicit models we cannot evaluate $q_\theta(\mathbf{x})$, and thus cannot use maximum likelihood training. As implicit models provide a sampling procedure, we instead are searching for learning principles that only use *samples from the model*.

The task of learning in implicit models is to determine, from two sets of samples, whether their distributions are close to each other and to quantify the distance between them. We can think of this as a ‘two sample’ or likelihood-free approach to learning by comparison. There are many ways of doing this, including using distributional divergences or distances through binary classification, the method of moments, and other approaches. Figure 26.2 shows an overview of different approaches for learning by comparison.

26.2.1 Guiding principles

We are looking for objectives $\mathcal{D}(p^*, q)$ that satisfy the following requirements:

1. Provide guarantees about learning the data distribution: $\operatorname{argmin}_q \mathcal{D}(p^*, q) = p^*$.
2. Can be evaluated only using samples from the data and model distribution.
3. Are computationally cheap to evaluate.

Many distributional distances and divergences satisfy the first requirement, since by definition they satisfy the following:

$$\mathcal{D}(p^*, q) \geq 0; \quad \mathcal{D}(p^*, q) = 0 \iff p^* = q \quad (26.3)$$

Many distributional distances and divergences, however, fail to satisfy the other two requirements: they cannot be evaluated only using samples — such as the KL divergence, or are computationally intractable — such as the Wasserstein distance. The main approach to overcome these challenges is to *approximate the desired quantity through optimization* by introducing a comparison model, often called a **discriminator** or a **critic** D , such that:

$$\mathcal{D}(p^*, q) = \operatorname{argmax}_D \mathcal{F}(D, p^*, q) \quad (26.4)$$

where \mathcal{F} is a functional that depends on p^* and q only through samples. For the cases we discuss, both the model and the critic are parametric with parameters θ and ϕ respectively; instead of optimizing over distributions or functions, we optimize with respect to parameters. For the critic, this results in the optimization problem $\operatorname{argmax}_\phi \mathcal{F}(D_\phi, p^*, q_\theta)$. For the model parameters θ , the exact objective $\mathcal{D}(p^*, q_\theta)$ is replaced with the tractable approximation provided through the use of D_ϕ .

A convenient approach to ensure that $\mathcal{F}(D_\phi, p^*, q_\theta)$ can be estimated using only samples from the model and the unknown data distribution is to depend on the two distributions only in expectation:

$$\mathcal{F}(D_\phi, p^*, q_\theta) = \mathbb{E}_{p^*(\mathbf{x})} f(\mathbf{x}, \phi) + \mathbb{E}_{q_\theta(\mathbf{x})} g(\mathbf{x}, \phi) \quad (26.5)$$

where f and g are real valued functions whose choice will define \mathcal{F} . In the case of implicit generative models, this can be rewritten to use the sampling path $\mathbf{x} = G_\theta(\mathbf{z})$, $\mathbf{z} \sim q(\mathbf{z})$:

$$\mathcal{F}(D_\phi, p^*, q_\theta) = \mathbb{E}_{p^*(\mathbf{x})} f(\mathbf{x}, \phi) + \mathbb{E}_{q(\mathbf{z})} g(G_\theta(\mathbf{z}), \phi) \quad (26.6)$$

which can be estimated using Monte Carlo estimation

$$\mathcal{F}(D_\phi, p^*, q_\theta) \approx \frac{1}{N} \sum_{i=1}^N f(\hat{\mathbf{x}}_i, \phi) + \frac{1}{M} \sum_{i=1}^M g(G_\theta(\hat{\mathbf{z}}_i), \phi); \quad \hat{\mathbf{x}}_i \sim p^*(\mathbf{x}); \quad \hat{\mathbf{z}}_i \sim q(\mathbf{z}) \quad (26.7)$$

Next, we will see how to instantiate these guiding principles in order to find the functions f and g and thus the objective \mathcal{F} which can be used to train implicit models: class probability estimation (Section 26.2.2), bounds on f -divergences (Section 26.2.3), integral probability metrics (Section 26.2.4), and moment matching (Section 26.2.5).

26.2.2 Density ratio estimation using binary classifiers

One way to compare two distributions p^* and q_θ is to compute their density ratio $r(\mathbf{x}) = \frac{p^*(\mathbf{x})}{q_\theta(\mathbf{x})}$. The distributions are the same if and only if the ratio is 1 everywhere in the support of q_θ . Since we cannot evaluate the densities of implicit models, we must instead develop techniques to compute the density ratio from samples alone, following the guiding principles established above.

Fortunately, we can use the trick from Section 2.7.5 which converts density estimation into a binary classification problem to write

$$\frac{p^*(\mathbf{x})}{q_\theta(\mathbf{x})} = \frac{D(\mathbf{x})}{1 - D(\mathbf{x})} \quad (26.8)$$

where $D(\mathbf{x})$ is the discriminator or critic which is trained to distinguish samples coming from p^* vs q_θ .

For parametric classification, we can learn discriminators $D_\phi(\mathbf{x}) \in [0, 1]$ with parameters ϕ . Using knowledge and insight about probabilistic classification, we can learn the parameters by minimizing any proper scoring rule [GR07b] (see also Section 14.2.1). For the familiar Bernoulli log-loss (or binary cross entropy loss), we obtain the objective:

$$\begin{aligned} V(q_\theta, p^*) &= \arg \max_{\phi} \mathbb{E}_{p(\mathbf{x}|y)p(y)} [y \log D_\phi(\mathbf{x}) + (1 - y) \log(1 - D_\phi(\mathbf{x}))] \\ &= \arg \max_{\phi} \mathbb{E}_{p(\mathbf{x}|y=1)p(y=1)} \log D_\phi(\mathbf{x}) + \mathbb{E}_{p(\mathbf{x}|y=0)p(y=0)} \log(1 - D_\phi(\mathbf{x})) \end{aligned} \quad (26.9)$$

$$= \arg \max_{\phi} \frac{1}{2} \mathbb{E}_{p^*(\mathbf{x})} \log D_\phi(\mathbf{x}) + \frac{1}{2} \mathbb{E}_{q_\theta(\mathbf{x})} \log(1 - D_\phi(\mathbf{x})). \quad (26.10)$$

Loss	Objective Function ($D := D(\mathbf{x}; \phi) \in [0, 1]$)
Bernoulli loss	$\mathbb{E}_{p^*(\mathbf{x})}[\log D] + \mathbb{E}_{q_\theta(\mathbf{x})}[\log(1 - D)]$
Brier score	$\mathbb{E}_{p^*(\mathbf{x})}[-(1 - D)^2] + \mathbb{E}_{q_\theta(\mathbf{x})}[-D^2]$
Exponential loss	$\mathbb{E}_{p^*(\mathbf{x})}\left[\left(-\frac{1-D}{D}\right)^{\frac{1}{2}}\right] + \mathbb{E}_{q_\theta(\mathbf{x})}\left[\left(-\frac{D}{1-D}\right)^{\frac{1}{2}}\right]$
Misclassification	$\mathbb{E}_{p^*(\mathbf{x})}[-\mathbb{I}[D \leq 0.5]] + \mathbb{E}_{q_\theta(\mathbf{x})}[-\mathbb{I}[D > 0.5]]$
Hinge loss	$\mathbb{E}_{p^*(\mathbf{x})}\left[-\max\left(0, 1 - \log\frac{D}{1-D}\right)\right] + \mathbb{E}_{q_\theta(\mathbf{x})}\left[-\max\left(0, 1 + \log\frac{D}{1-D}\right)\right]$
Spherical	$\mathbb{E}_{p^*(\mathbf{x})}[\alpha D] + \mathbb{E}_{q_\theta(\mathbf{x})}[\alpha(1 - D)]; \quad \alpha = (1 - 2D + 2D^2)^{-\frac{1}{2}}$

Table 26.1: Proper scoring rules that can be maximized in class probability-based learning of implicit generative models. Based on [ML16].

The same procedure can be extended beyond the Bernoulli log-loss to other proper scoring rules used for binary classification, such as those presented in Table 26.1, adapted from [ML16]. The optimal discriminator D is $\frac{p^*(\mathbf{x})}{p^*(\mathbf{x}) + q_\theta(\mathbf{x})}$, since:

$$\frac{p^*(\mathbf{x})}{q_\theta(\mathbf{x})} = \frac{D^*(\mathbf{x})}{1 - D^*(\mathbf{x})} \implies D^*(\mathbf{x}) = \frac{p^*(\mathbf{x})}{p^*(\mathbf{x}) + q_\theta(\mathbf{x})} \quad (26.11)$$

By substituting the optimal discriminator into the scoring rule (26.10), we can show that the objective V can also be interpreted as the minimization of the Jensen-Shannon divergence.

$$V^*(q_\theta, p^*) = \frac{1}{2}\mathbb{E}_{p^*(\mathbf{x})}[\log \frac{p^*(\mathbf{x})}{p^*(\mathbf{x}) + q_\theta(\mathbf{x})}] + \frac{1}{2}\mathbb{E}_{q_\theta(\mathbf{x})}[\log(1 - \frac{p^*(\mathbf{x})}{p^*(\mathbf{x}) + q_\theta(\mathbf{x})})] \quad (26.12)$$

$$= \frac{1}{2}\mathbb{E}_{p^*(\mathbf{x})}[\log \frac{p^*(\mathbf{x})}{\frac{p^*(\mathbf{x}) + q_\theta(\mathbf{x})}{2}}] + \frac{1}{2}\mathbb{E}_{q_\theta(\mathbf{x})}[\log(\frac{q_\theta(\mathbf{x})}{\frac{p^*(\mathbf{x}) + q_\theta(\mathbf{x})}{2}})] - \log 2 \quad (26.13)$$

$$= \frac{1}{2}D_{\text{KL}}\left(p^* \parallel \frac{p^* + q_\theta}{2}\right) + \frac{1}{2}D_{\text{KL}}\left(q_\theta \parallel \frac{p^* + q_\theta}{2}\right) - \log 2 \quad (26.14)$$

$$= JSD(p^*, q_\theta) - \log 2 \quad (26.15)$$

where JSD denotes the Jensen-Shannon divergence:

$$JSD(p^*, q_\theta) = \frac{1}{2}D_{\text{KL}}\left(p^* \parallel \frac{p^* + q_\theta}{2}\right) + \frac{1}{2}D_{\text{KL}}\left(q_\theta \parallel \frac{p^* + q_\theta}{2}\right) \quad (26.16)$$

This establishes a connection between *optimal* binary classification and distributional divergences. By using binary classification, we were able to compute the distributional divergence using only samples, which is the important property needed for learning implicit generative models; as expressed in the guiding principles (Section 26.2.1), we have turned an intractable estimation problem — how to estimate the JSD divergence, into an optimization problem — how to learn a classifier which can be used to approximate that divergence.

We would like to train the parameters θ of generative model to minimize the divergence:

$$\min_{\theta} JSD(p^*, q_\theta) = \min_{\theta} V^*(q_\theta, p^*) + \log 2 \quad (26.17)$$

$$= \min_{\theta} \frac{1}{2}\mathbb{E}_{p^*(\mathbf{x})} \log D^*(\mathbf{x}) + \frac{1}{2}\mathbb{E}_{q_\theta(\mathbf{x})} \log(1 - D^*(\mathbf{x})) + \log 2 \quad (26.18)$$

1 Since we do not have access to the optimal classifier D^* but only to the neural approximation D_ϕ
 2 obtained using the optimization in (26.10) , this results in a min-max optimization problem:
 3

$$4 \quad \min_{\theta} \max_{\phi} \frac{1}{2} \mathbb{E}_{p^*(\mathbf{x})} [\log D_\phi(\mathbf{x})] + \frac{1}{2} \mathbb{E}_{q_\theta(\mathbf{x})} [\log(1 - D_\phi(\mathbf{x}))] \quad (26.19)$$

5 By replacing the generating procedure (26.1) in (26.19) we obtain the objective in terms of the
 6 latent variables \mathbf{z} of the implicit generative model:
 7

$$8 \quad \min_{\theta} \max_{\phi} \frac{1}{2} \mathbb{E}_{p^*(\mathbf{x})} [\log D_\phi(\mathbf{x})] + \frac{1}{2} \mathbb{E}_{q(\mathbf{z})} [\log(1 - D_\phi(G_\theta(\mathbf{z})))], \quad (26.20)$$

9 which recovers the definition proposed in the original GAN paper [Goo+14]. The core principle
 10 behind GANs is to train a discriminator, in this case a binary classifier, to approximate a distance
 11 or divergence between the model and data distributions, and to then train the generative model to
 12 minimize this approximation of the divergence or distance.

13 Beyond the use of the Bernoulli scoring rule used above, other scoring rules have been used to
 14 train generative models via min-max optimization. The Brier scoring rule, which under discriminator
 15 optimality conditions can be shown to correspond to minimizing the Pearson χ^2 divergence via
 16 similar arguments as the ones shown above has lead to LS-GAN [Mao+17]. The hinge scoring rule
 17 has become popular [Miy+18b; BDS18], and under discriminator optimality conditions corresponds
 18 to minimizing the total variational distance [NWJ+09].

