

10 Generative Models in Healthcare

So far, we have learned about basic neural networks, namely DNN, RNN, CNN, and GNN. These networks are designed to extract feature information from various data types like structured EHR data, clinical notes, medical images, and physiological signals. They predict targets of interest such as disease diagnosis or medication recommendation. These models are supervised, and one needs sufficient annotated data (i.e., labeled data) for training. The output of these models is usually simple such as binary labels of disease presence/absence or numerical values like length of hospital stay.

In this chapter, we will learn about **generative models**, a new type of neural network that generates realistic data from scratch. These models can produce new images, videos, and sequences. The output of generative models is more complex (i.e., a high dimensional data sample) compared to supervised models (i.e., a scalar). After model training, generative models can create realistic new data based on random input or predefined input labels. The training process usually doesn't require annotated labeled data but utilizes cleverly constructed supervision from the data alone.

Generative models have the potential to revolutionize healthcare in various ways, such as personalized medicine [116], medical imaging [117], drug discovery [118], clinical trials [119, 120], rare diseases [121], genomics [122], electronic health records (EHR) [123], EEG [124] and ECG [125] analysis. For example, medical datasets are often sensitive due to privacy and legal constraints and cannot be made public. With generative models, it is possible to publish synthetic but realistic data that have similar distributions as private data, while still meeting privacy and legal compliance [123]. The newly generated data can also serve to augment the original datasets, thereby enriching the volume and diversity of the original datasets.

In recent years, significant advancements have been made in generative models in the fields of vision and natural language processing. OpenAI released a powerful question-answering dialog system called GPT-3.5¹ in late 2022, followed by GPT-4 in 2023. These large language models (LLMs) have revolutionized research in these fields, as well as many industries. By querying the LLM interface, one can get decent and human-like text or code answers. Since then, many open-sourced LLM pre-trained models have been released, such as Llama 2² from Meta, which can be fine-tuned for various downstream language or vision-language tasks. In addition, diffusion models

¹ <https://platform.openai.com/docs/models>, GPT stands for Generated Pretrained Transformer

² <https://llama.meta.com/llama2>

have impressed many people in the field of image generation. These models translate prompt text into images with various styles, such as Stable Diffusion 2.1³ from Stability and Midjourney⁴. The generated images or videos are difficult to distinguish from real ones. Recently, OpenAI’s popular Sora⁵ can even generate 60-second videos from text prompts without post-hoc editing.

Now, let’s learn about generative models. This chapter will introduce three generative models: Generative Adversarial Networks (GAN), developed by Goodfellow et al. [126, 127], Variational Auto-Encoder (VAE), proposed by Kingma et al. [128], and recently popular Diffusion Models [129, 130]. Towards the end of the chapter, we will demonstrate how to use PyHealth to implement GAN, VAE, and latent diffusion models (LDM) to generate synthetic chest X-ray images.

10.1 Variational Auto-encoder (VAE)

We introduce Variational Auto-encoder (VAE) in this section. As the background, let us start with the auto-encoder (AE) models proposed by Hinton and Salakhutdinov [131] in 2006. AE models use the *bottleneck layer* between encoder and decoder to learn compressed low-dimensional features for complex high-dimensional data. Note that AE models are not generative models but a nonlinear generalization of principal component analysis (PCA).

Motivated from AE, variational auto-encoder (VAE), proposed by Kingma et al [128] in 2013, follows a similar AE architecture which further imposes distribution constraints for the *bottleneck layer*. Different from AE and other AE variants, VAE model can generate new data that are similar to the original training data.

10.1.1 Simple auto-encoder (AE)

Auto-encoder (AE) is an unsupervised low-dimensional feature learning model. Like principle component analysis (PCA) or matrix factorization (MF), AE learns to compress complex data into a low-dimensional representation with non-linear functions (while PCA and MF use linear transformation). For example, given thousands of chest X-ray images, the AE model could learn to obtain meaningful image vector representations without label annotations. It mainly optimizes the reconstruction loss.

Encoder-decoder architecture

In Figure 10.1, the AE model consists of two main components: the *encoder* and the *decoder*. The encoder receives the input data (in this example, a chest X-ray image) and produces a low-dimensional embedding (also known as the bottleneck layer). The decoder takes this low-dimensional embedding and attempts to reconstruct the original

³ <https://huggingface.co/spaces/stabilityai/stable-diffusion>

⁴ <https://docs.midjourney.com/>

⁵ <https://openai.com/sora>

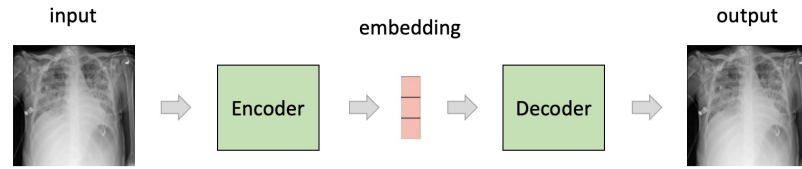


Figure 10.1 Illustration of an auto-encoder (AE) model

image. The entire model is trained using either the mean square error (MSE) or KL-divergence between the original and reconstructed images. Figure 10.2 provides an illustration of the process.

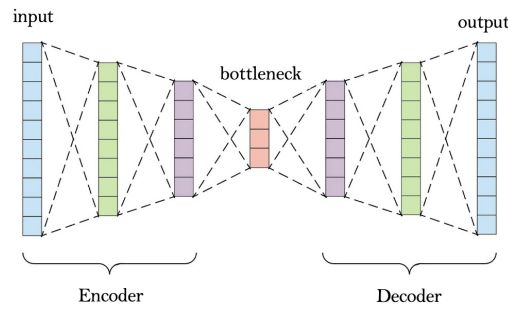


Figure 10.2 Illustration of auto-encoder (AE) architecture

Let's define the following notations: \mathbf{X} represents the original input data, $\text{Encoder}(\cdot)$ is the encoder model, \mathbf{h} is the bottleneck layer, $\text{Decoder}(\cdot)$ is the decoder model, and $\tilde{\mathbf{X}}$ is the reconstructed image. The encoder model can be any type of model, such as a convolutional neural network for images. It takes the original data as input and generates an embedding \mathbf{h} , such as a 256-dimensional vector. In other applications, like modeling electronic health record (EHR) data, the encoder can be an RNN model that models a sequence of patient visits. Formally, the encoder can be represented as:

$$\mathbf{h} = \text{Encoder}(\mathbf{X}). \quad (10.1)$$

The decoder model takes the bottleneck embedding \mathbf{h} as input and generates a reconstructed data sample. In this example for X-ray data, the decoder uses transposed convolutions to perform the opposite operation of the encoder, such as `torch.nn.ConvTranspose2d`⁶. Starting from the bottleneck embedding \mathbf{h} , it gradually upscales the representation back to the original input size. Transposed convolutional layers help in this process by increasing the spatial resolution of their input feature maps, effectively "reversing" the effect of convolutional layers used in the encoder.

⁶ <https://pytorch.org/docs/stable/generated/torch.nn.ConvTranspose2d.html>

$$\tilde{\mathbf{X}} = \text{Decoder}(\mathbf{h}). \quad (10.2)$$

Note that, in AE models, the encoder or the decoder can have multiple layers. The reconstruction is the MSE loss.

$$\mathcal{L} = \|\mathbf{X} - \tilde{\mathbf{X}}\|_F^2. \quad (10.3)$$

```

1  # Assume x is a chest X-ray image
2
3  class AE(nn.Module):
4      def __init__(self):
5          super(AE, self).__init__()
6
7          # two-layer DNN as encoder
8          self.encoder1 = nn.Linear(28*28, 128)
9          self.encoder2 = nn.Linear(128, 32)
10
11         # two-layer DNN as decoder
12         self.decoder1 = nn.Linear(32, 128)
13         self.decoder2 = nn.Linear(128, 28*28)
14
15     def encoder(self, x):
16         h = self.encoder2(torch.relu(self.encoder1(x)))
17         return h
18
19     def decoder(self, h):
20         x_hat = self.decoder2(torch.relu(self.decoder1(h)))
21         return x_hat
22
23     def forward(self, x):
24         bottleneck = self.encoder(x)
25         x_hat = self.decoder(torch.relu(bottleneck))
26         # the input, reconstructed images, and the bottleneck
27         return x, x_hat, bottleneck

```

In the code, both the encoder and the decoder use a two-layer deep neural network with ReLU as the activation.

Downstream Applications: Afterward, the encoder model can be used as pre-trained encoders in some downstream tasks, such as building a classification layer on top of the encoder for supervised learning tasks.

10.1.2 Variants of AE models

Simple AE models can easily overfit the training data, especially on small data sets. Several follow-up works are proposed to address the overfitting issue in AE. We briefly introduce some of those AE variants.

Denoising autoencoder

Vincent et al. proposed the de-noising auto-encoder (DAE) in 2008 [132]. The DAE model works by adding random noise to the input data and then requiring the auto-

encoder model to reconstruct the original data samples. This approach is based on the observation that humans are able to recognize objects or scenes even when they are partially occluded or corrupted [133]. The DAE architecture is similar to that of the AE architecture, with the only difference being that the input data is the noise-corrupted sample and the model still aims to reconstruct the original data. Some common noises include Gaussian noise and salt-and-pepper noise (i.e., setting a small percentage of pixels to minimum or maximum values). Check out Figure 10.3 for a visual representation of the DAE architecture. In healthcare, DAEs have been used to model patient records to create a latent patient representation [134].

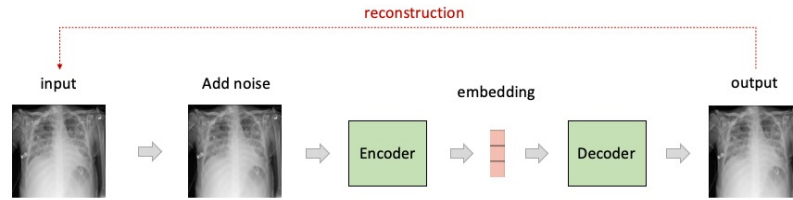


Figure 10.3 Illustration of de-noising auto-encoder (DAE) architecture

Sparse autoencoder

Adding "sparse" constraint on the hidden activation could also mitigate the AE overfitting problem and improve model robustness. Makhzani et al. [135] proposed the sparse auto-encoder (SAE), which forces the model to only keep k largest hidden units at the same time and set all other hidden units to zero. Mathematically, the SAE optimizes the following objective function:

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n L(\mathbf{X}^{(i)}, \hat{\mathbf{X}}^{(i)}) + \Omega(\mathbf{h}^{(i)})$$

where L is the reconstruction loss, Ω is a sparsity penalty term, and $\mathbf{h}^{(i)}$ is the hidden representation of the i -th sample. The sparsity penalty Ω encourages the hidden units to be sparse, i.e., only a small number of units are active for each input sample.

10.1.3 Variational Auto-encoder (VAE)

The autoencoder (AE) architecture, while powerful for learning compact representations of complex data, cannot generate new samples from the learned latent distribution. This limitation is addressed by the variational autoencoder (VAE), introduced by [128] in 2014, which builds upon the AE structure but introduces randomness and additional distribution constraints on the middle bottleneck layer.

Specifically, VAE enforces the bottleneck layer to follow a predefined standard distribution, such as the *multi-variate Gaussian distribution*. By doing so, VAE enables

the generation of new samples by sampling noise from the Gaussian distribution and using the decoder to map these noise to the data space.