19 The connection between proper scoring rules and distributional divergences allows the construction
 20 of convergence guarantees for the learning criteria above, under infinite capacity of the discriminator
 21 and generator: since the minimizer of distributional divergence is the true data distribution (Equa-
 22 tion 26.3), if the discriminator is optimal and the generator has enough capacity, it will learn the
 23 data distribution. In practice however, this assumption will not hold, as discriminators are rarely
 24 optimal; we will discuss this at length in Section 26.3.
 25

26.2.3 Bounds on f -divergences

26 As we saw with the appearance of the Jensen-Shannon divergence in the previous section, we can
 27 consider directly using a measure of distributional divergence to derive methods for learning in
 28 implicit models. One general class of divergences are the f -divergences (Section 2.7.1) defined as:
 29

$$30 \quad \mathcal{D}_f [p^*(\mathbf{x}) \| q_\theta(\mathbf{x})] = \int q_\theta(\mathbf{x}) f \left(\frac{p^*(\mathbf{x})}{q_\theta(\mathbf{x})} \right) d\mathbf{x} \quad (26.21)$$

31 where f is a convex function such that $f(1) = 0$. For different choices of f , we can recover known
 32 distributional divergences such as the KL, reverse KL, and Jensen-Shannon divergence. We discuss
 33 such connections in Section 2.7.1, and provide a summary in Table 26.2.
 34

35 To evaluate Equation (26.21) we will need to evaluate the density of the data $p^*(\mathbf{x})$ and the model
 36 $q_\theta(\mathbf{x})$, neither of which are available. In the previous section we overcame the challenge of evaluating
 37 the density ratio by transforming it into a problem of binary classification. In this section, we will
 38 instead look towards the role of lower bounds on f -divergences, which is an approach for tractability
 39 that is also used for variational inference (Chapter 10).
 40

41

Divergence	f	f^\dagger	Optimal Critic
KL	$u \log u$	e^{u-1}	$1 + \log r(\mathbf{x})$
Reverse KL	$-\log u$	$-1 - \log(-u)$	$-1/r(\mathbf{x})$
JSD	$u \log u - (u+1) \log \frac{u+1}{2}$	$-\log(2-e^u)$	$\frac{2}{1+1/r(\mathbf{x})}$
Pearson χ^2	$(u-1)^2$	$\frac{1}{4}u^2 + u$	$\left(\sqrt{r(\mathbf{x})} - 1\right) \sqrt{1/r(\mathbf{x})}$

Table 26.2: Standard divergences as f divergences for various choices of f . The optimal critic is written as a function of the density ratio $r(\mathbf{x}) = \frac{p^*(\mathbf{x})}{q_\theta(\mathbf{x})}$.

f -divergences have a widely-developed theory in convex analysis and information theory. Since the function f in Equation (26.21) is convex, we know that we can find a tangent that bounds it from below. The variational formulation of the f -divergence is [NWJ10b; NCT16c]:

$$\mathcal{D}_f[p^*(\mathbf{x}) \| q_\theta(\mathbf{x})] = \int q_\theta(\mathbf{x}) f\left(\frac{p^*(\mathbf{x})}{q_\theta(\mathbf{x})}\right) d\mathbf{x} \quad (26.22)$$

$$= \int q_\theta(\mathbf{x}) \sup_{t: \mathcal{X} \rightarrow \mathbb{R}} \left[t(\mathbf{x}) \frac{p^*(\mathbf{x})}{q_\theta(\mathbf{x})} - f^\dagger(t(\mathbf{x})) \right] d\mathbf{x} \quad (26.23)$$

$$= \int \sup_{t: \mathcal{X} \rightarrow \mathbb{R}} p^*(\mathbf{x}) t(\mathbf{x}) - q_\theta(\mathbf{x}) f^\dagger(t(\mathbf{x})) d\mathbf{x} \quad (26.24)$$

$$\geq \sup_{t \in \mathcal{T}} \mathbb{E}_{p^*(\mathbf{x})}[t(\mathbf{x})] - \mathbb{E}_{q_\theta(\mathbf{x})}[f^\dagger(t(\mathbf{x}))]. \quad (26.25)$$

In the second line we use the result from convex analysis, discussed Supplementary Section 6.3, that re-expresses the convex function f using $f(u) = \sup_t ut - f^\dagger(t)$, where f^\dagger is the convex conjugate of the function f , and t is a parameter we optimize over. Since we apply f at $u = \frac{p^*(\mathbf{x})}{q_\theta(\mathbf{x})}$ for all $\mathbf{x} \in \mathcal{X}$, we make the parameter t be a function $t(\mathbf{x})$. The final inequality comes from replacing the supremum over all functions from the data domain \mathcal{X} to \mathbb{R} with the supremum over a family of functions \mathcal{T} (such as the family of functions expressible by a neural network architecture), which might not be able to capture the true supremum. The function t takes the role of the discriminator or critic.

The final expression in Equation (26.25) follows the general desired form of Equation 26.5: it is the difference of two expectations, and these expectations can be computed by Monte Carlo estimation using only samples, as in Equation (26.7); despite starting with an objective (Equation 26.21) which contravened the desired principles for training implicit generative models, variational bounds have allowed us to construct an approximation which satisfies all desiderata.

Using bounds on the f -divergence, we obtain an objective (26.25) that allows learning both the generator and critic parameters. We use a critic D with parameters ϕ to estimate the bound, and then optimize the parameters θ of the generator to minimize the approximation of the f -divergence provided by the critic (we replace t above with D_ϕ , to retain standard GAN notation):

$$\min_{\theta} \mathcal{D}_f(p^*, q_\theta) \geq \min_{\theta} \max_{\phi} \mathbb{E}_{p^*(\mathbf{x})}[D_\phi(\mathbf{x})] - \mathbb{E}_{q_\theta(\mathbf{x})}[f^\dagger(D_\phi(\mathbf{x}))] \quad (26.26)$$

$$= \min_{\theta} \max_{\phi} \mathbb{E}_{p^*(\mathbf{x})}[D_\phi(\mathbf{x})] - \mathbb{E}_{q(\mathbf{z})}[f^\dagger(D_\phi(G_\theta(\mathbf{z})))] \quad (26.27)$$

This approach to train an implicit generative model leads to f -GANs [NCT16c]. It is worth noting that there exists an equivalence between the scoring rules in the previous section and bounds on

¹ ² ³ ⁴ ⁵ ⁶ ⁷ ⁸ ⁹ ¹⁰ ¹¹ ¹² ¹³ ¹⁴ ¹⁵ ¹⁶ ¹⁷ ¹⁸ ¹⁹ ²⁰ ²¹ ²² ²³ ²⁴ ²⁵ ²⁶ ²⁷ ²⁸ ²⁹ ³⁰ ³¹ ³² ³³ ³⁴ ³⁵ ³⁶ ³⁷ ³⁸ ³⁹ ⁴⁰ ⁴¹ ⁴² ⁴³ ⁴⁴ ⁴⁵ ⁴⁶ ⁴⁷

f-divergences [RW11]: for each scoring rule we can find an *f*-divergence that leads to the same training criteria and the same min-max game of Equation 26.27. An intuitive way to grasp the connection between *f*-divergences and proper scoring rules is through their use of density ratios: in both cases the optimal critic approximates a quantity directly related to the density ratio (see Table 26.2 for *f*-divergences and Equation (26.11) for scoring rules).

26.2.4 Integral probability metrics

Instead of comparing distributions by using their ratio as we did in the previous two sections, we can instead study their difference. A general class of measure of difference is given by the Integral Probability Metrics (Section 2.7.2) defined as:

$$I_{\mathcal{F}}(p^*(\mathbf{x}), q_\theta(\mathbf{x})) = \sup_{f \in \mathcal{F}} |\mathbb{E}_{p^*(\mathbf{x})} f(\mathbf{x}) - \mathbb{E}_{q_\theta(\mathbf{x})} f(\mathbf{x})|. \quad (26.28)$$

The function f is a test or witness function that will take the role of the discriminator or critic. To use IPMs we must define the class of real valued, measurable functions \mathcal{F} over which the supremum is taken, and this choice will lead to different distances, just as choosing different convex functions f leads to different *f*-divergences. Integral probability metrics are distributional distances: beyond satisfying the conditions for distributional divergences $\mathcal{D}(p^*, q) \geq 0$; $\mathcal{D}(p^*, q) = 0 \iff p^* = q$ (Equation (26.3)), they are also symmetric $\mathcal{D}(p, q) = \mathcal{D}(q, p)$ and satisfy the triangle inequality $\mathcal{D}(p, q) \leq \mathcal{D}(p, r) + \mathcal{D}(r, q)$.

Not all function families satisfy these conditions of create a valid distance $I_{\mathcal{F}}$. To see why consider the case where $\mathcal{F} = \{z\}$ where z is the function $z(\mathbf{x}) = 0$. This choice of \mathcal{F} entails that regardless of the two distributions chosen, the value in Equation 26.28 would be 0, violating the requirement that distance between two distributions be 0 only if the two distributions are the same. A popular choice of \mathcal{F} for which $I_{\mathcal{F}}$ satisfies the conditions of a valid distributional distance is the set of 1-Lipschitz functions, which leads to the Wasserstein distance [Vil08]:

$$W_1(p^*(\mathbf{x}), q_\theta(\mathbf{x})) = \sup_{f: \|f\|_{\text{Lip}} \leq 1} \mathbb{E}_{p^*(\mathbf{x})} f(\mathbf{x}) - \mathbb{E}_{q_\theta(\mathbf{x})} f(\mathbf{x}) \quad (26.29)$$

We show an example of a Wasserstein critic in Figure 26.3a. The supremum over the set of 1-Lipschitz functions is intractable for most cases, which again suggests the introduction of a learned critic:

$$W_1(p^*(\mathbf{x}), q_\theta(\mathbf{x})) = \sup_{f: \|f\|_{\text{Lip}} \leq 1} \mathbb{E}_{p^*(\mathbf{x})} f(\mathbf{x}) - \mathbb{E}_{q_\theta(\mathbf{x})} f(\mathbf{x}) \quad (26.30)$$

$$\geq \max_{\phi: \|D_\phi\|_{\text{Lip}} \leq 1} \mathbb{E}_{p^*(\mathbf{x})} D_\phi(\mathbf{x}) - \mathbb{E}_{q_\theta(\mathbf{x})} D_\phi(\mathbf{x}), \quad (26.31)$$

where the critic D_ϕ has to be regularized to be 1-Lipschitz (various techniques for Lipschitz regularization via gradient penalties or spectral normalization methods have been used [ACB17; Gul+17]). As was the case with *f*-divergences, we replace an intractable quantity which requires a supremum over a class of functions with a bound obtained using a subset of this function class, a subset which can be modeled using neural networks.

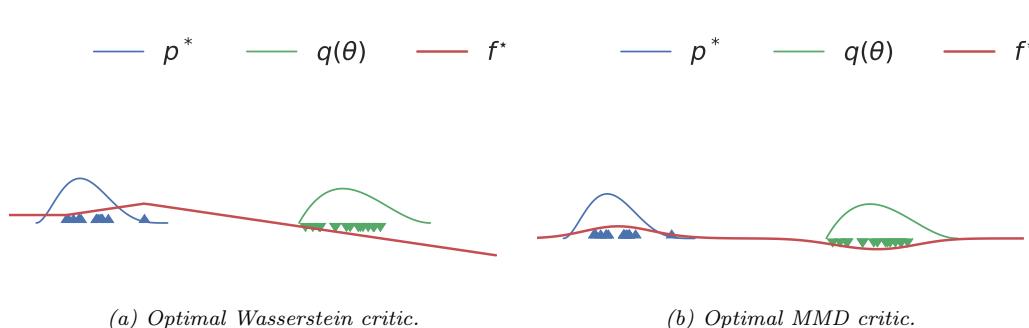


Figure 26.3: Optimal critics in Integral Probability Metrics (IPMs). Generated by `ipm_divergences.ipynb`

To train a generative model, we again introduce a min max game:

$$\min_{\theta} W_1(p^*(x), q_\theta(x)) \geq \min_{\theta} \max_{\phi: \|D_\phi\| \leq 1} \mathbb{E}_{p^*(x)} D_\phi(x) - \mathbb{E}_{q_\theta(x)} D_\phi(x) \quad (26.32)$$

$$= \min_{\boldsymbol{\theta}} \max_{\phi: \|D_\phi\|_{\text{op}} \leq 1} \mathbb{E}_{p^*(\mathbf{x})} D_\phi(\mathbf{x}) - \mathbb{E}_q(\mathbf{z}) D_\phi(G_{\boldsymbol{\theta}}(\mathbf{z})) \quad (26.33)$$

This leads to the popular WassersteinGAN [ACB17].

If we replace the choice of function family \mathcal{F} to that of functions in an RKHS (Section 18.3.7.1) with norm one, we obtain the **maximum mean discrepancy (MMD)** discussed in Section 2.7.3:

$$\text{MMD}(p^*(\mathbf{x}), q_\theta(\mathbf{x})) = \sup_{f: \|f\|_1, \dots, \|f\|_\infty = 1} \mathbb{E}_{p^*(\mathbf{x})} f(\mathbf{x}) - \mathbb{E}_{q_\theta(\mathbf{x})} f(\mathbf{x}). \quad (26.34)$$

We show an example of an MMD critic in Figure 26.3b. It is often more convenient to use the square MMD loss [LSZ15; DRG15], which can be evaluated using the kernel \mathcal{K} (Section 18.3.7.1):

$$\text{MMD}^2(p^*, q_0) = \mathbb{E}_{\pi \sim p^*} \mathbb{E}_{\pi' \sim p^*} \mathcal{K}(x, x') - 2\mathbb{E}_{\pi \sim p^*} \mathbb{E}_{\pi' \sim q_0} \mathcal{K}(x, u) + \mathbb{E}_{\pi \sim p^*} \mathbb{E}_{\pi' \sim q_0} \mathcal{K}(u, u') \quad (26.35)$$

$$= \mathbb{E}_{p^*(\mathbf{x})} \mathbb{E}_{p^*(\mathbf{x}')} \mathcal{K}(\mathbf{x}, \mathbf{x}') - 2 \mathbb{E}_{p^*(\mathbf{x})} \mathbb{E}_{q(\mathbf{z})} \mathcal{K}(\mathbf{x}, G_{\theta}(\mathbf{z})) + \mathbb{E}_{q(\mathbf{z})} \mathbb{E}_{q(\mathbf{z}')} \mathcal{K}(G_{\theta}(\mathbf{z}), G_{\theta}(\mathbf{z}')) \quad (26.36)$$

The MMD can be directly used to learn a generative model, often called a generative matching network [LSZ15].

$$\min_{\theta} \text{MMD}^2(p^*, q_{\theta}) \quad (26.37)$$

The choice of kernel is important. Using a fixed or predefined kernel such as a radial basis function (RBF) kernel might not be appropriate for all data modalities, such as high dimensional images. Thus we are looking for a way to learn a feature function ζ such that $\mathcal{K}(\zeta(\mathbf{x}), \zeta(\mathbf{x}'))$ is a valid kernel; luckily, we can use that for any characteristic kernel $\mathcal{K}(\mathbf{x}, \mathbf{x}')$ and injective function ζ , $\mathcal{K}(\zeta(\mathbf{x}), \zeta(\mathbf{x}'))$ is also a characteristic kernel. While this tells us that we can use feature functions in the MMD objective, it does not tell us how to learn the features. In order to ensure that the learned features are sensitive to differences between the data distribution $p^*(\mathbf{x})$ and the model distribution $q_\theta(\mathbf{x})$, the

1 kernel parameters are trained to *maximize* the square MMD. This again casts the problem into a
2 familiar min max objective by learning the projection ζ with parameters ϕ [Li+17b]:
3

$$\min_{\theta} \text{MMD}_{\zeta}^2(p_D, q_{\theta}) \quad (26.38)$$

$$\begin{aligned} &= \min_{\theta} \max_{\phi} \mathbb{E}_{p^*(\mathbf{x})} \mathbb{E}_{p^*(\mathbf{x}')} \mathcal{K}(\zeta_{\phi}(\mathbf{x}), \zeta_{\phi}(\mathbf{x}')) \\ &\quad - 2 \mathbb{E}_{p^*(\mathbf{x})} \mathbb{E}_{q_{\theta}(\mathbf{y})} \mathcal{K}(\zeta_{\phi}(\mathbf{x}), \zeta_{\phi}(\mathbf{y})) \\ &\quad + \mathbb{E}_{q_{\theta}(\mathbf{y})} \mathbb{E}_{q_{\theta}(\mathbf{y}')} \mathcal{K}(\zeta_{\phi}(\mathbf{y}), \zeta_{\phi}(\mathbf{y}')) \end{aligned} \quad (26.39)$$

11 where ζ_{ϕ} is regularized to be injective, though this is sometimes relaxed [Bin+18]. Unlike the
12 Wasserstein distance and f -divergences, Equation (26.39) can be estimated using Monte Carlo
13 estimation, without requiring a lower bound on the original objective.
14

15 26.2.5 Moment matching

16 More broadly than distances defined by integral probability metrics, for a set of test statistics s , one
17 can define a **moment matching** criteria [Pea36], also known as the method of moments:
18

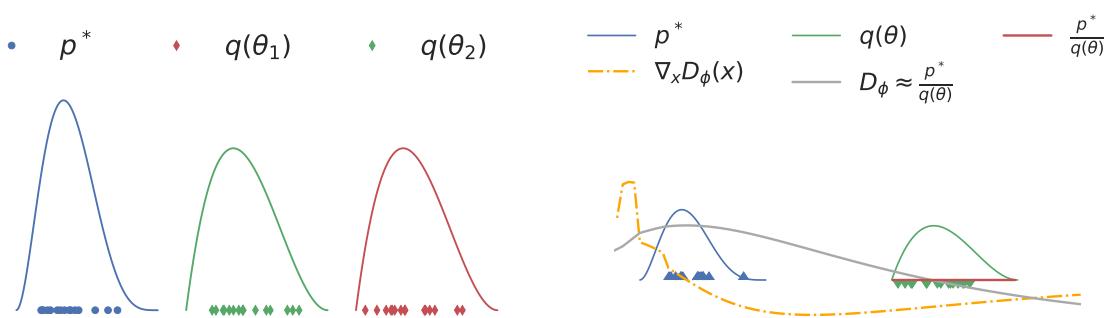
$$\min_{\theta} \left\| \mathbb{E}_{p^*(\mathbf{x})} s(\mathbf{x}) - \mathbb{E}_{q_{\theta}(\mathbf{x})} s(\mathbf{x}) \right\|_2^2 \quad (26.40)$$

19 where $m(\theta) = \mathbb{E}_{q_{\theta}(\mathbf{x})} s(\mathbf{x})$ is the *moment function*. The choice of statistic $s(\mathbf{x})$ is crucial, since as with
20 distributional divergences and distances, we would like to ensure that if the objective is minimized
21 and reaches the minimal value 0, the two distributions are the same $p^*(\mathbf{x}) = q_{\theta}(\mathbf{x})$. To see that not
22 all functions s satisfy this requirement consider the function $s(\mathbf{x}) = \mathbf{x}$: simply matching the means of
23 two distributions is not sufficient to match higher moments (such as variance). For likelihood based
24 models the score function $s(\mathbf{x}) = \log q_{\theta}(\mathbf{x})$ satisfies the above requirement and leads to a consistent
25 estimator [Vaa00], but this choice of s is not available for implicit generative models.
26

27 This motivates the search for other approaches of integrating the method of moments for implicit
28 models. The MMD can be seen as a moment matching criteria, by matching the means of the two
29 distributions after lifting the data into the feature space of an RKHS. But moment matching can go
30 beyond integral probability metrics: Ravuri et al. [Rav+18] show that one can *learn* useful moments by
31 using s as the set of features containing the gradients of a trained discriminator classifier D_{ϕ} together
32 with the features of the learned critic: $s_{\phi}(\mathbf{x}) = [\nabla_{\phi} D_{\phi}(\mathbf{x}), h_1(\mathbf{x}), \dots, h_n(\mathbf{x})]$ where $h_1(\mathbf{x}), \dots, h_n(\mathbf{x})$
33 are the hidden activations of the learned critic. Both features and gradients are needed: the gradients
34 $\nabla_{\phi} D_{\phi}(\mathbf{x})$ are required to ensure the estimator for the parameters θ is consistent, since the number
35 of moments $s(\mathbf{x})$ needs to be larger than the number of parameters θ , which will be true if the critic
36 will have more parameters than the model; the features $h_i(\mathbf{x})$ are added since they have been shown
37 empirically to improve performance, thus showcasing the importance of the choice of test statistics s
38 used to train implicit models.
39

40 26.2.6 On density ratios and differences

41 We have seen how density ratios (Sections 26.2.2 and 26.2.3) and density differences (Section 26.2.4)
42 can be used to define training objectives for implicit generative models. We now explore some of the
43 distinctions between using ratios and differences for learning by comparison, as well as explore the
44



(a) Failure of the KL divergence to distinguish between distributions with non-overlapping support:
 $D_{\text{KL}}(p^* \parallel q_{\theta_1}) = D_{\text{KL}}(p^* \parallel q_{\theta_2}) = \infty$, despite q_{θ_2} being closer to p^* than q_{θ_1} .

(b) The density ratio $\frac{p^*}{q_{\theta}}$ used by the KL divergence and a smooth estimate given by an MLP, together with the gradient it provides with respect to the input variable.

Figure 26.4: The KL divergence cannot provide learning signal for distributions without overlapping support (left), while the smooth approximation given by a learned decision surface like an MLP can (right). Generated by [ipm_divergences.ipynb](#)

effects of using approximations to these objectives using function classes such as neural networks has on these distinctions.

One often stated downside of using divergences that rely on density ratios (such as f -divergences) is their poor behavior when the distributions p^* and q_θ do not have overlapping support. For non-overlapping support, the density ratio $\frac{p^*}{q_\theta}$ will be ∞ in the parts of the space where $p^*(\mathbf{x}) > 0$ but $q_\theta(\mathbf{x}) = 0$, and 0 otherwise. In that case, the $D_{\text{KL}}(p^* \parallel q_\theta) = \infty$ and the $JSD(p^*, q_\theta) = \log 2$, regardless of the value of θ . Thus f -divergences cannot distinguish between different model distributions when they do not have overlapping support with the data distribution, as visualized in Figure 26.4a. This is in contrast with difference based methods such as IPMs such as the Wasserstein distance and the MMD, which have smoothness requirements built in the definition of the method, by constraining the norm of the critic (Equations (26.29) and (26.34)). We can see the effect of these constraints in Figure 26.3: both the Wasserstein distance and the MMD provide useful signal in the case of distributions with non-overlapping support.

While the definition of f -divergences relies on density ratios (Equation (26.21)), we have seen that to train implicit generative models we use approximations to those divergences obtained using a parametric critic D_ϕ . If the function family of the critic used to approximate the divergence (via the bound or class probability estimation) contains only smooth functions, it will not be able to model the sharp true density ratio, which jumps from 0 to ∞ , but it can provide a smooth approximation. We show an example in Figure 26.4b, where we show the density ratio for two distributions without overlapping support and an approximation provided by an MLP trained to approximate the KL divergence using Equation 26.25. Here, the smooth decision surface provided by the MLP can be used to train a generative model while the underlying KL divergence cannot be; the learned MLP provides the gradient signal on how to move distribution mass to areas with more density under the data distribution, while the KL divergence provides a zero gradient almost everywhere in the

space. This ability of approximations to f -divergences to overcome non-overlapping support issues is a desirable property of generative modeling training criteria, as it allows models to learn the data distribution regardless of initialization [Fed+18]. Thus while the case of non-overlapping support provides an important theoretical difference between IPMs and f -divergences, it is less significant in practice since bounds on f -divergences or class probability estimation are used with smooth critics to approximate the underlying divergence.

Some density ratio and density difference based approaches also share commonalities: bounds are used both for f -divergences (variational bounds in Equation 26.25) and for the Wasserstein distance (Equation (26.31)). These bounds to distributional divergence and distances have their own set of challenges: since the generator minimizes a lower bound of the underlying divergence or distance, minimizing this objective provides no guarantees that the divergence will decrease in training. To see this, we can look at Equation 26.26: its RHS can get arbitrarily low without decreasing the LHS, the divergence we are interested in minimizing; this is unlike variational *upper* bound on the KL divergence used to train variational autoencoders Chapter 21.

16

26.3 Generative adversarial networks

We have looked at different learning principles that do not require the use of explicit likelihoods, and thus can be used to train implicit models. These learning principles specify training criteria, but do not tell us how to *train* models or parameterize models. To answer these questions, we now look at algorithms for training implicit models, where the models (both the discriminator and generator) are deep neural networks; this leads us to generative adversarial networks (GANs). We cover how to turn learning principles into loss functions for training GANs (Section 26.3.1); how to train models using gradient descent (Section 26.3.2); how to improve GAN optimization (Section 26.3.4) and how to assess GAN convergence (Section 26.3.5).

27

26.3.1 From learning principles to loss functions

In Section 26.2 we discussed learning principles for implicit generative models: class probability estimation, bounds on f -divergences, integral probability metrics and moment matching. These principles can be used to formulate loss functions to train the model parameters θ and the critic parameters ϕ . Many of these objectives use **zero-sum losses** via a **min-max** formulation: the generator's goal is to minimize the same function the discriminator is maximizing. We can formalize this as:

$$\min \max V(\phi, \theta) \tag{26.41}$$

As an example, we recover the original GAN with the Bernoulli log-loss (Equation (26.19)) when

$$V(\phi, \theta) = \frac{1}{2} \mathbb{E}_{p^*(\mathbf{x})} [\log D_\phi(\mathbf{x})] + \frac{1}{2} \mathbb{E}_{q_\theta(\mathbf{x})} [\log(1 - D_\phi(\mathbf{x}))]. \tag{26.42}$$

The reason most of the learning principles we have discussed lead to zero-sum losses is due to their underlying structure: the critic maximizes a quantity in order to approximate a divergence or distance — such as an f -divergence or Integral Probability Metric — and the model minimizes this approximation to the divergence or distance. That need not be the case, however. Intuitively,

the discriminator training criteria needs to ensure that the discriminator can distinguish between data and model samples, while the generator loss function needs to ensure that model samples are indistinguishable from data according to the discriminator.