Assumption in bottleneck layer distribution

For a d -dimensional bottleneck layer, usually a multi-variate Gaussian distribution needs d parameters for the mean vector μ and $d \times d$ parameters for the co-variance matrix Σ . In VAE, the authors assume each dimension in the multi-variate Gaussian distribution be independent, which leads to a diagonal co-variance matrix, thereby reducing its number of parameters from $d \times d$ to d diagonal elements. This makes the model more efficient and easier to train.

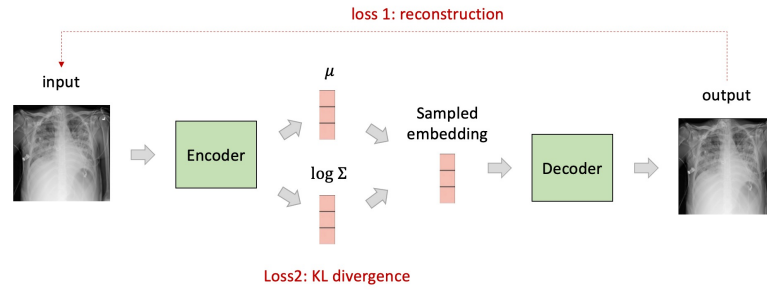


Figure 10.4 Illustration of variational auto-encoder (VAE) architecture. It has two objective functions: one is the MSE reconstruction loss between the original figures and the reconstructed figures, and another one is the KL-divergence between the standard normal distribution and the parameterized normal distribution.

An illustration of VAE is shown in Figure 10.4. We elaborate its forward propagation process below and explain the loss function later.

- **Step 1:** Given the high-dimensional input \mathbf{X} , the encoder of VAE first maps it to a mean vector $\mu \in \mathbb{R}^d$ and a co-variance diagonal vector $\log \Sigma \in \mathbb{R}^d$. Different from AE that directly map \mathbf{X} to the bottleneck layer, VAE essentially builds a parameterized multi-variate normal distribution at this step. Here, $\log \Sigma$ is a vector that contains the log-transformed diagonal elements from Σ (since diagonals of Σ is non-negative, we use log transformation).

$$\mu = \text{Encoder}_{\mu}(\mathbf{X}), \quad (10.4)$$

$$\log \Sigma = \text{Encoder}_{\Sigma}(\mathbf{X}). \quad (10.5)$$

- **Step 2:** Obtaining μ and $\log \Sigma$, we first apply a $\exp()$ on $\log \Sigma$ to restore the co-variance as Σ . Next, we want to get a sample from this parameterized distribution (μ, Σ) . However, directly sampling from there is in-differentiable and can break the loss back-propagation. The authors introduce a **re-parameterization trick** that samples a d -dimensional vector ϵ , where each element is independently sampled from a standard normal distribution. Then, re-scaling and translating ϵ

would generate a bottleneck embedding $\mathbf{h} \in \mathbb{R}^d$ that follows the parameterized distribution (μ, Σ) , and importantly \mathbf{h} is differentiable.

$$\epsilon_i \sim \mathcal{N}(0, 1), \quad i \in [0, 1, \dots, d-1], \quad (10.6)$$

$$\mathbf{h} = \mu + \epsilon \odot \exp(\log \Sigma)^{\frac{1}{2}} = \mu + \epsilon \odot \Sigma^{\frac{1}{2}}. \quad (10.7)$$

Here, \odot is the vector element-wise product.

- **Step 3:** Afterwards, the model uses a decoder to reconstruct the sample \mathbf{X} .

$$\tilde{\mathbf{X}} = \text{Decoder}(\mathbf{h}). \quad (10.8)$$

Objective functions

Now, we explain the loss function below, which contains two parts: one is the reconstruction loss, and another part is the distribution gap between (μ, Σ) and standard multi-variate Gaussian (where dimensions are all independent).

Optimizing the reconstruction loss. Similar to the AE model, reconstruction loss is defined by the MSE loss.

$$\mathcal{L}_1 = \|\mathbf{X} - \tilde{\mathbf{X}}\|_F^2. \quad (10.9)$$

Optimize KL-divergence loss. The new constraints on VAE is that on the middle bottleneck layer \mathbf{h} , we want the parameterized distribution follows standard Gaussian and use Kullback–Leibler (KL) divergence⁷ to measure the distribution gap.

Specifically, we denote the standard Gaussian as $p(\mathbf{h})$, which can be written down as

$$\mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix}\right). \quad (10.10)$$

Its probability density function (PDF) is

$$p(\mathbf{h}) = \frac{1}{(2\pi)^{\frac{d}{2}}} \exp\left(-\frac{\mathbf{h}^\top \mathbf{h}}{2}\right). \quad (10.11)$$

Our parameterized Gaussian $\mathcal{N}(\mu, \Sigma)$ is denoted as $q(\mathbf{h} \mid \mu, \Sigma)$, and here μ and Σ contain parameters. The PDF of this distribution could be written down as

$$q(\mathbf{h} \mid \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{d}{2}} \det(\Sigma)^{\frac{1}{2}}} \exp\left(-\frac{(\mathbf{h} - \mu)^\top \Sigma^{-1} (\mathbf{h} - \mu)}{2}\right). \quad (10.12)$$

⁷ https://en.wikipedia.org/wiki/Kullback%E2%80%93Leibler_divergence

The KL-divergence loss can be derived analytically,

$$\int_{\mathbf{h} \sim q} \ln \frac{p(\mathbf{h})}{q(\mathbf{h} | \mu, \Sigma)} d\mathbf{h} = \int q(\mathbf{h} | \mu, \Sigma) (\ln p(\mathbf{h}) - \ln q(\mathbf{h} | \mu, \Sigma)) d\mathbf{h} \quad (10.13)$$

$$= \underbrace{\int q(\mathbf{h} | \mu, \Sigma) \ln p(\mathbf{h}) d\mathbf{h}}_{\text{Part I}} - \underbrace{\int q(\mathbf{h} | \mu, \Sigma) \ln q(\mathbf{h} | \mu, \Sigma) d\mathbf{h}}_{\text{Part II}} \quad (10.14)$$

Here, we decompose the integral into two parts. Also, we know that $\mathbf{h} = \mu + \epsilon \odot \Sigma^{\frac{1}{2}}$. For Part I,

$$\int q(\mathbf{h} | \mu, \Sigma) \ln p(\mathbf{h}) d\mathbf{h} \quad (10.15)$$

$$= \int \frac{1}{(2\pi)^{\frac{d}{2}} \det(\Sigma)^{\frac{1}{2}}} \exp\left(-\frac{(\mathbf{h} - \mu)^{\top} \Sigma^{-1} (\mathbf{h} - \mu)}{2}\right) \ln \frac{1}{(2\pi)^{\frac{d}{2}}} \exp\left(-\frac{\mathbf{h}^{\top} \mathbf{h}}{2}\right) d\mathbf{h} \quad (10.16)$$

$$= \int \frac{1}{(2\pi)^{\frac{d}{2}} \det(\Sigma)^{\frac{1}{2}}} \exp\left(-\frac{(\mathbf{h} - \mu)^{\top} \Sigma^{-1} (\mathbf{h} - \mu)}{2}\right) \left(-\frac{d}{2} \ln(2\pi) - \frac{\mathbf{h}^{\top} \mathbf{h}}{2}\right) d\mathbf{h} \quad (10.17)$$

$$= -\frac{d}{2} \ln(2\pi) \int \frac{1}{(2\pi)^{\frac{d}{2}} \det(\Sigma)^{\frac{1}{2}}} \exp\left(-\frac{(\mathbf{h} - \mu)^{\top} \Sigma^{-1} (\mathbf{h} - \mu)}{2}\right) d\mathbf{h} \quad (10.18)$$

$$- \int \frac{\mathbf{h}^{\top} \mathbf{h}}{2} \frac{1}{(2\pi)^{\frac{d}{2}} \det(\Sigma)^{\frac{1}{2}}} \exp\left(-\frac{(\mathbf{h} - \mu)^{\top} \Sigma^{-1} (\mathbf{h} - \mu)}{2}\right) d\mathbf{h} \quad (10.19)$$

$$= -\frac{d}{2} \ln(2\pi) - (\mu^{\top} \mu + \|\Sigma\|_F^2). \quad (10.20)$$

For Part II, we have

$$\int q(\mathbf{h} | \mu, \Sigma) \ln q(\mathbf{h} | \mu, \Sigma) d\mathbf{h} \quad (10.21)$$

$$= \int \frac{1}{(2\pi)^{\frac{d}{2}} \det(\Sigma)^{\frac{1}{2}}} \exp\left(-\frac{(\mathbf{h} - \mu)^{\top} \Sigma^{-1} (\mathbf{h} - \mu)}{2}\right) \ln \frac{1}{(2\pi)^{\frac{d}{2}} \det(\Sigma)^{\frac{1}{2}}} \exp\left(-\frac{(\mathbf{h} - \mu)^{\top} \Sigma^{-1} (\mathbf{h} - \mu)}{2}\right) d\mathbf{h} \quad (10.22)$$

$$= \int \frac{1}{(2\pi)^{\frac{d}{2}} \det(\Sigma)^{\frac{1}{2}}} \exp\left(-\frac{(\mathbf{h} - \mu)^{\top} \Sigma^{-1} (\mathbf{h} - \mu)}{2}\right) \left(-\frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln(\det(\Sigma)) - \frac{(\mathbf{h} - \mu)^{\top} \Sigma^{-1} (\mathbf{h} - \mu)}{2}\right) d\mathbf{h} \quad (10.23)$$

$$= -\left(\frac{d}{2} \ln(2\pi) + \frac{1}{2} \ln(\det(\Sigma))\right) \int \frac{1}{(2\pi)^{\frac{d}{2}} \det(\Sigma)^{\frac{1}{2}}} \exp\left(-\frac{(\mathbf{h} - \mu)^{\top} \Sigma^{-1} (\mathbf{h} - \mu)}{2}\right) d\mathbf{h} \quad (10.24)$$

$$- \int \frac{(\mathbf{h} - \mu)^{\top} \Sigma^{-1} (\mathbf{h} - \mu)}{2} \frac{1}{(2\pi)^{\frac{d}{2}} \det(\Sigma)^{\frac{1}{2}}} \exp\left(-\frac{(\mathbf{h} - \mu)^{\top} \Sigma^{-1} (\mathbf{h} - \mu)}{2}\right) d\mathbf{h} \quad (10.25)$$

$$= -\frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln(\det(\Sigma)) - \frac{d}{2}. \quad (10.26)$$

Thus, the KL divergence is calculated by subtracting Part II from Part I, which

yields,

$$\int_{\mathbf{h} \sim q} \ln \frac{p(\mathbf{h})}{q(\mathbf{h} | \mu, \Sigma)} d\mathbf{h} = \frac{d}{2} + \frac{1}{2} \ln \det(\Sigma) - \mu^\top \mu - \|\Sigma\|_F^2. \quad (10.27)$$

We further add the reconstruction loss and the KL divergence together to be the final VAE loss function. Now, we will learn the PyTorch implementation of VAE on the MNIST dataset below. At the end of this chapter, we will learn to apply VAE on the patient chest X-ray dataset.