To construct a GAN that is not zero-sum, consider the zero-sum criteria in the original GAN (Equation 26.42), induced by the Bernoulli scoring rule. The discriminator tries to distinguish between data and model samples by classifying the data as real (label 1) and samples as fake (label 0), while the goal of the generator is to minimize the probability that the discriminator classifies its samples as fake: $\min_{\theta} \mathbb{E}_{q_{\theta}(\mathbf{x})} \log(1 - D_{\phi}(\mathbf{x}))$. An equally intuitive goal for the generator is to maximize the probability that the discriminator classifies its samples as real. While the difference might seem subtle, this loss, known as the “nonsaturating loss” [Goo+14], defined as $\mathbb{E}_{q_{\theta}(\mathbf{x})} - \log D_{\phi}(\mathbf{x})$, enjoys better gradient properties early in training, as shown in Figure 26.5: the non-saturating loss provides a stronger learning signal (via the gradient) when the generator is performing poorly, and the discriminator can easily distinguish its samples from data, i.e., $D(G(\mathbf{z}))$ is low; more on the gradients properties the saturating and non-saturating losses can be found in [AB17; Fed+18].

There exist many other GAN losses which are not zero-sum, including formulations of LS-GAN [Mao+17], GANs trained using the hinge loss [LY17], and RelativisticGANs [JM18]. We can thus generally write a GAN formulation as follows:

$$\min_{\phi} L_D(\phi, \theta); \quad \min_{\theta} L_G(\phi, \theta). \quad (26.43)$$

We recover the zero-sum formulations if $-L_D(\phi, \theta) = L_G(\phi, \theta) = V(\phi, \theta)$. Despite departing from the zero-sum structure, the nested form of the optimization remains in the general formulation, as we will discuss in Section 26.3.2.

The loss functions for the discriminator and generator, L_D and L_G respectively, follow the general form in Equation 26.5, which allows them to be used to efficiently train implicit generative models. The majority of loss functions considered here can thus be written as follows:

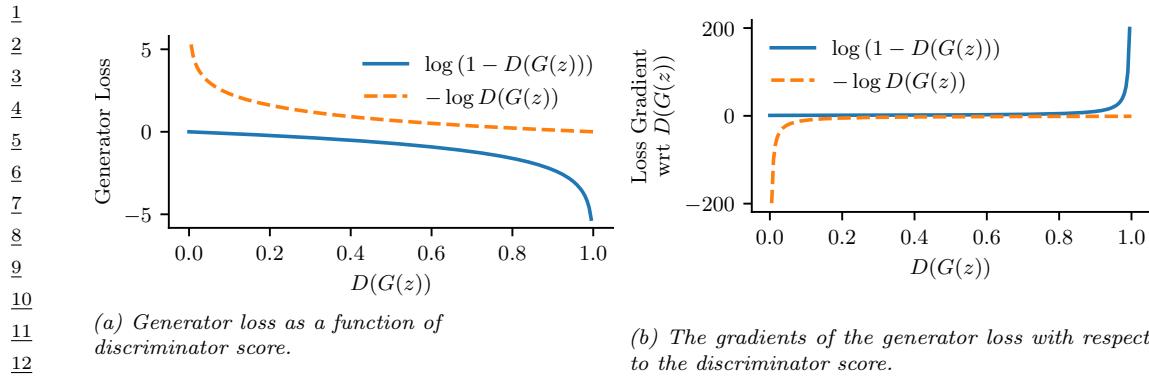
$$L_D(\phi, \theta) = \mathbb{E}_{p^*(\mathbf{x})} g(D_{\phi}(\mathbf{x})) + \mathbb{E}_{q_{\theta}(\mathbf{x})} h(D_{\phi}(\mathbf{x})) = \mathbb{E}_{p^*(\mathbf{x})} g(D_{\phi}(\mathbf{x})) + \mathbb{E}_{q(\mathbf{z})} h(D_{\phi}(G_{\theta}(\mathbf{z}))) \quad (26.44)$$

$$L_G(\phi, \theta) = \mathbb{E}_{q_{\theta}(\mathbf{x})} l(D_{\phi}(\mathbf{x})) = \mathbb{E}_{q(\mathbf{z})} l(D_{\phi}(G_{\theta}(\mathbf{z}))) \quad (26.45)$$

where $g, h, l : \mathbb{R} \rightarrow \mathbb{R}$. We recover the original GAN for $g(t) = -\log t$, $h(t) = -\log(1 - t)$ and $l(t) = \log(1 - t)$; the non-saturating loss for $g(t) = -\log t$, $h(t) = -\log(1 - t)$ and $l(t) = -\log(t)$; the Wasserstein distance formulation for $g(t) = t$, $h(t) = -t$ and $l(t) = t$; for f -divergences $g(t) = t$, $h(t) = -f^{\dagger}(t)$ and $l(t) = f^{\dagger}(t)$.

26.3.2 Gradient descent

GANs employ the learning principles discussed above in conjunction with gradient based learning for the parameters of the discriminator and generator. We assume a general formulation with a discriminator loss function $L_D(\phi, \theta)$ and a generator loss function $L_G(\phi, \theta)$. Since the discriminator is often introduced to approximate a distance or divergence $D(p^*, q_{\theta})$ (Section 26.2), for the generator to minimize a good approximation of that divergence one should solve the discriminator optimization fully for each generator update. That would entail that for each generator update one would first find the optimal discriminator parameters $\phi^* = \operatorname{argmin}_{\phi} L_D(\phi, \theta)$ in order to perform a gradient update given by $\nabla_{\theta} L_G(\phi^*, \theta)$. Fully solving the inner optimization problem $\phi^* = \operatorname{argmin}_{\phi} L_D(\phi, \theta)$



(a) Generator loss as a function of discriminator score.

(b) The gradients of the generator loss with respect to the discriminator score.

Figure 26.5: Saturating $\log(1 - D(G(z)))$ vs non-saturating $-\log D(G(z))$ loss functions. The non-saturating loss provides stronger gradients when the discriminator is easily detecting that generated samples are fake.
Generated by [gan_loss_types.ipynb](#)

for each optimization step of the generator is computationally prohibitive, which motivates the use of alternating updates: performing a few gradient steps to update the discriminator parameters, followed by a generator update. Note that when updating the discriminator, we keep the generator parameters fixed, and when updating the generator, we keep the discriminator parameters fixed. We show a general algorithm for these alternative updates in Algorithm 26.1.

Algorithm 26.1: General GAN training algorithm with alternating updates

```

26 1 Initialize  $\phi, \theta$ 
27 2 for each training iteration do
28 3   for  $K$  steps do
29 4     Update the discriminator parameters  $\phi$  using the gradient  $\nabla_\phi L_D(\phi, \theta)$ ;
30 5     Update the generator parameters  $\theta$  using the gradient  $\nabla_\theta L_G(\phi, \theta)$ 
31 6 Return  $\phi, \theta$ 

```

We are thus interested in computing $\nabla_\phi L_D(\phi, \theta)$ and $\nabla_\theta L_G(\phi, \theta)$. Given the choice of loss functions follows the general form in Equations 26.44 and 26.45 both for the discriminator and generator, we can compute the gradients that can be used for training. To compute the discriminator gradients, we write:

$$\nabla_\phi L_D(\phi, \theta) = \nabla_\phi [\mathbb{E}_{p^*(\mathbf{x})} g(D_\phi(\mathbf{x})) + \mathbb{E}_{q_\theta(\mathbf{x})} h(D_\phi(\mathbf{x}))] \quad (26.46)$$

$$= \mathbb{E}_{p^*(\mathbf{x})} \nabla_\phi g(D_\phi(\mathbf{x})) + \mathbb{E}_{q_\theta(\mathbf{x})} \nabla_\phi h(D_\phi(\mathbf{x})) \quad (26.47)$$

where $\nabla_\phi g(D_\phi(\mathbf{x}))$ and $\nabla_\phi h(D_\phi(\mathbf{x}))$ can be computed via backpropagation, and each expectation can be estimated using Monte Carlo estimation. For the generator, we would like to compute the gradient:

$$L_G(\phi, \theta) = \nabla_\theta \mathbb{E}_{q_\theta(\mathbf{x})} l(D_\phi(\mathbf{x})) \quad (26.48)$$

Here we cannot change the order of differentiation and integration since the distribution under the integral depends on the differentiation parameter θ . Instead, we will use that $q_\theta(\mathbf{x})$ is the distribution induced by an implicit generative model (also known as the “reparameterization trick”, see Section 6.3.5):

$$\nabla_\theta L_G(\phi, \theta) = \nabla_\theta \mathbb{E}_{q_\theta(\mathbf{x})} l(D_\phi(\mathbf{x})) = \nabla_\theta \mathbb{E}_{q(\mathbf{z})} l(D_\phi(G_\theta(\mathbf{z}))) = \mathbb{E}_{q(\mathbf{z})} \nabla_\theta l(D_\phi(G_\theta(\mathbf{z}))) \quad (26.49)$$

and again use Monte Carlo estimation to approximate the gradient using samples from the prior $q(\mathbf{z})$. Replacing the choice of loss functions and Monte Carlo estimation in Algorithm 26.1 leads to Algorithm 26.2, which is often used to train GANs.

Algorithm 26.2: GAN training algorithm

```

1 Initialize  $\phi, \theta$ 
2 for each training iteration do
3   for  $K$  steps do
4     Sample minibatch of  $M$  noise vectors  $\mathbf{z}_m \sim q(\mathbf{z})$ 
5     Sample minibatch of  $M$  examples  $\mathbf{x}_m \sim p^*(\mathbf{x})$ 
6     Update the discriminator by performing stochastic gradient descent using this gradient:
7        $\nabla_\phi \frac{1}{M} \sum_{m=1}^M [g(D_\phi(\mathbf{x}_m)) + \nabla_\phi h(D_\phi(G_\theta(\mathbf{z}_m)))]$ .
8     Sample minibatch of  $M$  noise vectors  $\mathbf{z}_m \sim q(\mathbf{z})$ 
9     Update the generator by performing stochastic gradient descent using this gradient:
10     $\nabla_\theta \frac{1}{M} \sum_{m=1}^M l(D_\phi(G_\theta(\mathbf{z}_m)))$ .
11
12  Return  $\phi, \theta$ 

```

26.3.3 Challenges with GAN training

Due to the adversarial game nature of GANs the optimizing dynamics of GANs are both hard to study in theory, and to stabilize in practice. GANs are known to suffer from **mode collapse**, a phenomenon where the generator converges to a distribution which does not cover not all the modes (peaks) of the data distribution, thus the model underfits the distribution. We show an example in Figure 26.6: while the data is a mixture of Gaussians with 16 modes, the model converges only to a few modes. Alternatively, another problematic behavior is **mode hopping**, where the generator “hops” between generating different modes of the data distribution. An intuitive explanation for this behavior is as follows: if the generator becomes good at generating data from one mode, it will generate more from that mode. If the discriminator cannot learn to distinguish between real and generated data in this mode, the generator has no incentive to expand its support and generate data from other modes. On the other hand, if the discriminator eventually learns to distinguish between the real and generated data inside this mode, the generator can simply move (hop) to a new mode, and this game of cat and mouse can continue.

While mode collapse and mode hopping are often associated with GANs, many improvements have made GAN training more stable, and these behaviors more rare. These improvements include using large batch sizes, increasing the discriminator neural capacity, using discriminator and generator regularization, as well as more complex optimization methods.

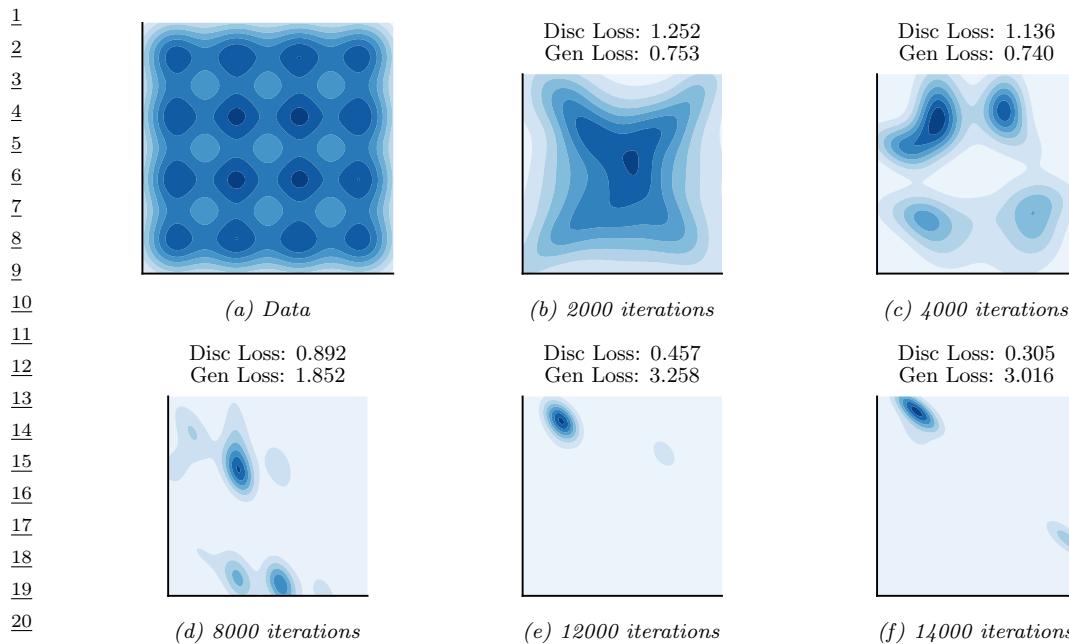


Figure 26.6: Illustration of mode collapse and mode hopping in GAN training. (a) The dataset, a mixture of 16 Gaussians in 2 dimensions. (b-f) Samples from the model after various amounts of training. Generated by [gan_mixture_of_gaussians.ipynb](#).

26.3.4 Improving GAN optimization

Hyperparameter choices such as the choice of momentum can be crucial when training GANs, with lower momentum values being preferred compared to the usual high momentum used in supervised learning. Algorithms such as Adam [KB14a] provide a great boost in performance [RMC16a]. Many other optimization methods have been successfully applied to GANs, such as those which target variance reduction [Cha+19c]; those which backpropagate through gradient steps, thus ensuring that generator does well against the discriminator *after it has been updated* [Met+16]; or using a local bilinear approximation of the two player game [SA19]. While promising, these advanced optimization methods tend to have a higher computational cost, making them harder to scale to large models or large datasets compared to less efficient optimization methods.

26.3.5 Convergence of GAN training

The challenges with GAN optimization make it hard to quantify when convergence has occurred. In Section 26.2 we saw how global convergence guarantees can be provided under optimality conditions for multiple objectives constructed starting with different distributional divergences and distances: if the discriminator is optimal, the generator is minimizing a distributional divergence or distance between the data and model distribution, and thus under infinite capacity and perfect optimization can learn the data distribution. This type of argument has been used since the original GAN paper [Goo+14]

to connect GANs to standard objectives in generative models, and obtain the associated theoretical guarantees. From a game theory perspective, this type of convergence guarantee provides an existence proof of a global Nash equilibrium for the GAN game, though under strong assumptions. A Nash equilibrium is achieved when both players (the discriminator and generator) would incur a loss if they decide to act by changing their parameters. Consider the original GAN defined by the objective in Equation 26.19; then $q_\theta = p^*$ and $D_\phi(\mathbf{x}) = \frac{p^*(\mathbf{x})}{p^*(\mathbf{x}) + q_\theta(\mathbf{x})} = \frac{1}{2}$ is a global Nash equilibrium, since for a given q_θ , the ratio $\frac{p^*(\mathbf{x})}{p^*(\mathbf{x}) + q_\theta(\mathbf{x})}$ is the optimal discriminator (Equation 26.11), and given an optimal discriminator, the data distribution is the optimal generator as it is the minimizer of the Jensen-Shannon divergence (Equation 26.15).

While these global theoretical guarantees provide useful insights about the GAN game, they do not account for optimization challenges that arise with accounting for the optimization trajectories of the two players, or for neural network parameterization since they assume infinite capacity both for the discriminator and generator. In practice GANs do not decrease a distance or divergence at every optimization step [Fed+18] and global guarantees are difficult to obtain when using optimization methods such as gradient descent. Instead, the focus shifts towards local convergence guarantees, such as reaching a local Nash equilibrium. A local Nash equilibrium requires that both players are at a local, not global minimum: a local Nash equilibrium is a stationary point (the gradients of the two loss functions are zero, i.e. $\nabla_\phi L_D(\phi, \theta) = \mathbf{0}$ and $\nabla_\theta L_G(\phi, \theta) = \mathbf{0}$), and the eigenvalues of the Hessian of each player ($\nabla_\phi \nabla_\phi L_D(\phi, \theta)$ and $\nabla_\theta \nabla_\theta L_G(\phi, \theta)$) are non-negative; for a longer discussion on Nash equilibria in continuous games, see [RBS16]. For the general GAN game, it is not guaranteed that a local Nash equilibrium always exists [FO20], and weaker conditions such as stationarity or locally stable stationarity have been studied [Ber+19]; other equilibrium definitions inspired by game theory have also been used [JNJ20; HLC19].

To motivate why convergence analysis is important in the case of GANs, we visualize an example of a GAN that does not converge trained with gradient descent. In DiracGAN [MGN18a] the data distribution $p^*(\mathbf{x})$ is the Dirac delta distribution with mass at zero. The generator is modeling a Dirac delta distribution with parameter θ : $G_\theta(z) = \theta$ and the discriminator is a linear function of the input with learned parameter ϕ : $D_\phi(x) = \phi x$. We also assume a GAN formulation where $g = h = -l$ in the general loss functions L_D and L_G defined above, see Equations (26.44) and (26.45). This results in the zero-sum game given by:

$$L_D = \mathbb{E}_{p^*(x)} - l(D_\phi(x)) + \mathbb{E}_{q_\theta(x)} - l(D_\phi(x)) = -l(0) - l(\theta\phi) \quad (26.50)$$

$$L_G = \mathbb{E}_{p^*(x)} l(D_\phi(x)) + \mathbb{E}_{q_\theta(x)} l(D_\phi(x)) = +l(0) + l(\theta\phi) \quad (26.51)$$

where l depends on the GAN formulation used ($l(z) = -\log(1 + e^{-z})$ for instance). The unique equilibrium point is $\theta = \phi = 0$. We visualize the DiracGAN problem in Figure 26.7 and show that DiracGANs with alternating gradient descent (Algorithm 26.1) do not reach the equilibrium point, but instead takes a circular trajectory around the equilibrium.

There are two main theoretical approaches taken to understand GAN convergence behavior around an equilibrium: by analyzing either the discrete dynamics of gradient descent, or the underlying continuous dynamics of the game using approaches such as stability analysis. To understand the difference between the two approaches, consider the discrete dynamics defined by gradient descent

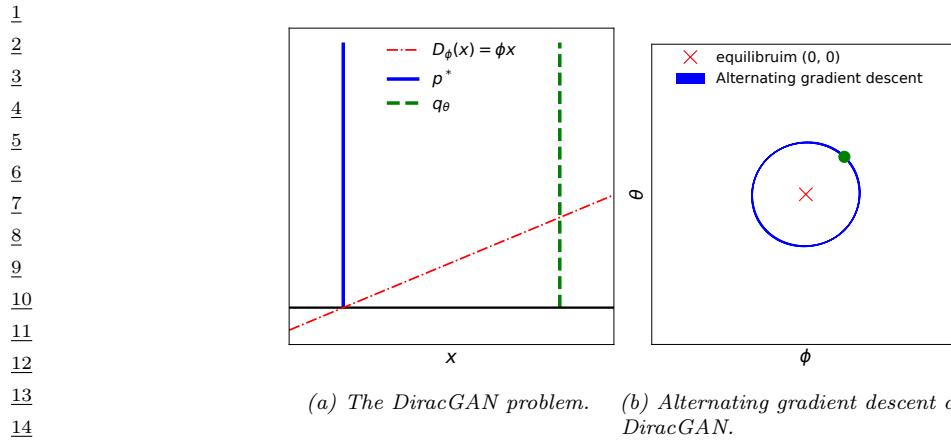


Figure 26.7: Visualizing divergence using a simple GAN: DiracGAN. Generated by [dirac_gan.ipynb](#)

with learning rates αh and λh , either via alternating updates (as we have seen in Algorithm 26.1):

$$\phi_t = \phi_{t-1} - \alpha h \nabla_\phi L_D(\phi_{t-1}, \theta_{t-1}), \quad (26.52)$$

or simultaneous updates, where instead of alternating the gradient updates between the two players, they are both updated simultaneously:

$$\phi_t = \phi_{t-1} - \alpha h \nabla_\phi L_D(\phi_{t-1}, \theta_{t-1}), \quad (26.54)$$

$$\theta_t = \theta_{t-1} - \lambda h \nabla_\theta L_G(\phi_t, \theta_{t-1}) \quad (26.55)$$

The above dynamics of gradient descent are obtained using Euler numerical integration from the ODEs that describes the game dynamics of the two players:

$$\dot{\phi} = -\nabla_\phi L_D(\phi, \theta), \quad (26.56)$$

$$\dot{\theta} = -\nabla_\theta L_G(\phi, \theta) \quad (26.57)$$

One approach to understand the behavior of GANs is to study these underlying ODEs, which, when discretized, result in the gradient descent updates above, rather than directly studying the discrete updates. These ODEs can be used for stability analysis to study the behavior around an equilibrium. This entails finding the eigenvalues of the Jacobian of the game

$$J = \begin{bmatrix} -\nabla_\phi \nabla_\phi L_D(\phi, \theta) & -\nabla_\theta \nabla_\phi L_D(\phi, \theta) \\ -\nabla_\phi \nabla_\theta L_G(\phi, \theta) & -\nabla_\theta \nabla_\theta L_G(\phi, \theta) \end{bmatrix} \quad (26.58)$$

evaluated at a stationary point (i.e., where $\nabla_\phi L_D(\phi, \theta) = 0$ and $\nabla_\theta L_G(\phi, \theta) = 0$). If the eigenvalues of the Jacobian all have negative real parts, then the system is asymptotically stable around the equilibrium; if at least one eigenvalue has positive real part, the system is unstable around the

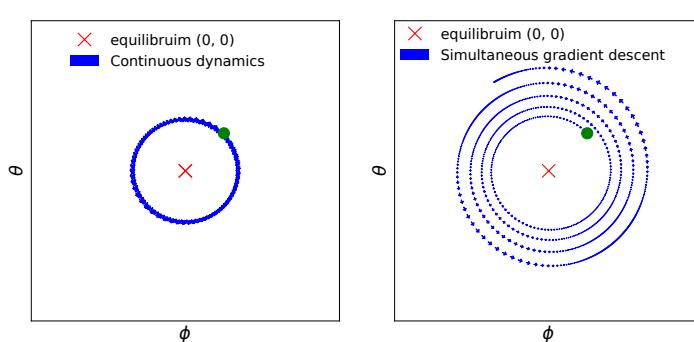


Figure 26.8: Continuous (left) and discrete dynamics (right) take different trajectories in DiracGAN. Generated by [dirac_gan.ipynb](#)

equilibrium. For the DiracGAN, the Jacobian evaluated at the equilibrium $\theta = \phi = 0$ is:

$$J = \begin{bmatrix} \nabla_\phi \nabla_\phi l(\theta\phi) + l(0) & \nabla_\theta \nabla_\phi l(\theta\phi) + l(0) \\ -\nabla_\phi \nabla_\theta (l(\theta\phi) + l(0)) & -\nabla_\theta \nabla_\theta (l(\theta\phi) + l(0)) \end{bmatrix} = \begin{bmatrix} 0 & l'(0) \\ -l'(0) & 0 \end{bmatrix} \quad (26.59)$$

where eigenvalues of this Jacobian are $\lambda_{\pm} = \pm il'(0)$. This is interesting, as the real parts of the eigenvalues are both 0; this result tells us that there is no asymptotic convergence to an equilibrium, but linear convergence could still occur. In this simple case we can reach the conclusion that convergence does not occur as we observe that there is a preserved quantity in this system, as $\theta^2 + \phi^2$ does not change in time (Figure 26.8, left):

$$\frac{d(\theta^2 + \phi^2)}{dt} = 2\theta \frac{d\theta}{dt} + 2\phi \frac{d\phi}{dt} = -2\theta l'(\theta\phi)\phi + 2\phi l'(\theta\phi)\theta = 0.$$

Using stability analysis to understand the underlying continuous dynamics of GANs around an equilibrium has been used to show that explicit regularization can help convergence [NK17; Bal+18]. Alternatively, one can directly study the updates of simultaneous gradient descent shown in Equations 26.54 and 26.55. Under certain conditions, [MNG17b] prove that GANs trained with simultaneous gradient descent reach a local Nash equilibrium. Their approach relies on assessing the convergence of series of the form $F^k(\mathbf{x})$ resulting from the repeated application of gradient descent update of the form $F(\mathbf{x}) = \mathbf{x} + hG(\mathbf{x})$, where h is the learning rate. Since the function F depends on the learning rate h , their convergence results depend on the size of the learning rate, which is not the case for continuous time approaches.

Both continuous and discrete approaches have been useful in understanding and improving GAN training; however, both approaches still leave a gap between our theoretical understanding and the most commonly used algorithms to train GANs in practice, such as alternating gradient descent or more complex optimizers used in practice, like Adam. Far from only providing different proof techniques, these approaches can reach different conclusions about the convergence of a GAN: we show an example in Figure 26.8, where we see that simultaneous gradient descent and the continuous dynamics behave differently when a large enough learning rate is used. In this case, the discretization error — the difference between the behavior of the continuous dynamics in Equations 26.56 and 26.57

and the gradient descent dynamics in Equations 26.54 and 26.55 — makes the analysis of gradient descent using continuous dynamics reach the wrong conclusion about DiracGAN [Ros+21]. This difference in behavior has been a motivator to train GANs with higher order numerical integrators such as RungeKutta4, which to more closely follow the underlying continuous system compared to gradient descent [Qin+20].

While optimization convergence analysis is an indispensable step in understanding GAN training and has led to significant practical improvements, it is worth noting that ensuring converge to an equilibrium does not ensure the model has learned a good fit of the data distribution. The loss landscape determined by the choice of L_D and L_G , as well as the parameterization of the discriminator and generator can lead to equilibria which do not capture the data distribution. The lack of distributional guarantees provided by game equilibria showcases the need to complement convergence analysis with work looking at the effect of gradient based learning in this game setting on the learned distribution.