10.1.4 PyTorch Implementation of VAE

Step 1: Loading the MNIST dataset

We first load the training and test data from online MNIST sources. We could see that there are 60k training data and 10k test data samples, and each is a digit image of 28×28 size. An illustration of the first five training images of MNIST is shown in Figure 10.5.

```

1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4 import torchvision
5 import cv2
6 import numpy as np
7
8 # we will use the first GPU on the machine for training
9 # user could also change into "cpu" if no GPU is available
10 device = torch.device('cuda:0')
11
12 # for training set, we load the data and label
13 trainset = torchvision.datasets.MNIST(\
14     root='./data', train = True, download = True)
15 xtrain = trainset.data.numpy()
16 ytrain = trainset.targets.numpy()
17
18 # for test set, we load the data and label
19 testset = torchvision.datasets.MNIST(
20     root='./data', train = False, download = True)
21 xtest = testset.data.numpy()
22 ytest = testset.targets.numpy()
23
24 print (xtrain.shape, ytrain.shape)
25 print (xtest.shape, ytest.shape)
26 """
27 (60000, 28, 28) (60000,)
28 (10000, 28, 28) (10000,)
29 """

```

We then wrap them up as the data loaders. To make life easier, we transform the pixel range from (0, 255) to binary {0, 1} by setting a threshold 128.

```

1 # transform pixel value (0, 255) into {0,1}
2 xtrain = np.where(xtrain > 128, 1, 0)
3 xtrain = xtrain.astype(np.float32)

```

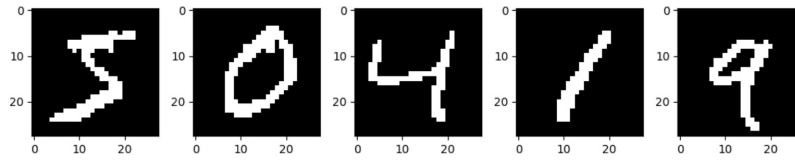


Figure 10.5 Illustration of the first 5 MNIST training samples

```

4
5 # transform pixel value (0, 255) into {0,1}
6 xtest = np.where(xtest > 128, 1, 0)
7 xtest = xtest.astype(np.float32)
8
9 batch_size = 256
10 trainloader = torch.utils.data.DataLoader(
11     [[xtrain[i], ytrain[i]] for i in range(len(ytrain))], \
12     shuffle=True, batch_size=batch_size)
13 testloader = torch.utils.data.DataLoader(
14     [[xtest[i], ytest[i]] for i in range(len(ytest))], \
15     shuffle=False, batch_size=batch_size)

```

Step 2: build the VAE model

We then build the VAE model, which has an encoder and a decoder. The encoder is a two-layer neural network with tanh as the activation, and it has two heads, `self.mu` is used for predicting the μ vector, and `self.log_sigma2` is used for predicting the $\ln \Sigma$ vector. The decoder is also a two-layer neural network with tanh as the activation while we finally use Sigmoid to re-scale the output. Since the MNIST dataset has all 0 or 1 as pixel values, we use binary cross entropy loss instead of MSE loss. Here, $d = 8$ is used as the bottleneck dimension.

```

1 class VAE(nn.Module):
2     def __init__(self):
3         super(VAE, self).__init__()
4
5         # DNN as encoder
6         self.encoder1 = nn.Linear(28*28, 128)
7         self.mu = nn.Linear(128, 8)
8         self.log_sigma2 = nn.Linear(128, 8)
9
10        # DNN as decoder
11        self.decoder1 = nn.Linear(8, 128)
12        self.decoder2 = nn.Linear(128, 28*28)
13
14        def encoder(self, x):
15            h = torch.tanh(self.encoder1(x))
16            return self.mu(h), torch.sqrt(torch.exp(self.log_sigma2(h)))
17
18        def sampling(self, mu, std): # Reparameterization trick
19            eps = torch.randn_like(std)
20            return mu + eps * std
21

```

```

22     def decoder(self, z):
23         h = torch.tanh(self.decoder1(z))
24         return torch.sigmoid(self.decoder2(h))
25
26     def forward(self, x):
27         mu, std = self.encoder(x.view(-1, 28*28))
28         z = self.sampling(mu, std)
29         return self.decoder(z), mu, std
30
31     @staticmethod
32     def loss_func(x_hat, x, mu, std):
33         # define the reconstruction loss
34         ERR = F.binary_cross_entropy(x_hat, x, reduction='sum')
35         # define the KL divergence loss
36         KLD = -0.5 * torch.sum(1 + torch.log(std**2) - mu**2 - std**2)
37         return ERR + KLD, ERR, KLD
38
39 model = VAE()
40 model.to(device) # move model to GPU

```

In the above scripts, the VAE model has an encoder layer `self.encoder1()` with the `tanh()` activation that maps the input images into hidden embedding **h**. From **h**, the `self.h()` function will generate the mean of Gaussian distribution while the `self.log_sigma2()` will generate $\log \Sigma$ (we will use `exp()` to restore the standard deviation Σ of the Gaussian distribution). Let us look at the `self.forward()` function, after obtaining the mean (μ) and standard deviation (Σ), we use the `self.sampling()` method to obtain a differentiable sample from the Gaussian distribution by $\mu + \epsilon \odot \Sigma^{\frac{1}{2}}$ and decode the it into a reconstructed image.

Step 3: Model Training and Evaluation

We use the standard training pipeline and evaluate the model on test data after each training epoch. We could observe that both the training and test error (i.e., reconstruction error).

Ideally the KL divergence should increase as well but it goes down gradually, which is a common issue in practical implementation of VAE model. This issue could be addressed by carefully tuning the hyperparameters or using learning rate scheduling as discussed in a Microsoft Research (MSR) blog ⁸. Later on chest X-ray application, we will show that the latent diffusion model (LDM) could fix the problem and improve the generation quality.

```

1 optimizer = torch.optim.AdamW(model.parameters(), lr = 0.002)
2
3 rec_loss, kl_loss = [], []
4 for epoch in range(50):
5
6     """ model training """
7     model.train()
8     cur_rec_loss, cur_kl_loss = [], []

```

⁸ <https://www.microsoft.com/en-us/research/blog/less-pain-more-gain-a-simple-method-for-vae-training-with-less-of-that-kl-vanishing-agony/>

```

9     for batch_idx, (data, _) in enumerate(trainloader):
10         data = data.to(device)
11         # get reconstruction, mean and std vectors
12         rec, mu, std = model(data)
13         # calculate the loss, reconstruction error, and KL divergence
14         loss, err, kl = model.loss_func(rec, data.reshape(-1, 28*28), mu, std)
15
16         # loss backprop
17         optimizer.zero_grad()
18         loss.backward()
19         optimizer.step()
20
21         cur_rec_loss.append(err.item())
22         cur_kl_loss.append(kl.item())
23
24     rec_loss.append(np.mean(cur_rec_loss))
25     kl_loss.append(np.mean(cur_kl_loss))
26
27     """ model evaluation """
28     with torch.no_grad():
29         test_loss = []
30         for batch_idx, (data, _) in enumerate(testloader):
31             data = data.to(device)
32             rec, mu, std = model(data)
33             _, loss, _ = model.loss_func(rec, data.reshape(data.shape[0], -1),
34                                         mu, std)
35
36             test_loss.append(loss.item())
37
38     if epoch % 10 == 0:
39         print (f"-- epoch {epoch} --, train loss: {np.mean(cur_rec_loss)}, \
40               train KL: {np.mean(cur_kl_loss)}, test loss: {np.mean(test_loss)}")
41
42     """
43     -- epoch 0 --, train loss: 45738.48, train KL: 2987.92, test loss: 33421.22
44     -- epoch 10 --, train loss: 24595.57, train KL: 4055.006, test loss: 24020.88
45     -- epoch 20 --, train loss: 23028.66, train KL: 4170.09, test loss: 22486.60
46     -- epoch 30 --, train loss: 22330.19, train KL: 4242.63, test loss: 21889.96
47     -- epoch 40 --, train loss: 21925.85, train KL: 4286.16, test loss: 21725.06
48     """

```

Step 4: Generating new digit images

After training the model, we can generate synthetic digit images. One could sample 8-dimensional Gaussian noise and inputs them into the VAE decoder for obtaining the synthesized images. Figure 10.6 shows 10 images generated by VAE.

```

1 plt.figure(figsize=(10, 3))
2
3 model.eval()
4 with torch.no_grad():
5     for l in range(10):
6         plt.subplot(2, 5, l+1)
7
8         # randomly generate a 8-dim bottleneck vector
9         x = np.random.normal(0, 1, 8)
10        x = x.astype(np.float32)

```

```

11         x = torch.from_numpy(x).to(device).unsqueeze(0)
12
13         # use decoder to generate the digit image
14         rec = model.decoder(x).detach().cpu().numpy()
15         rec = rec.reshape((1, 28, 28))[0]
16         plt.imshow(rec, cmap="gray")
17
18     plt.tight_layout()
19     plt.show()

```

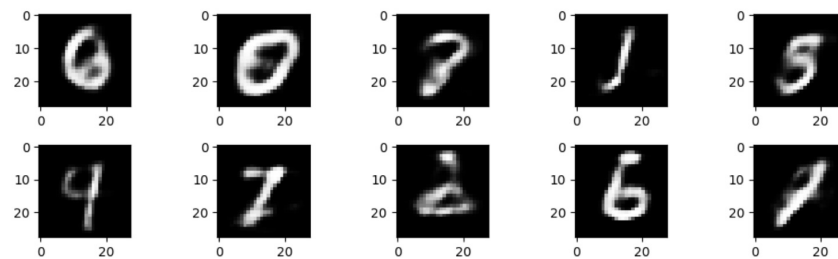


Figure 10.6 Ten digit images generated by VAE

Step 5: Conditional VAE

The VAE model can generate new digit images from random Gaussian vectors, but currently we can't control which digit is generated. To controllably generate class-specific digit images, we could build a ConditionalVAE model by adding the class label information as part of the input. The ConditionalVAE model is similar to the VAE model, and the only difference is in the `self.decoder` part. We additionally learn embedding vectors for each label (0-9), stored in `self.digit`, and this label information is added to the bottleneck layer z before feeding it into the decoder model.

```

1  class ConditionalVAE(nn.Module):
2      def __init__(self):
3          super(ConditionalVAE, self).__init__()
4
5          # for each digit, we learn an 8-dim vector, initialized randomly
6          # they are learned together with the encoder and decoder model
7          self.digit = nn.Embedding(10, 8)
8
9          # DNN as encoder
10         self.encoder1 = nn.Linear(28*28, 128)
11         self.mu = nn.Linear(128, 8)
12         self.log_sigma2 = nn.Linear(128, 8)
13
14         # DNN as decoder
15         self.decoder1 = nn.Linear(8, 128)
16         self.decoder2 = nn.Linear(128, 28*28)
17
18     def encoder(self, x):
19         h = torch.tanh(self.encoder1(x))

```

```

20         return self.mu(h), torch.sqrt(torch.exp(self.log_sigma2(h)))
21
22     def sampling(self, mu, std): # Reparameterization trick
23         eps = torch.randn_like(std)
24         return mu + eps * std
25
26     def decoder(self, z, l):
27         # add label embedding to random vector
28         z = z + self.digit(l)
29         h = torch.tanh(self.decoder1(z))
30         return torch.sigmoid(self.decoder2(h))
31
32     def forward(self, x, l):
33         mu, std = self.encoder(x.view(-1, 28*28))
34         z = self.sampling(mu, std)
35         return self.decoder(z, l), mu, std
36
37     @staticmethod
38     def loss_func(x_hat, x, mu, std):
39         # define the reconstruction loss
40         ERR = F.binary_cross_entropy(x_hat, x, reduction='sum')
41         # define the KL divergence loss
42         KLD = -0.5 * torch.sum(1 + torch.log(std**2) - mu**2 - std**2)
43         return ERR + KLD, ERR, KLD
44
45 model = ConditionalVAE()
46 model.to(device) # move model to GPU

```

Figure 10.7 shows an illustration of five synthetic "0" and five synthetic "1" generated using this method.

```

1 plt.figure(figsize=(10, 3))
2
3 model.eval()
4 with torch.no_grad():
5     for idx, l in enumerate([0 for i in range(5)] + [1 for i in range(5)]):
6         plt.subplot(2, 5, idx+1)
7
8         # randomly generate a 8-dim bottleneck vector
9         x = np.random.normal(0, 1, 8)
10        x = x.astype(np.float32)
11        x = torch.from_numpy(x).to(device).unsqueeze(0)
12        l = torch.LongTensor([l]).to(device).unsqueeze(0)
13
14        # use decoder to generate the digit image
15        rec = model.decoder(x, l).detach().cpu().numpy()
16        rec = rec.reshape((1, 28, 28))[0]
17        plt.imshow(rec, cmap="gray")
18
19 plt.tight_layout()
20 plt.show()

```

A concrete notebook including the implementation of VAE and ConditionalVAE could be found in this public repository⁹.