15

26.4 Conditional GANs

17

We have thus far discussed how to use implicit generative models to learn a true unconditional distribution $p^*(\mathbf{x})$ from which we only have samples. It is often useful, however, to be able to learn conditional distributions of the from $p^*(\mathbf{x}|\mathbf{y})$. This requires having **paired data**, where each input \mathbf{x}_n is paired with a corresponding set of covariates \mathbf{y}_n , such as a class label, or a set of attributes or words, so $\mathcal{D} = \{(\mathbf{x}_n, \mathbf{y}_n) : n = 1 : N\}$, as in standard supervised learning. The conditioning variable can be discrete, like a class label, or continuous, such as an embedding which encodes other information. Conditional generative models are appealing since we can specify that we want the generated sample to be associated with conditioning information y , making them very amenable to real world applications, as we discuss in Section 26.7.

To be able to learn implicit conditional distributions $q_\theta(\mathbf{x}|\mathbf{y})$, we require datasets that specify the conditioning information associated with data, and we have to adapt the model architectures and loss functions. In the GAN case, changing the loss function for the generative model can be done by changing the critic, since the critic is part of the loss function of the generator; it is important for the critic to provide learning signal accounting for conditioning information, by penalizing a generator which provides realistic samples but which ignore the provided conditioning.

If we do not change the form of the min-max game, but provide the conditioning information to the two players, a **conditional GAN** can be created from the original GAN game [MO14]:

$$\min_{\theta} \max_{\phi} \frac{1}{2} \mathbb{E}_{p(\mathbf{y})} \mathbb{E}_{p^*(\mathbf{x}|\mathbf{y})} [\log D_\phi(\mathbf{x}, \mathbf{y})] + \frac{1}{2} \mathbb{E}_{p(\mathbf{y})} \mathbb{E}_{q_\theta(\mathbf{x}|\mathbf{y})} [\log(1 - D_\phi(\mathbf{x}, \mathbf{y}))] \quad (26.60)$$

In the case of implicit latent variable models, the embedding information becomes an additional input to the generator, together with the latent variable \mathbf{z} :

$$\min_{\theta} \max_{\phi} \mathcal{L}(\theta, \phi) = \frac{1}{2} \mathbb{E}_{p(\mathbf{y})} \mathbb{E}_{p^*(\mathbf{x}|\mathbf{y})} [\log D_\phi(\mathbf{x}, \mathbf{y})] + \frac{1}{2} \mathbb{E}_{p(\mathbf{y})} \mathbb{E}_{q(\mathbf{z})} [\log(1 - D_\phi(\mathcal{G}_\theta(\mathbf{z}, \mathbf{y}), \mathbf{y}))] \quad (26.61)$$

For discrete conditioning information such as labels, one can also add a new loss function, by training a critic which does not only learn to distinguish between real and fake data, but learns to classify both data and generated samples as pertaining to one of the K classes provided in the

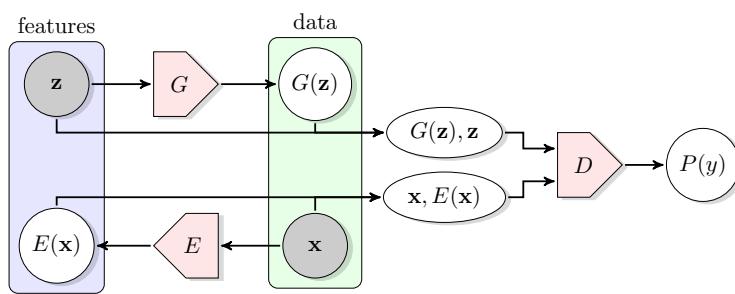


Figure 26.9: Learning an implicit posterior using an adversarial approach, as done in BiGAN. From Figure 1 of [DKD16]. Used with kind permission of Jeff Donahue.

dataset [OOS17]:

$$\mathcal{L}_c(\theta, \phi) = - \left[\frac{1}{2} \mathbb{E}_{p(y)} \mathbb{E}_{p^*(x|y)} [\log D_\phi(y|x)] + \frac{1}{2} \mathbb{E}_{p(y)} \mathbb{E}_{q_\theta(x|y)} [\log(D_\phi(y|x))] \right] \quad (26.62)$$

Note that while we could have two critics, one unsupervised critic and one supervised which maximizes the equation above, in practice the same critic is used, to aid shaping the features used in both decision surfaces. Unlike the adversarial nature of the unsupervised game, it is in the interest of both players to minimize the classification loss \mathcal{L}_c . Thus together with the adversarial dynamics provided by \mathcal{L} , the two players are trained as follows:

$$\max_{\phi} \mathcal{L}(\theta, \phi) - \mathcal{L}_c(\theta, \phi) \quad \min_{\theta} \mathcal{L}(\theta, \phi) + \mathcal{L}_c(\theta, \phi) \quad (26.63)$$

In the case of conditional latent variable models, the latent variable controls the sample variability *inside* the mode specified by the conditioning information. In early conditional GANs, the conditioning information was provided as additional input to the discriminator and generator, for example by concatenating the conditioning information to the latent variable z in the case of the generator; it has been since observed that it is important to provide the conditioning information at various layers of the model, both for the generator and the discriminator [DV+17; DSK16] or use a projection discriminator [MK18].

26.5 Inference with GANs

Unlike other latent variable models such as variational autoencoders, GANs do not define an inference procedure associated with the generative model. To deploy the principles behind GANs to find a posterior distribution $p(z|x)$, multiple approaches have been taken, from combining GANs and variational autoencoders via hybrid methods [MNG17a; Sri+17; Lar+16; Mak+15b] to constructing inference methods catered to implicit variable models [Dum+16; DKD16; DS19]. An overview of these methods can be found in [Hus17b].

GAN based methods which perform inference and learn **implicit posterior distribution** $p(z|x)$ introduce changes to the GAN algorithm to do so. An example of such a method is **BiGAN** (bidirectional GAN) [DKD16] or **ALI** (adversarially learned inference) [Dum+16], which trains an

¹ implicit parameterized encoder E_ζ to map input \mathbf{x} to latent variables \mathbf{z} . To ensure consistency between
² the encoder E_ζ and the generator G_θ , an adversarial approach is introduced with a discriminator
³ D_ϕ learning to distinguish between pairs of data and latent samples: D_ϕ learns to consider pairs
⁴ $(\mathbf{x}, E_\zeta(\mathbf{x}))$ with $\mathbf{x} \sim p^*$ as real, while $(G_\theta(\mathbf{z}), \mathbf{z})$ with $\mathbf{z} \sim q(\mathbf{z})$ is considered fake. This approach,
⁵ shown in Figure 26.9, ensures that the joint distributions are matched, and thus the marginal
⁶ distribution $q_\theta(\mathbf{x})$ given by G_θ should learn $p^*(\mathbf{x})$, while the conditional distribution $p_\zeta(\mathbf{z}|\mathbf{x})$ given by
⁷ E_ζ should learn $q_\theta(\mathbf{z}|\mathbf{x}) = \frac{q_\theta(\mathbf{x}, \mathbf{z})}{q_\theta(\mathbf{x})} \propto q_\theta(\mathbf{x}|\mathbf{z})q(\mathbf{z})$. This joint GAN loss can be used both to train the
⁸ generator G_θ and the encoder E_ζ , without requiring a reconstruction loss common in other inference
⁹ methods. While not using a reconstruction loss, this objective retains the property that under
¹⁰ global optimality conditions the encoder and decoder are inverses of each other: $E_\theta(G_\zeta(\mathbf{z})) = \mathbf{z}$ and
¹¹ $G_\zeta(E_\theta(\mathbf{x})) = \mathbf{x}$. (See also Section 21.2.4 for a discussion of how VAEs learn to ensure $p^*(\mathbf{x})p_\zeta(\mathbf{z}|\mathbf{x})$
¹² matches $p(\mathbf{z})p_\theta(\mathbf{x}|\mathbf{z})$ using an explicit model of the data.)

¹⁴

¹⁵

¹⁶

¹⁷ 26.6 Neural architectures in GANs

¹⁸

¹⁹ We have so far discussed the learning principles, algorithms, and optimization methods that can
²⁰ be used to train implicit generative models parameterized by deep neural networks. We have not
²¹ discussed, however, the importance of the choice of neural network architectures for the model and
²² the critic, choices which have fueled the progress in GAN generation since their conception. We will
²³ look at a few case studies which show the importance of information about data modalities into
²⁴ the critic and the generator (Section 26.6.1), employing the right inductive biases (Section 26.6.2),
²⁵ incorporating attention in GAN models (Section 26.6.3), progressive generation (Section 26.6.4),
²⁶ regularization (Section 26.6.5), and using large scale architectures (Section 26.6.6).

²⁷

²⁸

²⁹

³⁰ 26.6.1 The importance of discriminator architectures

³¹

³² Since the discriminator or critic is rarely optimal — either due to the use of alternating gradient
³³ descent or the lack of capacity of the neural discriminator — GANs do not perform distance or
³⁴ divergence minimization in practice. Instead, the critic acts as part of a **learned loss function**
³⁵ for the model (the generator). Every time the critic is updated, the loss function for the generative
³⁶ model changes; this is in stark contrast with divergence minimization such maximum likelihood
³⁷ estimation, where the loss function stays the same throughout the training of the model. Just as
³⁸ learning features of data instead of handcrafting them is a reason for the success of deep learning
³⁹ methods, learning loss functions advanced the state of the art of generative modeling. Critics that
⁴⁰ take data modalities into account — such as convolutional critics for images and recurrent critics
⁴¹ for sequential data such as text or audio — become part of data modality dependent loss functions.
⁴² This in turn provides modality-specific learning signal to the model, for example by penalizing blurry
⁴³ images and encouraging sharp edges, which is achieved due to the convolutional parameterization of
⁴⁴ the critic. Even within the same data modality, changes to critic architectures and regularization
⁴⁵ have been one of the main drivers in obtaining better GANs, since they affect the generator’s loss
⁴⁶ function, and thus also the *gradients of the generator* and have a strong effect on optimization.

⁴⁷

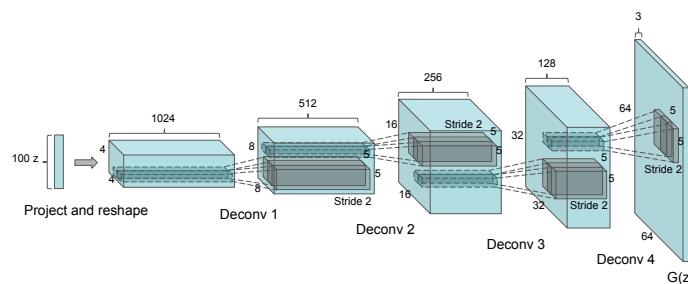


Figure 26.10: DCGAN convolutional generator. From Figure 1 of [RMC15]. Used with kind permission of Alec Radford.

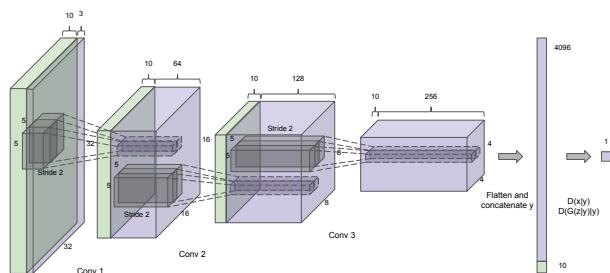


Figure 26.11: DCGAN convolutional discriminant. From Figure 1 of [RMC15]. Used with kind permission of Alec Radford.

26.6.2 Architectural inductive biases

While the original GAN paper used convolutions only sparingly, deep convolutional GAN (**DCGAN**) [RMC15] performed an extensive study on what architectures are most useful for GAN training, resulting in a set of useful guidelines that led to a substantial boost in performance. Without changing the learning principles behind GANs, DCGAN was able to obtain better results on image data by using convolutional generators (Figure 26.10) and critics, using BatchNormalization for both the generator and critic, replacing pooling layers with strided convolutions, using ReLU activation networks in the generator, and LeakyReLU activations in the discriminator. Many of these principles are still in use today, for larger architectures and with various loss functions. Since DCGAN, residual convolutional layers have become a key staple of both models and critics for image data [Gul+17], and recurrent architectures are used for sequence data such as text [SSG18b; Md+19].

26.6.3 Attention in GANs

Attention mechanisms are explained in detail in Section 16.2.7. In this section, we discuss how to use them for both the GAN generator and discriminator; this is called the self attention GAN or **SAGAN** model [Zha+19c]. The advantage of self attention is that it ensures that both discriminator and generator have access to a *global* view of other units of the same layer, unlike convolutional



15 *Figure 26.12: Attention queries used by a SAGAN model, showcasing the global span of attention. Each row*
16 *first shows the input image and a set of color coded query locations in the image. The subsequent images*
17 *show the attention maps corresponding to each query location in the first image, with the query color coded*
18 *location being shown, and arrows from it to the attention map are used to highlight the most attended regions.*
19 *From Figure 1 of [Zha+19c]. Used with kind permission of Han Zhang.*

20
21

22 layers. This is illustrated in Figure 26.12, which visualizes the global span of attention: query points
23 can attend to various other areas in the image.