⁹ <https://github.com/sunlabuiuc/pyhealth-book/tree/main/chap8-GenAI/notebook>

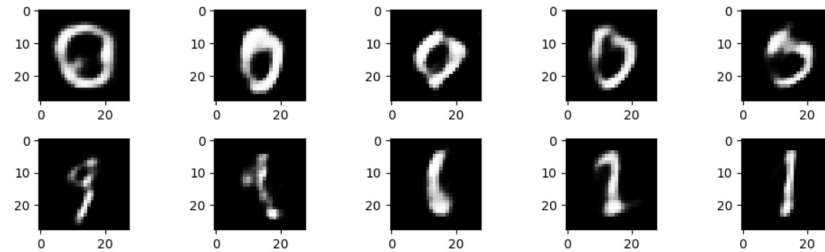


Figure 10.7 Five “0” and Five “1” digit images generated by ConditionalVAE

10.2 Generative Adversarial Network (GAN)

Generative adversarial networks (GAN) were proposed by Goodfellow [126, 127] as a new way of training synthetic data generators. The GAN model is composed of two parts: a generator (G) and a discriminator (D).

The generator (G) starts by taking a random noise from standard distributions (such as multi-variate Gaussian) and converting it into a "fake" data sample. The goal of the GAN model is to ensure that the generated samples look similar to the real samples. To achieve this, the discriminator (D), which is a supervised learning classifier, takes a data sample as input and determines whether it is real or fake. If the discriminator (D) is strong enough and cannot tell the difference between the real and fake samples, then the generator (G) is proven to produce highly realistic images.

Figure 10.8 depicts the GAN training process: the discriminator tries to differentiate between real and fake images while the generator tries to fool the discriminator by creating realistic fake images.

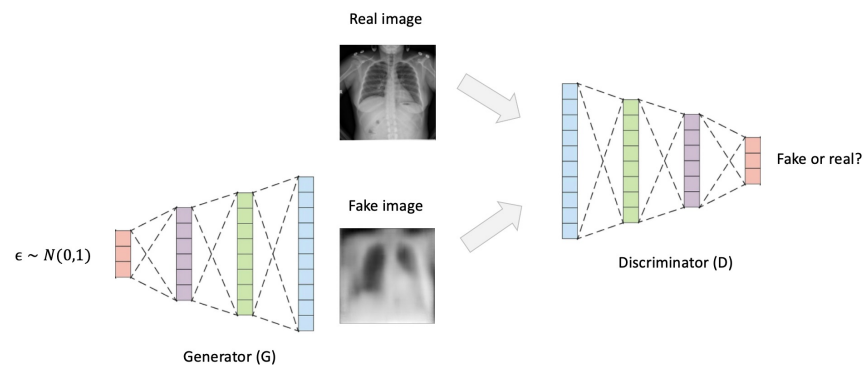


Figure 10.8 Illustration of generative adversarial network (GAN)

10.2.1 GAN training process within generator (G) and discriminator (D)

Formally, we assume the generator is $G_\theta(\cdot)$ with θ as the parameters and the discriminator is $D_\phi(\cdot)$ with ϕ as the parameters. The model works as follows,

- Generate a random sample ϵ from a standard distribution, such as 8-dimensional multi-variate Gaussian distribution with independent dimensions.

$$\epsilon \sim \mathcal{N}(0, \mathbf{I}). \quad (10.28)$$

- The generator (G) will encode the random noise into a fake sample \mathbf{X}_f , the same size as the real images.

$$\mathbf{X}_f = G_\theta(\epsilon). \quad (10.29)$$

- Collect a dataset with fake samples and real samples. Fake samples \mathbf{X}_f have label $y = 0$, and real samples \mathbf{X}_r have label $y = 1$. Use the discriminator $D(\cdot)$ to predict whether the data sample is fake or real. The output of discriminator is within $(0, 1)$ range since it passes through the Softmax function.

$$\hat{y} = D_\phi(\mathbf{X}) \in (0, 1). \quad (10.30)$$

- Calculate the binary cross entropy (BCE) loss $\mathcal{L}(y, \hat{y})$. To minimize this BCE loss, essentially, we will maximize $\log \hat{y} = \log D_\phi(\mathbf{X})$ for real samples and maximize $\log(1 - \hat{y}) = \log(1 - D_\phi(\mathbf{X}))$ for fake samples.

$$\mathcal{L}(y, \hat{y}) = -y \log \hat{y} - (1 - y) \log(1 - \hat{y}). \quad (10.31)$$

To train a good discriminator D , we maximize the following quantity, including the real image part and the fake image part.

$$\max_{\phi} \mathbb{E}_{\mathbf{X}_r \sim p_r} [\log D_\phi(\mathbf{X}_r)] + \mathbb{E}_{\mathbf{X}_f \sim \theta} [\log(1 - D_\phi(\mathbf{X}_f))], \quad (10.32)$$

where $\mathbf{X}_r \sim p_r$ means the real samples are obtained from true data distribution, and $\mathbf{X}_f \sim \theta$ means that the fake ones are generated from the generator with parameter θ . It can also be written as

$$\max_{\phi} \mathbb{E}_{\mathbf{X}_r \sim p_r} [\log D_\phi(\mathbf{X}_r)] + \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I})} [\log(1 - D_\phi(G_\theta(\epsilon)))]. \quad (10.33)$$

However, **for training a good generator G** , we want the fake image to look as real as possible, which means fixing the discriminator parameters ϕ , we want to minimize the following.

$$\min_{\theta} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I})} [\log(1 - D_\phi(G_\theta(\epsilon)))]. \quad (10.34)$$

Overall, the goal of the GAN training follows a min-max principle. At each iteration, we train the discriminator first (maximizing the following quantity over ϕ by fixing θ) and then train the generator (minimizing the same quantity over θ by fixing ϕ).

$$\min_{\theta} \max_{\phi} \mathbb{E}_{\mathbf{X}_r \sim p_r} [\log D_\phi(\mathbf{X}_r)] + \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I})} [\log(1 - D_\phi(G_\theta(\epsilon)))]. \quad (10.35)$$

To improve the model stability, researchers usually use another form, which has the same global optimal solutions but enjoys better numerical properties.

$$\min_{\theta} \max_{\phi} \mathbb{E}_{\mathbf{X}_r \sim p_r} [\log D_{\phi}(\mathbf{X}_r)] - \mathbb{E}_{\epsilon \sim \mathcal{N}(0,1)} [\log(D_{\phi}(G_{\theta}(\epsilon)))]. \quad (10.36)$$

10.2.2 PyTorch Implementation of GAN

We will learn how to implement a GAN model from scratch with PyTorch. In the code, we use a three-layer neural network for the discriminator and another three-layer neural network for the generator in GAN. The dataset will still be MNIST. We show the application on chest X-ray images at the end of the chapter.

Define the GAN model

We define the GAN model class below. In the model, the discriminator uses a ReLU as the activation function and the generator also uses ReLU as the activation function. The model has three functions: `self.discriminate()` function takes a 28×28 image as input and generates a probability between (0, 1) while 0 indicating “fake” and 1 indicating “real”. The `self.generate_fake()` function will sample `n_samples` random noise and use the generator to generate `n_samples` fake images.

```

1  class GAN(nn.Module):
2      def __init__(self):
3          super(GAN, self).__init__()
4
5          # three layer neural network as discriminator
6          self.discriminator = nn.Sequential(
7              nn.Linear(28*28, 256),
8              nn.LeakyReLU(negative_slope=0.01),
9              nn.Linear(256, 128),
10             nn.LeakyReLU(negative_slope=0.01),
11             nn.Linear(128, 1),
12             nn.Sigmoid()
13         )
14
15         self.generator = nn.Sequential(
16             nn.Linear(8, 128),
17             nn.LeakyReLU(negative_slope=0.01),
18             nn.Linear(128, 256),
19             nn.LeakyReLU(negative_slope=0.01),
20             nn.Linear(256, 28*28),
21         )
22
23     def discriminate(self, x):
24         y = self.discriminator(x)
25         return y
26
27     def generate_fake(self, n_samples, device):
28         # sample 8-dimensional vectors as random noise
29         eps = torch.randn(n_samples, 8).to(device)
30         # generate images based on the random vectors
31         fake_images = torch.sigmoid(self.generator(eps))
32         return fake_images

```

```

33 model = GAN()
34
35 model.to(device) # move model to GPU

```

GAN Model Training and Evaluation

Updating the GAN model is a two-step process with binary cross-entropy (BCE) loss: we first optimize the discriminator and then optimize the generator.

To begin, we get a real batch of MNIST samples from the training loader, called `real_imgs`. Then, we ask the generator to create fake images of the same size, called `fake_imgs`. We label the real images as "1" and the fake images as "0". Next, we ask the discriminator to calculate the classification loss, and minimizing this loss would improve the discriminator's performance. We use the optimizer `opt_D` to update the discriminator.

In the second step, we generate another random batch of fake images, named `fake_imgs` and label them as "1" since at this step, our goal is to let the discriminator believes the generated images are real. We ask the discriminator to calculate the BCE loss and then minimize it. This step helps us improve the generator based on the current discriminator. After loss back-propagation, the optimizer `opt_G` will then update the generator.

Our ultimate goal is to build a better discriminator that can distinguish between real and fake images and an even better generator that can fool this discriminator. The discriminator and generator are optimized alternately to achieve this goal in an adversarial way.

```

1  # the loss function
2  loss = torch.nn.BCELoss()
3
4  # optimizers
5  opt_G = torch.optim.AdamW(model.generator.parameters(), lr=2e-4)
6  opt_D = torch.optim.AdamW(model.discriminator.parameters(), lr=2e-4)
7
8  # store the loss curve of D and G
9  curve_D, curve_G = [], []
10 for epoch in range(50):
11
12     # store the loss over one epoch
13     curve_G.append(0)
14     curve_D.append(0)
15
16     for data, label in trainloader:
17         """ train discriminator """
18         opt_D.zero_grad()
19
20         batch_size = data.shape[0]
21         # get real and fake images
22         real_imgs = data.view(batch_size, -1).to(device)
23         fake_imgs = model.generate_fake(batch_size, device)
24         # set their labels (real gets ones, and fake gets zeros)
25         ones = torch.ones(batch_size, 1).to(device)
26         zeros = torch.zeros(batch_size, 1).to(device)

```

```

27
28     real_loss = loss(model.discriminator(real_imgs), ones)
29     fake_loss = loss(model.discriminator(fake_imgs), zeros)
30     loss_D = (real_loss + fake_loss) / 2
31
32     loss_D.backward()
33     opt_D.step()
34
35     """ train generator """
36     opt_G.zero_grad()
37     # get fake images and label it as ones
38     fake_imgs = model.generate_fake(batch_size, device)
39     loss_G = loss(model.discriminator(fake_imgs), ones)
40
41     loss_G.backward()
42     opt_G.step()
43
44     curve_G[-1] += loss_G.item()
45     curve_D[-1] += loss_D.item()
46
47     print (f"epoch: {epoch} --- \
48           loss of G: {curve_G[-1]}, loss of D: {curve_D[-1]}")
49
50     """
51     epoch: 0 --- loss of G: 746.3518468141556, loss of D: 43.85571137443185
52     epoch: 1 --- loss of G: 1275.6765661239624, loss of D: 5.175452238880098
53     epoch: 2 --- loss of G: 1571.3696827888489, loss of D: 4.034340920858085
54     epoch: 3 --- loss of G: 1540.1672163009644, loss of D: 1.1458109847735614
55     ...
56     epoch: 44 --- loss of G: 2625.360191345215, loss of D: 1.2749904468219029
57     epoch: 45 --- loss of G: 2535.6209321022034, loss of D: 1.6519437322713202
58     epoch: 46 --- loss of G: 2483.5613226890564, loss of D: 2.927129815390799
59     epoch: 47 --- loss of G: 2380.45946598053, loss of D: 2.1951384948333725
60     epoch: 48 --- loss of G: 2434.2216053009033, loss of D: 2.180281783337705
61     epoch: 49 --- loss of G: 2226.480115890503, loss of D: 2.546948796516517
62     """

```

Generating synthetic digit images with GAN

As demonstrated in Figure 10.9, similar to a VAE model, a trained GAN model can be utilized to generate synthetic digit images. However, in this particular instance, the quality of images produced by the GAN model is not as good as those generated by the VAE model. It's important to note that this is not a general statement about GAN, as there are numerous examples of successful GAN models that produce realistic data samples.

```

1 plt.figure(figsize=(10, 3))
2
3 model.eval()
4 with torch.no_grad():
5     for i in range(10):
6         plt.subplot(2, 5, i+1)
7
8         # use decoder to generate the digit image
9         rec = model.generate_fake(1, device).detach().cpu().numpy()

```

```

10     rec = rec.reshape((1, 28, 28))[0]
11     plt.imshow(rec, cmap="gray")
12
13 plt.tight_layout()
14 plt.show()

```

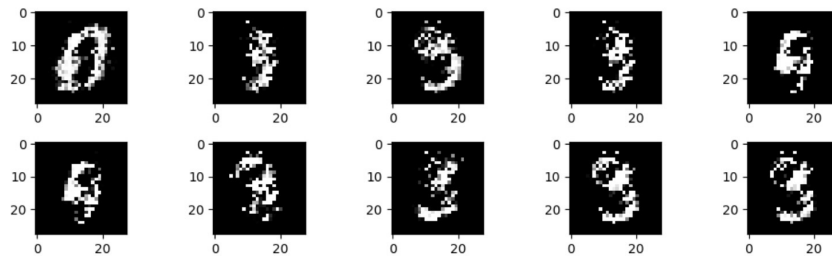


Figure 10.9 Ten “fake” digit images generated by GAN

Conditional GAN

We can also create a model called ConditionalGAN that generates images based on specific labels instead of random digits. The ConditionalGAN model differs from the simple GAN model in two ways: (i) the label embedding is added to the random noise ϵ before being fed into the generator, and (ii) instead of using binary cross-entropy loss, ConditionalGAN uses cross-entropy loss over various classes, including class 0 to class 9, and an additional class for fake images.

To train the discriminator, we assign real images to class labels and use the additional class for fake images. When training the generator, we set the desired class labels for the fake images and minimize the cross-entropy loss. You can find the implementation of the ConditionalGAN model below.

```

1  class ConditionalGAN(nn.Module):
2      def __init__(self):
3          super(ConditionalGAN, self).__init__()
4
5          # embedding for each class, initialized randomly
6          # it is learned together with the generator and discriminator
7          self.digit = nn.Embedding(10, 8)
8
9          # three layer neural network as discriminator
10         # output a 11-dim vector represents the probability of each class
11         # including the additional fake class
12         self.discriminator = nn.Sequential(
13             nn.Linear(28*28, 256),
14             nn.LeakyReLU(negative_slope=0.01),
15             nn.Linear(256, 128),
16             nn.LeakyReLU(negative_slope=0.01),
17             nn.Linear(128, 11),
18         )
19
20         # three layer neural network as generator

```

```

21         self.generator = nn.Sequential(
22             nn.Linear(8, 128),
23             nn.LeakyReLU(negative_slope=0.01),
24             nn.Linear(128, 256),
25             nn.LeakyReLU(negative_slope=0.01),
26             nn.Linear(256, 28*28),
27         )
28
29     def discriminate(self, x):
30         y = self.discriminator(x)
31         return y
32
33     def generate_fake(self, labels, device):
34         # generate random vectors
35         eps = torch.randn(len(labels), 8).to(device)
36         # get the class label embeddings
37         digit_emb = self.digit(labels.to(device))
38         # add these two and pass to generator
39         fake_images = torch.sigmoid(self.generator(eps + digit_emb))
40         return fake_images
41
42 model = ConditionalGAN()
43 model.to(device) # move model to GPU

```

We skip the training process of ConditionalGAN here as it is similar to training of GAN. Figure 10.10 shows the digit images from 0 to 9 generated by ConditionalGAN model, the quality of which seems better than GAN model.

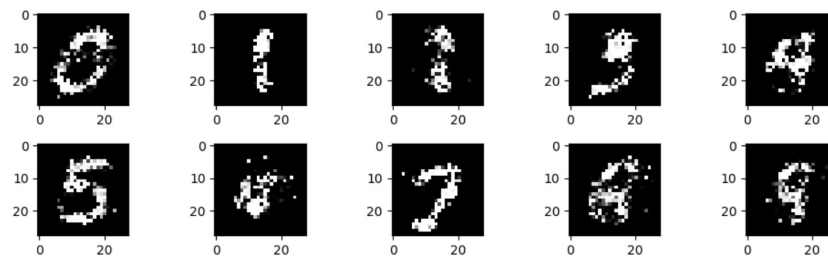


Figure 10.10 Ten digit images generated by ConditionalGAN

A concrete notebook including implementations of GAN and ConditionalGAN could be found in this public repository¹⁰.

Practical improvements of GAN models

The original training procedure of Generative Adversarial Networks (GANs) often runs into issues such as instability and mode collapse, where the generator produces a limited variety of outputs. To address these challenges, Arjovsky and his colleagues introduced a new type of Generative Adversarial Networks (GANs), called Wasserstein GAN (WGAN) [136]. The WGAN model replaces the discriminator with a critic that

¹⁰ <https://github.com/sunlabuiuc/pyhealth-book/tree/main/chap8-GenAI/notebook>

evaluates the generated samples based on the Earth Mover’s Distance instead of the KL divergence used in the traditional GANs. The WGAN model is known for its stable and diverse training results.

The Earth Mover’s distance (EMD), also known as the Wasserstein distance, in the context of WGAN, is mathematically formulated as the minimum cost of transporting mass to transform one distribution into the other. For two probability distributions, \mathbb{P}_r (real data distribution) and \mathbb{P}_g (generated data distribution), the Wasserstein distance is defined as:

$$W(\mathbb{P}_r, \mathbb{P}_g) = \inf_{\gamma \in \Pi(\mathbb{P}_r, \mathbb{P}_g)} \mathbb{E}_{(x,y) \sim \gamma} [\|x - y\|].$$

Here, \inf denotes the infimum over all joint distributions γ whose marginals are \mathbb{P}_r and \mathbb{P}_g , and $\Pi(\mathbb{P}_r, \mathbb{P}_g)$ represents the set of all possible joint distributions (couplings) of \mathbb{P}_r and \mathbb{P}_g . The expectation $\mathbb{E}_{(x,y) \sim \gamma} [\|x - y\|]$ calculates the average cost of moving the mass from x to y , representing the amount of work required to transform one distribution into the other.

As for the implementation, WGAN and GAN differ in that:

- WGAN does not use Sigmoid function in the discriminator.
- WGAN uses Wasserstein loss ¹¹ to train the discriminator and the generator.
- WGAN clips the discriminator weights to a limited range after each mini-batch update (such as in $[-0.01, 0.01]$).
- In WGAN, the discriminator is trained more times than the generator (e.g., 5 times).
- WGAN uses RMSProp ¹² as the optimizer with a small learning rate and no momentum.

10.3 Diffusion Models

A diffusion model is a generative deep learning technique used for creating new data, such as images, by learning to reverse a gradual noising process. The model is trained to denoise data by learning the reverse process of gradually adding noise to the data. The key steps are:

1. **Forward process.** Take a training dataset and gradually add noise to the data in a series of steps, resulting in a sequence of increasingly noisy versions of the data. This step does not need parameters.
2. **Reverse process.** Train a neural network to reverse this noising process by learning to predict less noisy versions of the data given the noisy input. This steps require a parameterized step-wise denoising model.
3. **Generation.** To generate new data, start with random noise and iteratively apply the trained reverse process model to remove the noise, eventually resulting in a new sample similar to the training data.

¹¹ https://en.wikipedia.org/wiki/Wasserstein_GAN

¹² <https://pytorch.org/docs/stable/generated/torch.optim.RMSprop.html>

Diffusion models have shown impressive results in generating high-quality images, audio, and other types of data. They are known for their stability and the ability to produce diverse and realistic outputs [137, 138, 139]. These are a new type of generative model inspired by a physical process. Notable works of diffusion models include diffusion probabilistic models [140], noise-conditioned score network [141], and denoising diffusion probabilistic models [129]. Next, we will introduce diffusion models in detail.

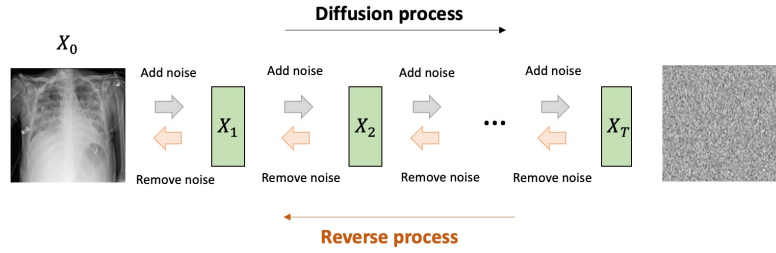


Figure 10.11 Illustration of the diffusion and reverse process in the diffusion model

10.3.1 Forward process of diffusion models

Let's use the example of generating a chest X-ray image to understand the diffusion process. In Figure 10.11, we start with a real chest X-ray image, denoted as \mathbf{X}_0 , and add a small amount of noise step by step, generating $\mathbf{X}_1, \dots, \mathbf{X}_T$. Each \mathbf{X}_t contains more Gaussian noise than its predecessor. Finally, we arrive at \mathbf{X}_T , which contains 100% noise, following a standard Gaussian distribution.