24 The self-attention mechanism for convolutional features reshaped to $\mathbf{h} \in \mathbb{R}^{C \times N}$ is defined by
25 $\mathbf{f} = W_f \mathbf{h}$, $\mathbf{g} = W_g \mathbf{h}$, $\mathbf{S} = \mathbf{f}^T \mathbf{g}$, where $W_f \in \mathbb{R}^{C' \times C}$, $W_g \in \mathbb{R}^{C' \times C}$, where $C' \leq C$ is a hyperparameter.
26 From $\mathbf{S} \in \mathbb{R}^{N \times N}$, a probability row matrix β is obtained by applying the softmax operator for each
27 row, which is then used to *attend* to a linear transformation of the features $\mathbf{o} = W_o(W_h \mathbf{h})\beta^T \in \mathbb{R}^{C \times N}$,
28 using learned operators $W_h \in \mathbb{R}^{C' \times C}$, $W_o \in \mathbb{R}^{C \times C'}$. An output is then created by $\mathbf{y} = \gamma \mathbf{o} + \mathbf{h}$, where
29 $\gamma \in \mathbb{R}$ is a learned parameter.

30 Beyond providing global signal to the players, it is worth noting the flexibility of the self attention
31 mechanism. The learned parameter γ ensures that the model can decide not to use the attention
32 layer, and thus adding self attention does not restrict the set of possible models an architecture
33 can learn. Moreover, self attention significantly increases the number of parameters of the model
34 (each attention layer introduced 4 learned matrices \mathbf{W}_f , \mathbf{W}_g , \mathbf{W}_h , \mathbf{W}_o), an approach that has been
35 observed as a fruitful way to improve GAN training.

36

37 26.6.4 Progressive generation

38

39 One of the first successful approaches to generating higher resolution, color images from a GAN
40 is via an *iterative* process, by first generating a lower dimensional sample, and then using that as
41 conditioning information to generate a higher dimensional sample, and repeating the process until the
42 desired resolution is reached. **LapGAN** [DCF+15] uses a Laplacian pyramid as the iterative building
43 block, by first upsampling the lower dimensional samples using a simple upsampling operation, such
44 as smoothed upsampling, and then using a conditional generator to produce a residual to be added
45 to the upsampled version to produce the higher resolution sample. In turn, this higher resolution
46 sample can then be provided to another LapGAN layer to produce another, even higher resolution
47

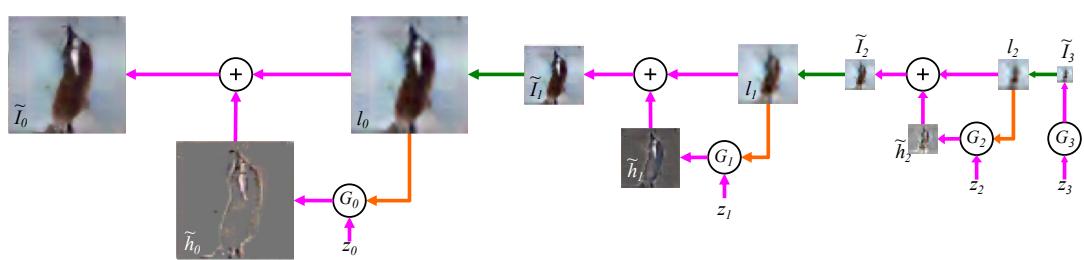


Figure 26.13: LapGAN generation algorithm: the generation process starts with a low dimension sample, which gets upscaled and residually added to the output of a generator at a higher resolution. The process gets repeated multiple times. From Figure 1 of [DCF+15]. Used with kind permission of Emily Denton.

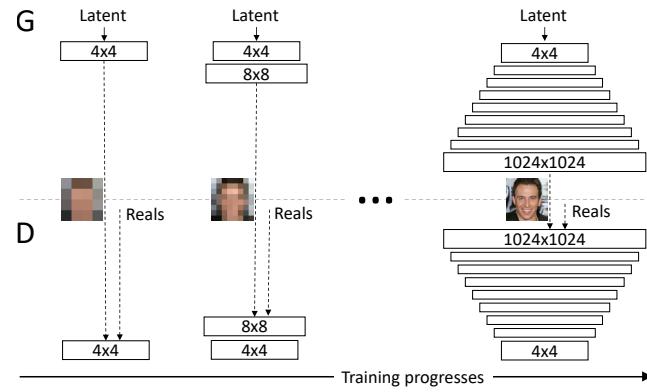


Figure 26.14: ProgressiveGAN training algorithm. The input to the discriminator at the bottom of the figure is either a generated image, or a real image (denotes as ‘Reals’ in the figure) at the corresponding resolution. From Figure 1 of [Kar+18]. Used with kind permission of Tero Karras.

sample, and so on — this process is shown in Figure 26.13. In LapGAN, a different generator and critic are trained for each iterative block of the model; in ProgressiveGAN [Kar+18] the lower resolution generator and critic are “grown”, by becoming part of the generator and critic used to learn to generate higher resolution samples. The higher resolution generator is obtained by adding new layers on top of the last layer of the lower resolution generator. A residual connection between an upscaled version of the lower dimensional sample and the output of the newly created higher resolution generator is added, which is annealed from 0 to 1 in training — transitioning from using the upscaled version of the lower dimensional sample early in training, to only using the sample of the higher resolution generator at the end of training. A similar change is done to the discriminator. Figure 26.14 shows the growing generator and discriminators in ProgressiveGAN training.

26.6.5 Regularization

Regularizing both the discriminator and the generator has by now a long tradition in GAN training. Regularizing GANs can be justified from multiple perspectives: theoretically, as it has been shown to

1 be tied to convergence analysis [MGN18b]; empirically, as it has been shown to help performance and
 2 stability in practice [RMC15; Miy+18c; Zha+19c; BDS18]; and intuitively, as it can be used to avoid
 3 overfitting in the discriminator and generator. Regularization approaches include adding noise to the
 4 discriminator input [AB17], adding noise to the discriminator and generator hidden features [ZML16],
 5 using BatchNorm for the two players [RMC15], adding dropout in the discriminator [RMC15],
 6 spectral normalization [Miy+18c; Zha+19c; BDS18], and gradient penalties (penalizing the norm
 7 of the discriminator gradient with respect to its input $\|\nabla_{\mathbf{x}} D_{\phi}(\mathbf{x})\|^2$) [Arb+18; Fed+18; ACB17;
 8 Gul+17]. Often regularization methods help training regardless of the type of loss function used,
 9 and have been shown to have effects both on training performance as well as the stability of the
 10 GAN game. However, improving stability and improving performance in GAN training can be at
 11 odds with each other, since too much regularization can make the models very stable but reduce
 12 performance [BDS18].
 13

14

15

16 26.6.6 Scaling up GAN models

17 By combining many of the architectural tricks discussed thus far — very large residual networks, self
 18 attention, spectral normalization both in the discriminator and the generator, BatchNormalization
 19 in the generator — one can train GANs to generating diverse, high quality data, as done with
 20 BigGAN [BDS18], StyleGAN [Kar+20c], and alias-free GAN [Kar+21]. Beyond combining carefully
 21 chosen architectures and regularization, creating large scale GANs also require changes in optimization,
 22 with large batch sizes being a key component. This furthers the view that the key components of the
 23 GAN game — the losses, the parameterization of the models, and the optimization method — have
 24 to be viewed collectively rather than in isolation.
 25

26

27

28 26.7 Applications

29 The ability to generate new plausible data enables a wide range of applications for GANs. This
 30 section will look at a set of applications that aim to demonstrate the breadth of GANs across different
 31 data modalities, such as images (Section 26.7.1), video (Section 26.7.2), audio (Section 26.7.3), and
 32 text (Section 26.7.4), and include applications such as imitation learning (Section 26.7.5), domain
 33 adapation (Section 26.7.6), and art (Section 26.7.7).
 34

35

36

37 26.7.1 GANs for image generation

38 The most widely studied application area is in image generation. We focus on the translation of one
 39 image to another using either paired or unpaired datasets. There are many other topics related to
 40 image GANs that we do not cover, and a more complete overview can be found in other sources,
 41 such as [Goo16] for the theory and [Bro19] for the practice. A JAX notebook which uses a small
 42 pretrained GAN to generate some face images can be found at [GAN_JAX_CelebA_demo.ipynb](#).
 43 We show the progression of quality in sample generation of faces using GANs in Figure 26.15. There
 44 is also increasing need to consider the generation of images with regards to the potential risks they
 45 can have when used in other domains, which involve discussions of synthetic media and **deep fakes**,
 46 and sources for discussion include [Bru+18; Wit].
 47



Figure 26.15: Increasingly realistic synthetic faces generated by different kinds of GAN, specifically (from left to right): original GAN [Goo+14], DCGAN [RMC15], CoupledGAN [LT16], ProgressiveGAN [Kar+18], StyleGAN [KLA19]. Used with kind permission of Ian Goodfellow. An online demo, which randomly generates face images using StyleGAN, can be found at <https://thispersondoesnotexist.com>.

26.7.1.1 Conditional image generation

Class-conditional image generation using GANs has become a very fruitful endeavor. BigGAN [BDS18] carries out class-conditional generation of ImageNet samples across a variety of categories, from dogs and cats to volcanoes and hamburgers. StyleGAN [KLA19] is able to generate high quality images of faces at high resolution by learning a conditioning style vector and the ProgressiveGAN architecture discussed in Section 26.6.4. By learning the conditioning vector they are able to generate samples which interpolate between the styles of other samples, for example by preserving coarser style elements such as pose or face shape from one sample, and smaller scale style elements such as hair style from another; this provides fine grained control over the style of the generated images.

26.7.1.2 Paired image-to-image generation

We have discussed in Section 26.4 how using paired data of the form $(\mathbf{x}_n, \mathbf{y}_n)$ can be used to build conditional generative models of $p(\mathbf{x}|\mathbf{y})$. In some cases, the conditioning variable \mathbf{y} has the same size and shape as the output variable \mathbf{x} . The resulting model $p_\theta(\mathbf{x}|\mathbf{y})$ can then be used to perform **image to image translation**, as illustrated in Figure 26.16, where \mathbf{y} is drawn from the **source domain**, and \mathbf{x} from the **target domain**. Collecting paired data of this form can be expensive, but in some cases, we can acquire it automatically. One such example is image colorization, where a paired dataset can easily be obtained by processing color images into grayscale images (see e.g., [Jas]).

A conditional GAN used for paired image-to-image translation was proposed in [Iso+17], and is known as the **pix2pix** model. It uses a U-net style architecture for the generator, as used for semantic segmentation tasks. However, they replace the batch normalization layers with instance normalization, as in neural style transfer.

For the discriminator, pix2pix uses a **patchGAN** model, that tries to classify local patches as being real or fake (as opposed to classifying the whole image). Since the patches are local, the discriminator is forced to focus on the style of the generated patches, and ensure they match the statistics of the target domain. A patch-level discriminator is also faster to train than a whole-image discriminator, and gives a denser feedback signal. This can produce results similar to Figure 26.16

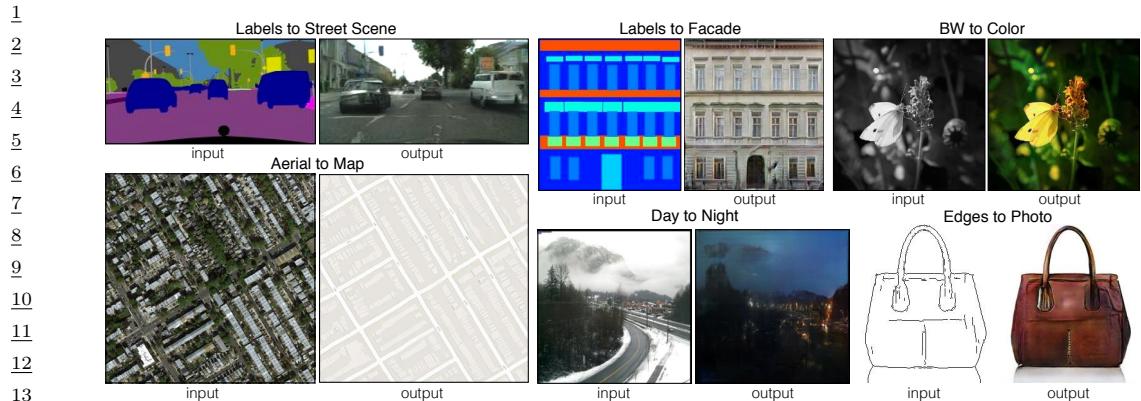


Figure 26.16: Example results on several image-to-image translation problems as generated by the pix2pix conditional GAN. From Figure 1 of [Iso+17]. Used with kind permission of Phillip Isola.

(depending on the dataset).

26.7.1.3 Unpaired image-to-image generation

A major drawback of conditional GANs is the need to collect paired data. It is often much easier to collect **unpaired data** of the form $\mathcal{D}_x = \{\mathbf{x}_n : n = 1 : N_x\}$ and $\mathcal{D}_y = \{\mathbf{y}_n : n = 1 : N_y\}$. For example, \mathcal{D}_x might be a set of daytime images, and \mathcal{D}_y a set of night-time images; it would be impossible to collect a paired dataset in which exactly the same scene is recorded during the day and night (except using a computer graphics engine, but then we wouldn't need to learn a generator).

We assume that the datasets \mathcal{D}_x and \mathcal{D}_y come from the marginal distributions $p(\mathbf{x})$ and $p(\mathbf{y})$ respectively. We would then like to fit a joint model of the form $p(\mathbf{x}, \mathbf{y})$, so that we can compute conditionals $p(\mathbf{x}|\mathbf{y})$ and $p(\mathbf{y}|\mathbf{x})$ and thus translate from one domain to another. This is called **unsupervised domain translation**.