The degree of noise in each step is controlled by factors $\{\beta_t \in (0, 1)\}_{t=1}^T$. We model each step in the diffusion process by Gaussian distribution sampling, starting from \mathbf{X}_0 . Specifically, we have:

$$\mathbf{X}_t \sim q(\mathbf{X}_t | \mathbf{X}_{t-1}) = \mathcal{N}(\sqrt{1 - \beta_t} \mathbf{X}_{t-1}, \beta_t \mathbf{I}), \quad (10.37)$$

Mean and Variance of \mathbf{X}_t

Assume the initial sample is normalized, meaning that $\mathbb{E}[\mathbf{X}_0] = 0$ and $\mathbb{E}[(\mathbf{X}_0)^2] = \mathbf{I}$. Then, sampled \mathbf{X}_t has the following statistics,

$$\text{mean} = \mathbb{E}[\mathbf{X}_t] = \mathbb{E}[\sqrt{1 - \beta_t} \mathbf{X}_{t-1}] = \sqrt{1 - \beta_t} \mathbb{E}[\mathbf{X}_{t-1}] = 0, \quad (10.38)$$

$$\text{variance} = \mathbb{E}[(\mathbf{X}_t)^2] = (\sqrt{1 - \beta_t})^2 \mathbb{E}[(\mathbf{X}_{t-1})^2] + \beta_t \mathbf{I} = \mathbf{I}. \quad (10.39)$$

Following Markov property, the whole diffusion process can be modeled as a product of the Gaussian distribution,

$$\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_T \sim p(\mathbf{X}_1, \dots, \mathbf{X}_T | \mathbf{X}_0) = \prod_{t=1}^T q(\mathbf{X}_t | \mathbf{X}_{t-1}). \quad (10.40)$$

As T becomes bigger (if $T \rightarrow \infty$), the features/patterns from the original data \mathbf{X}_0 gradually becomes indistinguishable and finally fully noisy, which is \mathbf{X}_T .

\mathbf{X}_t is Gaussian Distributed

Using the above diffusion process, the intermediate \mathbf{X}_t are all Gaussian distributed. Let us look at this. Given the initial \mathbf{X}_0 , the first noisy sample is generated by (following Equation 10.37),

$$\mathbf{X}_1 = \sqrt{1 - \beta_1} \mathbf{X}_0 + \sqrt{\beta_1} \epsilon_0 \sim \mathcal{N}(\sqrt{1 - \beta_1} \mathbf{X}_0, \beta_1 \mathbf{I}). \quad (10.41)$$

The second noisy sample is generated by (following Equation 10.37),

$$\mathbf{X}_2 = \sqrt{1 - \beta_2} \mathbf{X}_1 + \sqrt{\beta_2} \epsilon_1 \quad (10.42)$$

$$= \sqrt{1 - \beta_2} (\sqrt{1 - \beta_1} \mathbf{X}_0 + \sqrt{\beta_1} \epsilon_0) + \sqrt{\beta_2} \epsilon_1 \quad (10.43)$$

$$= \sqrt{1 - \beta_1} \sqrt{1 - \beta_2} \mathbf{X}_0 + \sqrt{1 - \beta_2} \sqrt{\beta_1} \epsilon_0 + \sqrt{\beta_2} \epsilon_1 \quad (10.44)$$

$$= \sqrt{(1 - \beta_1)(1 - \beta_2)} \mathbf{X}_0 + \sqrt{1 - (1 - \beta_1)(1 - \beta_2)} \tilde{\epsilon}_1 \quad (10.45)$$

$$\sim \mathcal{N}(\sqrt{(1 - \beta_1)(1 - \beta_2)} \mathbf{X}_0, (1 - (1 - \beta_1)(1 - \beta_2)) \mathbf{I}). \quad (10.46)$$

The last step is given by merging two identical independently distributed Gaussian noise. Since the merge of two Gaussian distribution $\mathcal{N}(\mu_1, \sigma_1^2)$ and $\mathcal{N}(\mu_2, \sigma_2^2)$ is $\mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$. Equivalently, we can get the third noisy sample as

$$\mathbf{X}_3 = \sqrt{1 - \beta_3} \mathbf{X}_2 + \sqrt{\beta_3} \epsilon_2 \quad (10.47)$$

$$= \sqrt{(1 - \beta_1)(1 - \beta_2)(1 - \beta_3)} \mathbf{X}_0 + \sqrt{1 - (1 - \beta_1)(1 - \beta_2)(1 - \beta_3)} \tilde{\epsilon}_2 \quad (10.48)$$

$$\sim \mathcal{N}(\sqrt{(1 - \beta_1)(1 - \beta_2)(1 - \beta_3)} \mathbf{X}_0, (1 - (1 - \beta_1)(1 - \beta_2)(1 - \beta_3)) \mathbf{I}). \quad (10.49)$$

and thus, the last noisy sample is

$$\mathbf{X}_T \sim \mathcal{N}\left(\sqrt{\prod_{t=1}^T (1 - \beta_t)} \mathbf{X}_0, \left(1 - \prod_{t=1}^T (1 - \beta_t)\right) \mathbf{I}\right). \quad (10.50)$$

In practice, the values of $\prod_{t=1}^T (1 - \beta_t)$ will gradually vanish over each timestep, reaching 0 at the final timestep. This allows for a smooth transition from the original data distribution to the noisy distribution, ensuring that \mathbf{X}_T is a fully random Gaussian distribution.

In other words, *the diffusion process does not need any learnable parameters*, and it provides a series of noisy samples $\{\mathbf{X}_t, t = 1, \dots, T\}$ where \mathbf{X}_0 is the original data, \mathbf{X}_t is more noisy than \mathbf{X}_{t-1} , and \mathbf{X}_T is asymptotically randomly gaussian distributed.

10.3.2 Reverse process of diffusion model

After understanding the forward process of the diffusion model gradually adding noise to the original input data, we will next describe the reverse process that gradually denoises the fully noisy image until we recover the initial data in a stochastic way. This

reverse process is motivated by gradient Langevin dynamics [142] where the concept is developed for modeling the molecular systems in physics.

Given that the noise applied at each step is Gaussian distributed, and the sampled intermediate data $\{\mathbf{X}_t\}$ is Gaussian distributed, we thus know the reconstructed intermediate data is also Gaussian distributed. Let us denote the recovering function during the process is $f_\theta(\mathbf{X}_{t-1} | \mathbf{X}_t)$ with parameters θ . Thus, the first step is to obtain a \mathbf{X}_{T-1} from \mathbf{X}_T .

$$\mathbf{X}_{T-1} \sim f_\theta(\mathbf{X}_{T-1} | \mathbf{X}_T) = \mathcal{N}(\mu_\theta(\mathbf{X}_T, T), \Sigma_\theta(\mathbf{X}_T, T)). \quad (10.51)$$

and step two is

$$\mathbf{X}_{T-2} \sim f_\theta(\mathbf{X}_{T-2} | \mathbf{X}_{T-1}) = \mathcal{N}(\mu_\theta(\mathbf{X}_{T-1}, T-1), \Sigma_\theta(\mathbf{X}_{T-1}, T-1)). \quad (10.52)$$

until

$$\mathbf{X}_0 \sim f_\theta(\mathbf{X}_0 | \mathbf{X}_1) = \mathcal{N}(\mu_\theta(\mathbf{X}_1, 1), \Sigma_\theta(\mathbf{X}_1, 1)). \quad (10.53)$$

Here, we split the parameterized functions f_θ into a function for predicting the mean value $\mu_\theta(\mathbf{X}_t, t)$ and another one for predicting the covariance matrix $\Sigma_\theta(\mathbf{X}_t, t)$. They both take the sample \mathbf{X}_t and the reverse step t as inputs. To train these two functions, we need to know the real distribution of \mathbf{X}_{t-1} from \mathbf{X}_t , which can be obtained from the posterior at the forward process.

$$q(\mathbf{X}_{t-1} | \mathbf{X}_t, \mathbf{X}_0) = \frac{q(\mathbf{X}_t, \mathbf{X}_{t-1} | \mathbf{X}_0)}{q(\mathbf{X}_t | \mathbf{X}_0)} = q(\mathbf{X}_t | \mathbf{X}_{t-1}, \mathbf{X}_0) \frac{q(\mathbf{X}_{t-1} | \mathbf{X}_0)}{q(\mathbf{X}_t | \mathbf{X}_0)}. \quad (10.54)$$

Now, we could use the posterior distribution $q(\mathbf{X}_{t-1} | \mathbf{X}_t, \mathbf{X}_0)$ as the ground truth to regulate the parameterized posterior $f_\theta(\mathbf{X}_{t-1} | \mathbf{X}_t)$ by KL divergence. We do not elaborate on the detailed math derivation of the objective, which could be found in Weng et al. [143].

Objective functions

Thus, the final objective function of the diffusion model contains two parts: one part is still the reconstruction loss, and another part is a series of KL divergence on $q(\mathbf{X}_{t-1} | \mathbf{X}_t, \mathbf{X}_0)$ and $f_\theta(\mathbf{X}_{t-1} | \mathbf{X}_t)$ to regulate the intermediate variables. Compared to VAE model, we find that the high-level difference here is that VAE uses one hidden variable \mathbf{h} (the bottleneck layer) while diffusion model uses a series of hidden variable \mathbf{X}^t for generation.

10.3.3 PyTorch Implementation of Diffusion Model

In this section, we will use the MNIST data set to show how to implement a simple diffusion model in PyTorch. The model explained in this section is a simpler version, where we do not normalize \mathbf{X}_t to follow $\mathcal{N}(0, \mathbf{I})$ and we use a deterministic form in the reverse process.

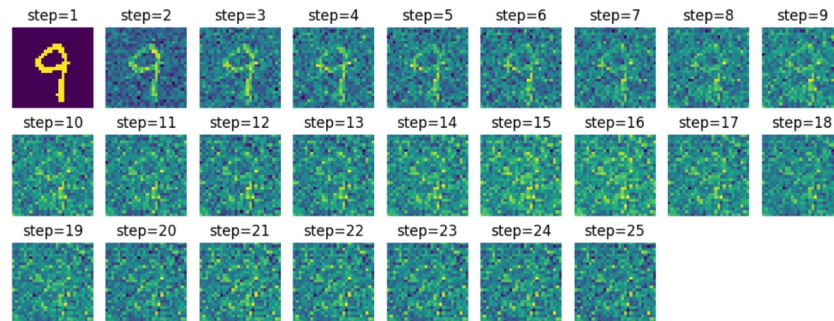


Figure 10.12 Illustration of Diffusion Process (25 steps) on MNIST Digit Image

Simpler diffusion process: adding noise

We first define a `unit_normal` sampler that has the following function `unit_normal.sample([shape])` to sample a Gaussian random variable with any shape. We get the first image from the first batch and its label, named `image` and `label`, and then gradually add noise to it. Given the original image X_0 , the diffusion process is shown in Figure 10.12.

```

1 sigma = 0.3
2
3 # Create a normal distribution sampler
4 unit_normal = torch.distributions.Normal(0, sigma)
5 unit_normal.loc = unit_normal.loc.to(device)
6 unit_normal.scale = unit_normal.scale.to(device)
7
8 # Demo: Diffusion process
9 images, labels = next(iter(trainloader))
10 image = images[0]
11 label = labels[0]
12
13 # Plot the images at different steps of diffusion process
14 plt.figure(figsize=(10, 4))
15 for i in range(diffusion_steps):
16     plt.subplot(3, 9, i + 1)
17     plt.imshow(image)
18     plt.axis('off')
19     plt.title(f"Step={i+1}")
20
21     # Sample noise from a normal distribution
22     noise = unit_normal.sample(image.shape).cpu()
23     # Add the noise to the image
24     image += noise
25
26 plt.tight_layout()

```

The simple diffusion process is explained below:

- **Step 0:** Initial image is X_0 , and we do not normalize it.