In general, this is an ill-posed problem, since there are an infinite number of different joint distributions that are consistent with a set of marginals $p(\mathbf{x})$ and $p(\mathbf{y})$. We can try, however, to learn a joint distribution such that samples from it satisfy additional constraints. For example, if G is a conditional generator that maps a sample from \mathcal{X} to \mathcal{Y} , and F maps a sample from \mathcal{Y} to \mathcal{X} , it is reasonable to require that these be inverses of each other, i.e., $F(G(\mathbf{x})) = \mathbf{x}$ and $G(F(\mathbf{y})) = \mathbf{y}$. This is called a **cycle consistency** loss [Zhu+17]. We can encourage G and F to satisfy this constraint by using a penalty term on the difference between the starting image and the image we get after going through this cycle:

$$\mathcal{L}_{\text{cycle}} = \mathbb{E}_{p(\mathbf{x})} \|F(G(\mathbf{x})) - \mathbf{x}\|_1 + \mathbb{E}_{p(\mathbf{y})} \|G(F(\mathbf{y})) - \mathbf{y}\|_1 \quad (26.64)$$

To ensure that the outputs of G are samples from $p(\mathbf{y})$ and those of F are samples from $p(\mathbf{x})$, we use a standard GAN approach, introducing discriminators D_X and D_Y , which can be done using any choice of GAN loss \mathcal{L}_{GAN} , as visualized in Figure 26.17. Finally, we can optionally check that applying the conditional generator to images from its own domain does not change them:

$$\mathcal{L}_{\text{identity}} = \mathbb{E}_{p(\mathbf{x})} \|\mathbf{x} - F(\mathbf{x})\|_1 + \mathbb{E}_{p(\mathbf{y})} \|\mathbf{y} - G(\mathbf{y})\|_1 \quad (26.65)$$

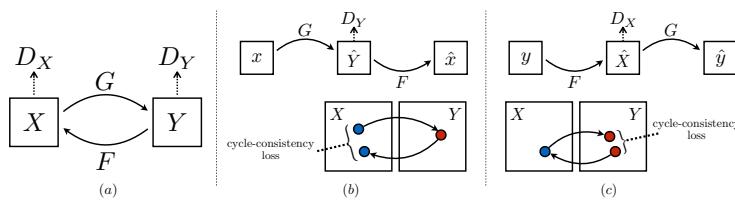


Figure 26.17: Illustration of the CycleGAN training scheme. (a) Illustration of the 4 functions that are trained. (b) Forwards cycle consistency from X back to X . (c) Backwards cycle consistency from Y back to Y . From Figure 3 of [Zhu+17]. Used with kind permission of Jun-Yan Zhu.

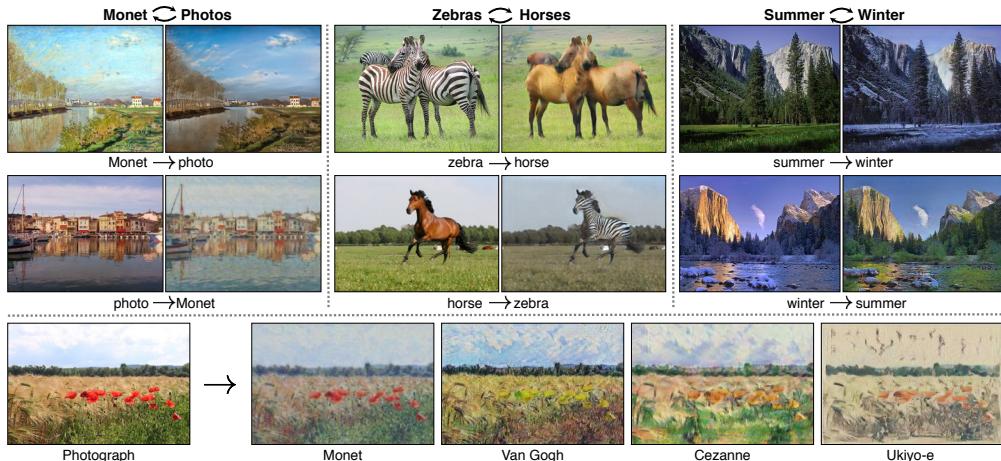


Figure 26.18: Some examples of unpaired image-to-image translation generated by the CycleGAN model. From Figure 1 of [Zhu+17]. Used with kind permission of Jun-Yan Zhu.

We can combine all three of these consistency losses to train the translation mappings F and G , using hyperparameters λ_1 and λ_2 :

$$\mathcal{L} = \mathcal{L}_{\text{GAN}} + \lambda_1 \mathcal{L}_{\text{cycle}} + \lambda_2 \mathcal{L}_{\text{identity}} \quad (26.66)$$

CycleGAN results on various datasets are shown in Figure 26.18. The bottom row shows how CycleGAN can be used for **style transfer**.

26.7.2 Video generation

The GAN framework can be expanded from individual images (frames) to videos; the techniques used to generate realistic images can also be applied to generate videos, with additional techniques required to ensure *spatio-temporal consistency*. Spatio-temporal consistency is obtained by ensuring that the discriminator has access to the real data and generated sequences in order, thus penalizing the generator when generating realistic individual frames without respecting temporal order [SMS17];

¹ [Sai+20](#); [CDS19](#); [Tul+18](#). Another discriminator can be employed to additionally ensure each frame
² is realistic [[Tul+18](#); [CDS19](#)]. The generator itself needs to have a temporal element, which is often
³ implemented through a recurrent component. As with images, the generation framework can be
⁴ expanded to video-to-video translation [[Ban+18](#); [Wan+18](#)], encompassing applications such as motion
⁵ transfer [[Cha+19a](#)].

8 26.7.3 Audio generation

10 Generative models have been demonstrated in the tasks of generating audio waveforms, as well
11 as for the task of text-to-speech (TTS) generation. Other types of generative models, such as
12 autoregressive models, such as WaveNet [oor+16] and WaveRNN [Kal+18b] have been developed for
13 these applications, although autoregressive models are difficult to parallelize over time since they
14 predict each time step of the audio sequentially and can be computationally expensive and too slow
15 to be used in practice. GANs provide an alternative approach for these tasks and other paths for
16 addressing these concerns.

Many different GAN architectures have been developed for audio-only generation, including generation of single note recordings from instruments by GANSynth, a vocoder model that uses GANs to generate magnitude spectrograms from mel-spectrograms [Eng+18], in voice conversion using a modified CycleGAN discussed above [Kan+20], and the direct generation of raw audio in WaveGAN [DMP18].

22 Initial work on GANs for TTS was developed [Yan+17] whose approach is similar to conditional
23 GANs for image generation (see Section 26.7.1.2), but uses 1d convolution instead of 2d. More
24 recent GANs such as GAN-TTS [Bi+19] use more advanced architectures and discriminators that
25 operate at multiple frequency scales that have performance that now matches the best performing
26 autoregressive models when assessed using mean opinion scores. In both the direct-audio generation,
27 the ability of GANs to allow faster generation and different types of context is the advantage that
28 makes them advantageous compared to other models.

30 26.7.4 Text generation

Similar to image and audio domains, there are several tasks for text data for which GAN-based approaches have been developed, including conditional text generation and text-style transfer. Text data are often represented as discrete values, at either the character level or the word-level, indicating membership within a set of a particular vocabulary size (alphabet size, or number of words). Due to the discrete nature of text, GAN models trained on text are *explicit*, since they explicitly model the probability distribution of the output, rather than modeling the sampling path. This is unlike most GAN models of continuous data such as images that we have discussed in the chapter so far, though explicit GANs of continuous data do exist [Die+19b].

40 The discrete nature of text is why maximum likelihood is one of the most common methods of
41 learning generative models of text. However, models trained with maximum likelihood are often
42 limited to autoregressive models, while like in the audio case, GANs make it possible to generate
43 text in a non-autoregressive manner, making other tasks possible, such as one-shot feedforward
44 generation [Gu+17].

45 The difficulty of generating discrete data such as text using GANs can be seen looking at their
46 loss function, such as in Equations (26.19), (26.21) and (26.28). GAN losses contain terms of
47

the form $\mathbb{E}_{q_\theta(\mathbf{x})} f(\mathbf{x})$, which we not only need to evaluate, but also backpropagate through, by computing $\nabla_\theta \mathbb{E}_{q_\theta(\mathbf{x})} f(\mathbf{x})$. In the case of implicit distributions given by latent variable models, we used the reparameterization trick to compute this gradient (Equation 26.49). In the discrete case, the reparameterization trick is not available and we have to look for other ways to estimate the desired gradient. One approach is to use the score function estimator, discussed in Section 6.3.4. However, the score function estimator exhibits high gradient variance, which can destabilize training. One common approach to avoid this issue is to pre-train the language model generator using maximum likelihood, and then to fine-tune with a GAN loss which gets backpropagated into the generator using the score-function estimator, as done by Sequence GAN [Yu+17], MaliGAN [Che+17], and RankGAN [Lin+17a]. While these methods spearheaded the use of GANs for text, they do not address the inherent instabilities of score function estimation and thus have to limit the amount of adversarial fine tuning to a small number of epochs and often use a small learning rate, keeping their performance close to that of the maximum-likelihood solution [SSG18a; Cac+18].

An alternative to maximum likelihood pretraining is to use other approaches to stabilize the score function estimator or to use continuous relaxations for backpropagation. ScratchGAN is a word-level model that uses large batch sizes and discriminator regularization to stabilize score function training (these techniques are the same that we have seen as stabilizers for training image GANs) [Md+19]. [Pre+17b] completely avoid the score function estimator and develop a character level model without pre-training, by using continuous relaxations and curriculum learning. These training approaches can also benefit from other architectural advances, e.g., [NPN19] showed that language GANs can benefit from complex architectures such as relation networks [San+17].

Finally, unsupervised text style transfer, mimicking image style transfer, have been proposed by [She+17; Fu+17] using adversarial classifiers to decode to a different style/language, or like [Pra+18] who trains different encoders, one per style, by combining the encoder of a pre-trained NMT and style classifiers, among other approaches.

26.7.5 Imitation learning

Imitation learning takes advantage of observations of expert demonstrations to learn action policies and reward functions of unknown environments by minimizing some form of discrepancy between the learned and the expert behaviors. There are many approaches available, including behavioral cloning [PPG91] that treats this problem as one of supervised learning, and inverse reinforcement learning [NR00b]. GANs are appealing for imitation learning since they provide a way to avoid the difficulty of designing good discrepancy functions for behaviors, and instead learn these discrepancy functions using a discriminator between trajectories generated by a learned agent and observed demonstrations.

This approach, known as generative adversarial imitation learning (**GAIL**) [HE16a] demonstrates the ability to use GANs for complex behaviors in high-dimensional environments. GAIL jointly learns a generator, which forms a stochastic policy, along with a discriminator that acts as a reward signal. Like we saw in the probabilistic development of GANs in the earlier sections, GAIL can also be generalized to multiple f -divergences, rather than the standard Jensen-Shannon divergence used as the standard loss in GANs. This has led to a family of other GAIL variants that use other f -divergences [Ke+19a; Fin+16; Bro+20c], including f -GAIL that aims to also learn the best f -divergence to use [Zha+20e], as well as new analytical insight into the computation and generalization of such approaches [Che+20b].

1 **26.7.6 Domain adaptation**

2 An important task in machine learning is to correct for shifts in the data distribution over time,
3 minimizing some measure of domain shift, as we discuss in Section 19.5.3. Like with the other
4 applications, GANs are popular as ways of avoiding the choice of distance or degree of shift. Both
5 the supervised and unsupervised approaches for image generation we reviewed earlier looked at pixel-
6 level domain adaptation models that perform distribution alignment in raw pixel space, translating
7 source data to the style of a target domain, as with pix2pix and CycleGAN. Extensions of these
8 approaches for the general problem of domain adaptation seek to do this not only in the observed
9 data space (e.g., with pixels), but also at the feature level. One general approach is domain-
10 adversarial training of neural networks [Gan+16b] or adversarial discriminative domain adaptation
11 (ADDA) [Tze+17]. The CyCADA approach of [Hof+18] extends CycleGAN by enforcing both
12 structural and semantic consistency during adaptation using a cycle-consistency loss and semantics
13 losses based on a particular visual recognition task. There are also many extensions that include
14 class conditional information [Tsa+18; Lon+18] or adaptation when the modes to be matched have
15 different frequencies in the source and target domains [BHC19].

16

17 **26.7.7 Design, art and creativity**

18 Generative models, particularly of images, have added to approaches in the more general area of
19 algorithmic art. The applications in image and audio generation with transfer can also be considered
20 aspects of artistic image generation. In these cases, the goal of training is not generalization, but to
21 create appealing images across different types of visual aesthetics [Sar18]. One example takes style
22 transfer GANs to create visual experiences, in which objects placed under a video are re-rendered
23 using other visual styles in real time [AFG19]. The generation ability has been used to explore
24 alternative designs and fabrics in fashion [Kat+19], and have now also become part of major drawing
25 software to provide new tools to support designers [Ado]. And beyond images, creative and artistic
26 expression using GANs include areas in music, voice, dance, and typography [AI 19].

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

PART V

Discovery