- **Step 1:** $\mathbf{X}_1 = \mathbf{X}_0 + \epsilon_1$, where $\epsilon_1 \sim \mathcal{N}(0, \Sigma)$ is a random noise. Σ is a diagonal matrix with each diagonal value equals to 0.3, i.e., $\Sigma = 0.3\mathbf{I}$.
- **Step 2:** $\mathbf{X}_2 = \mathbf{X}_1 + \epsilon_2 = \mathbf{X}_0 + \epsilon_1 + \epsilon_2$, where $\epsilon_2 \sim \mathcal{N}(0, \Sigma)$ is a another random noise. In fact, since ϵ_1, ϵ_2 are independent, \mathbf{X}_2 is equivalent to $\mathbf{X}_2 = \mathbf{X}_0 + \epsilon_1^*$, where $\epsilon_1^* \sim \mathcal{N}(0, 2\Sigma)$.
- **Step t:** Similarly, we have that $\mathbf{X}_t = \mathbf{X}_{t-1} + \epsilon_t = \mathbf{X}_0 + \epsilon_t^*$, where $\epsilon_t^* \sim \mathcal{N}(0, t\Sigma)$.

In sum, each intermediate image \mathbf{X}_t has two ways to construct: (i) adding a noise sample from $\mathcal{N}(0, \Sigma)$ to \mathbf{X}_{t-1} ; (ii) adding a noise sample from $\mathcal{N}(0, t\Sigma)$ to the original image \mathbf{X}_0 .

Simpler Reverse Process and Training: De-noising

In the reverse process, we start with a Gaussian noise sample $\tilde{\mathbf{X}}_{25} \sim \mathcal{N}(0, 25\Sigma)$, assuming we have 25 diffusion steps. Although the final sample at diffusion step 25 should be the original image plus the accumulated noise, we can treat $\tilde{\mathbf{X}}_{25}$ as the final diffusion result since the noise component significantly outweighs the original image.

The reverse process relies on a function f_θ that takes the t -th intermediate image \mathbf{X}_t , the diffusion step t , and the data label as inputs. The function aims to infer the additional noise $\tilde{\epsilon}_t$ at diffusion step t :

$$\tilde{\epsilon}_t = f_\theta(\mathbf{X}_t, t, \text{label}) \sim \mathcal{N}(0, \Sigma). \quad (10.55)$$

The output of f_θ is post-processed to follow the same distribution as the real noise ϵ_t , which is $\mathcal{N}(0, \Sigma)$.

Starting from the pure Gaussian noise $\tilde{\mathbf{X}}_{25}$, we generate the estimated noise $\tilde{\epsilon}_{25}$ using f_θ . We then approximate the previous image $\tilde{\mathbf{X}}_{24}$ by subtracting the estimated noise from $\tilde{\mathbf{X}}_{25}$:

$$\tilde{\mathbf{X}}_{24} = \tilde{\mathbf{X}}_{25} - \tilde{\epsilon}_{25}. \quad (10.56)$$

Similarly, we obtain the next noise estimate $\tilde{\epsilon}_{24}$ by feeding $(\tilde{\mathbf{X}}_{24}, 24, \text{label})$ to f_θ . We then perform the de-noising step to approximate $\tilde{\mathbf{X}}_{23}$:

$$\tilde{\mathbf{X}}_{23} = \tilde{\mathbf{X}}_{24} - \tilde{\epsilon}_{24}. \quad (10.57)$$

This process is repeated step by step until we finally obtain the denoised image $\tilde{\mathbf{X}}_0$.

Training Sample Construction

The training objective is to minimize the difference between the estimated noise $\tilde{\epsilon}_t$ and the real noise ϵ_t used in the forward diffusion process. By learning to predict and remove the noise at each step accurately, the model learns to generate realistic images from random Gaussian noise.

To optimize the function f_θ , we can use the following method to construct training samples.

- Get an image and the label from the training set, \mathbf{X}_0 and label .

- Randomly select a step $t \in \{1, 2, \dots, 25\}$ and generate a noise $\epsilon_t^* \sim \mathcal{N}(0, t\Sigma)$.
- Then, randomly generate a noise $\epsilon \sim \mathcal{N}(0, \Sigma)$.
- We let $\mathbf{X}_t = \mathbf{X}_0 + \epsilon_t^*$, and let $\mathbf{X}_{t+1} = \mathbf{X}_t + \epsilon$.
- The ground truth is \mathbf{X}_t , and we use \mathbf{X}_{t+1} to predict the noise.

$$\tilde{\epsilon} = f_{\theta}(\mathbf{X}_{t+1}, t, \text{label}). \quad (10.58)$$

- Then, the approximated de-noised version at step t is

$$\tilde{\mathbf{X}}_t = \mathbf{X}_{t+1} - \tilde{\epsilon}. \quad (10.59)$$

- We use mean square error (MSE) loss as the objective between the real noise and the predicted noise,

$$\mathcal{L} = \|\epsilon - \tilde{\epsilon}\|_F^2. \quad (10.60)$$

Next we illustrate the PyTorch implementation.

Load MNIST data

Similar to previous GAN and VAE examples, we load the MNIST training data as the training set.

```

1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4 import torchvision
5
6 import numpy as np
7 from tqdm import tqdm
8 import matplotlib.pyplot as plt
9
10 device = torch.device('cuda:0')
11
12 # for the training set, we load the data and label
13 trainset = torchvision.datasets.MNIST(\
14     root='./data', train = True, download = True)
15 xtrain = trainset.data.numpy()
16 ytrain = trainset.targets.numpy()
17
18 print (xtrain.shape, ytrain.shape)
19 """
20 (60000, 28, 28) (60000,)
21 """
22
23 # transform pixel range (0, 255) into {0,1} by splitting at pixel 128
24 xtrain = np.where(xtrain > 128, 1, 0) #pixel > 128 assigning to 1, otherwise 0.
25 xtrain = xtrain.astype(np.float32)
26
27 batch_size = 512
28 trainloader = torch.utils.data.DataLoader([xtrain[i], ytrain[i]] \
29     for i in range(len(ytrain)), shuffle=True, batch_size=batch_size)
30
31 # sigma is the noise std
32 sigma = 0.3
    
```

Define the Diffusion Model

We define the diffusion model f_θ as a two-layer neural network. In the `self.forward()` function, the model takes a batch of samples `x`, the step of the diffusion process `step`, and the label of the images as input `label`. Then, it concatenates the image 28×28 pixels, one hot `step` encoding, and one hot `label` encoding as the features into the two-layer neural network. The output is further normalized into $\mathcal{N}(0, \Sigma)$ distribution.

```

1  class Diffuser(nn.Module):
2      def __init__(self, diff_steps, n_labels):
3          super(Diffuser, self).__init__()
4
5          # the diffusion model is a two-layer neural network
6          self.diffusion = nn.Sequential(
7              nn.Linear(28*28 + diff_steps + n_labels, 28*28),
8              nn.ReLU(),
9              nn.Linear(28*28, 28*28),
10             )
11
12             self.diff_steps = diff_steps
13             self.n_labels = n_labels
14
15         def forward(self, x, step, label):
16
17             # one hot for diffusion step
18             step_encoding = F.one_hot(
19                 torch.tensor([step]*x.shape[0]), num_classes = self.diff_steps
20             ).to(device)
21
22             # one hot for label
23             labels_encoding = F.one_hot(label, \
24                 num_classes = self.n_labels).to(device)
25
26             # concat features and run diffusion process
27             emb = torch.cat([x, step_encoding, labels_encoding], 1)
28             out = self.diffusion(emb)
29
30             # normalize the output
31             MEAN = out.mean(dim=1, keepdim=True)
32             STD = out.std(dim=1, keepdim=True) + 1e-8
33             out = (out - MEAN) / STD * sigma
34             return out
35
36         diffusion_steps = 25
37         model = Diffuser(25, 10)
38         model.to(device) # move model to GPU

```

Model Training and Evaluation

As part of the model training process, we require an ϵ noise sampler, which we define directly on the device for faster execution. This sampler generates Gaussian distributed data for any given shape using the syntax `unit_normal.sample([shape])`.

In the following Python code snippet, we create a normal distribution sampler and transfer it to the CUDA device:

```

1  # normal distribution sampler (we could throw to cuda for once)
2  unit_normal = torch.distributions.Normal(0, sigma)
3  unit_normal.loc = unit_normal.loc.to(device)
4  unit_normal.scale = unit_normal.scale.to(device)

```

The evaluation function begins with a 100% noisy sample starting from random noise sampled from $\mathcal{N}(0, 25\Sigma)$ at step 25. The model then progressively predicts the per-step noise and removes it from the current intermediate images. We repeat this process ten times, once for each label, and then print out the resulting images for digits 0 to 9.

```

1  def evaluate(model, diffusion_steps):
2
3      # 0, 1, ..., 9
4      labels = torch.arange(10)
5
6      # random images (assume at step 25)
7      noisy_images = unit_normal.sample((len(labels), 28, 28)) \
8          * np.sqrt(diffusion_steps)
9
10     plt.figure(figsize=(10, 5))
11
12     # we run diffusion process for multiple steps
13     # to denoise the noisy images into a digit image
14     for diffusion_step in range(diffusion_steps):
15         step = diffusion_steps - diffusion_step - 1
16         with torch.no_grad():
17             noise = model(noisy_images.view(10, 28*28), step, labels)
18             noisy_images -= noise.view(len(labels), 28, 28)
19
20     for i in range(10):
21         plt.subplot(1, 10, i+1)
22         # print the digit one by one
23         plt.imshow(noisy_images[i].cpu())
24         plt.axis('off')
25     plt.show()
26
27 evaluate(model, diffusion_steps)

```

We show the generated synthetic digit images before training in Figure 10.13, which are basically random noises.

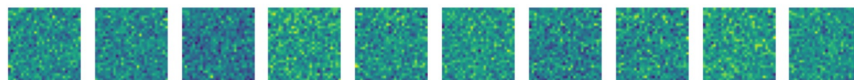


Figure 10.13 Illustration of Generated Digit Images from Diffusion Model (before training)

During the training process, we use Adam as the optimizer and MSE loss. We use 50 epochs to train the diffusion model. For each batch from the training data, we randomly shuffle the diffusion steps and use Equation (10.58) and Equation (10.60) to optimize the parameters. The model will be evaluated every 5 epochs.

```

1 optimizer = torch.optim.Adam(model.parameters(), 0.1)
2 criterion = nn.MSELoss()
3
4 for epoch in range(50):
5     epoch_loss = 0
6     for images, labels in tqdm(trainloader):
7         images = images.to(device).view(-1, 28*28)
8
9         for step in torch.randperm(diffusion_steps):
10             # get a sample at step: original image + sqrt(step) * normal sample
11             less_noisy = images + \
12                 unit_normal.sample(images.shape) * np.sqrt(step)
13             # get a normal sample
14             one_step_noise = unit_normal.sample(images.shape)
15             # get the sample at (step + 1)
16             more_noisy = less_noisy + one_step_noise
17
18             output = model(more_noisy, step, labels)
19             loss = criterion(output, one_step_noise)
20
21             optimizer.zero_grad()
22             loss.backward()
23             optimizer.step()
24
25             epoch_loss += loss.item()
26         print(f'Epoch: {epoch} \t Loss: {epoch_loss}')
27
28     if epoch % 5 == 0:
29         evaluate(model, diffusion_steps)

```

The final generated digit images drawn from `evaluate()` function are shown in Figure 10.14, which is of high quality.

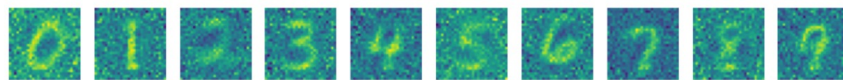


Figure 10.14 Illustration of Generated Digit Images from Diffusion Model

A concrete notebook including PyTorch implementations of diffusion model could be found in this public repository¹³.

10.4 Chest X-ray Image Generation with VAE, GAN, and Diffusion Model

Chest X-ray images are useful in diagnosing numerous conditions, including heart diseases and respiratory conditions. Generating synthetic chest X-ray images can be beneficial in improving the accuracy of machine learning models. These artificial images can augment small real datasets or allow other researchers to replicate research

¹³ <https://github.com/sunlabuiuc/pyhealth-book/tree/main/chap8-GenAI/notebook>

models without disclosing real patient information. Note that the example requires using pyhealth package in python.

In this section, we will demonstrate how to use the VAE, GAN, and diffusion models in generating chest X-ray images based on the PyHealth package. We will be using the dataset from the COVID-19 Radiography Database¹⁴.

10.4.1 Step 1: data loading

For chest X-ray dataset loading, we use the *COVID19CXRDataset* API which will process the COVID19 chest X-ray image dataset. This API only requires the root path to the dataset as the argument, and everything else is handled by the API already.

```
1 from pyhealth.datasets import COVID19CXRDataset
2
3 root = "/srv/local/data/COVID-19_Radiography_Dataset"
4 base_dataset = COVID19CXRDataset(root)
5 base_dataset.stat()
6
7 """
8 Statistics of COVID19CXRDataset:
9 Number of samples: 21165
10 Number of classes: 4
11 Class distribution: Counter({'Normal': 10192, 'Lung Opacity': 6012, \
12                               'COVID': 3616, 'Viral Pneumonia': 1345})
13 """
```

We use the `.stat()` function to show the dataset information. In the dataset, there are 21,165 sample images and they belong to 4 different categories: normal, lung opacity, covid, and viral pneumonia. The label is not used in unconditional VAE and GAM model training, while they can be used to train the conditional version, such as ConditionalVAE, ConditionalGAN, that generates the data sample under certain labels.

10.4.2 Step 2: define machine learning task

We start by using the `set_task()` function to process the dataset into a dictionary structure. Each sample is formatted as a dictionary, and the dataset's default task is to classify chest X-ray images into different categories. We use the default task to process the dataset.

```
1 base_dataset.default_task
2 # COVID19CXRClassification(task_name='COVID19CXRClassification', \
3 #                           input_schema={'path': 'image'}, output_schema={'label': 'label'})
4
5 sample_dataset = base_dataset.set_task()
```

Next, we add another sample transformation step to further clean up the image data. The transformation first converts the data into tensor format and scales the image pixel

¹⁴ <https://www.kaggle.com/datasets/tawsifurrahman/covid19-radiography-database>

intensity into the range of $[0, 1]$. We then use the “Lambda” function to transform all images into three channels. This custom lambda transformation checks if the input image x has three channels (i.e., RGB). If the image already has three channels, it is returned as is. Otherwise, if the image has a single channel (e.g., grayscale), it is repeated three times along the channel dimension to convert it to a three-channel image. We also resize the chest X-ray image from 228×228 to 128×128 to make it smaller for efficiency of the demonstration.

```

1  from torchvision import transforms
2
3  the transformation automatically normalize the pixel intensity into [0, 1]
4  transform = transforms.Compose([
5      # use three channels
6      transforms.Lambda(lambda x: x if x.shape[0] == 3 else x.repeat(3, 1, 1),
7      transforms.Resize((128, 128))),
8      transforms.Normalize(mean=[0.5862785803043838], std=[0.27950088968644304])
9  ])
10
11 def encode(sample):
12     sample["path"] = transform(sample["path"])
13     return sample
14
15 sample_dataset.set_transform(encode)

```

After transforming the data, we split it into three parts: training, validation, and test sets. The sizes of these sets are distributed in the ratio of 60% : 20% : 20% respectively. We obtained three data loaders and noticed that each chest X-ray image has a dimension of $1 \times 128 \times 128$. The size of the training set is 12,699, while the validation and test sets have a size of 4,233 each.

```

1  from pyhealth.datasets import split_by_visit, get_dataloader
2
3  # split dataset
4  train_dataset, val_dataset, test_dataset = split_by_visit(
5      sample_dataset, [0.6, 0.2, 0.2]
6  )
7  train_dataloader = get_dataloader(train_dataset, batch_size=256, shuffle=True)
8  val_dataloader = get_dataloader(val_dataset, batch_size=256, shuffle=False)
9  test_dataloader = get_dataloader(test_dataset, batch_size=256, shuffle=False)
10
11 data = next(iter(train_dataloader))
12 print (data["path"][0].shape)
13 print (
14     "loader size: train/val/test",
15     len(train_dataset),
16     len(val_dataset),
17     len(test_dataset),
18 )
19
20 """
21 torch.Size([1, 128, 128])
22 loader size: train/val/test 12699 4233 4233
23 """

```

10.4.3 Step 3: initialize the VAE model

PyHealth offers the VAE API, which has arguments similar to other models. However, in this case, the feature key and the label key are the same, and we use the regression mode. In VAE, the input is the data sample, and the reconstruction supervision is also the same data sample (image "path" in this example), which is treated as a regression task in PyHealth. The hidden dimension represents the dimension of the Gaussian distribution. There are two new arguments: (i) `input_channel` indicates the number of channels in the input, and (ii) `input_size` represents the size of the data (height or weight), and they must be the same. Modifying these two arguments will lead to different model architectures, such as the number of layers, etc.

```

1  from pyhealth.models import VAE
2
3  model = VAE(
4      dataset=sample_dataset,
5      input_channel=3,
6      input_size=128,
7      feature_keys=["path"],
8      label_key="path",
9      mode="regression",
10     hidden_dim = 256,
11 )

```

10.4.4 Step 4&5: model training and inference

To train the model, we use the *pyhealth.trainer* API. The learning task is treated as a regression task, and we record the KL divergence, mean square error (MSE), and mean absolute error (MAE) in this generation application. We monitor the change in KL divergence over each epoch and consider the epoch with the lowest KL divergence as the best trained model.

```

1  from pyhealth.trainer import Trainer
2
3  trainer = Trainer(model=model, device="cuda:0", \
4                  metrics=["kl_divergence", "mse", "mae"])
5  trainer.train(
6      train_dataloader=train_dataloader,
7      val_dataloader=val_dataloader,
8      epochs=20,
9      monitor = "kl_divergence",
10     monitor_criterion = 'min',
11     optimizer_params={"lr": 2e-3},
12 )

```

We use the evaluation function to assess the reconstruction error on the test data, which is shown below.

```

1  # evaluate on the test data
2  print(trainer.evaluate(test_dataloader))
3

```

```

4 """
5 {'kl_divergence': 0.092681274, 'mse': 2.3712183e-15, \
6   'mae': 3.5363602e-08, 'loss': 19134.9755565672}
7 """

```

10.4.5 Evaluate the reconstruction

We use the `inference()` function to reconstruct all the samples in the test set and then compare with the original data samples. As shown in Figure 10.15, the reconstructed image can restore the overall shape and some patterns of the original image, however, detailed information is missing.

```

1 X, X_rec, _ = trainer.inference(test_dataloader)
2
3 import matplotlib.pyplot as plt
4
5 idx = 4 # choose which image to visualize
6 plt.figure()
7 plt.subplot(1, 2, 1)
8 plt.imshow(X[idx].reshape(128, 128), cmap="gray")
9 plt.subplot(1, 2, 2)
10 plt.imshow(X_rec[idx].reshape(128, 128), cmap="gray")
11 plt.show()

```

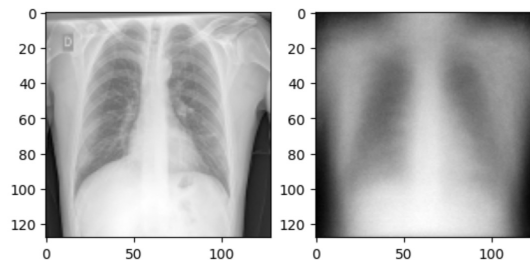


Figure 10.15 Compare the original data (left) and the reconstruction (right)

10.4.6 Generating new chest xray images

After training the VAE model, we can use the decoder to create new chest X-ray images that have a similar style to the ones used for training. To generate these images, we start by selecting a random vector from a multi-variant normal distribution. We have chosen to use 128 as the hidden gaussian dimension, which means we generate a 128-dimensional vector \mathbf{x} . We then input this vector into the decoder model, which reshapes the output into a 128x128 image.

```

1 import torch
2 import numpy as np

```

```

3 model = trainer.model
4
5
6 model.eval()
7 with torch.no_grad():
8     # we generate a random 64-dimensional multi-gaussian vector
9     x = np.random.normal(0, 1, 256)
10
11     x = x.astype(np.float32)
12     x = torch.from_numpy(x).to(trainer.device)\
13         .unsqueeze(0).unsqueeze(2).unsqueeze(3)
14     rec = model.decoder(x).detach().cpu().numpy() #convert to numpy array
15     rec = rec.reshape((3, 128, 128))[0] #reshape to an image
16     plt.imshow(rec, cmap="gray")
17     plt.show()

```

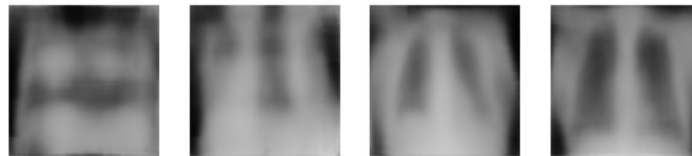


Figure 10.16 A synthesized chest X-ray image from VAE

Remarks

The generated images or the reconstructed images seem to have poor qualities. Improvements are possible with more training data, more training epochs, and better encoder and decoder models. Readers can try the pipeline on their own and tune the models in this repository¹⁵.

10.4.7 Use the GAN model for generation

In this section, we will be using the GAN model from PyHealth. This model will be used to generate Chest X-ray images. The GAN model requires alternative updating steps for training the discriminator and the generator separately. Therefore, we will create our own training pipeline on top of the PyHealth GAN model.

To initialize the GAN model, we set the input image channels to 3, the input image size to 128, and the hidden dimension to 256 (random vector size). The GAN model we are using is based on a deep CNN that uses 3-channel images with a resolution of 128. We will use the same train/val/test loaders from the VAE pipeline.

```

1 from pyhealth.models import GAN
2
3 model = GAN(
4     input_channel=3,
5     input_size=128,

```

¹⁵ <https://github.com/sunlabuiuc/pyhealth-book/tree/main/chap8-GenAI/notebook>

```
6     hidden_dim = 256,
7 )
```

Training GAN models

In the following steps, we will demonstrate how to train a GAN model. Firstly, we define the loss function for the discriminator and use two different optimizers for the generator and the discriminator.

We keep track of the loss of the generator and discriminator separately and store them in `curve_G` and `curve_D` respectively.

During the training, both real images (with label 1) and generated (fake, with label 0) images are used to train the discriminator. Then, we use the generated images to train the generator, while ignoring the loss function for the discriminators. Through this process, the discriminator learns to distinguish between real and fake images, and the generator becomes stronger in generating fake images that can fool the discriminator.

```
1 import torch
2 from tqdm import tqdm
3
4 device = torch.device("cuda:0")
5
6 # Loss function
7 loss = torch.nn.BCELoss()
8 opt_G = torch.optim.AdamW(model.generator.parameters(), lr=1e-3)
9 opt_D = torch.optim.AdamW(model.discriminator.parameters(), lr=1e-4)
10
11 model.to(device)
12
13 curve_D, curve_G = [], []
14
15 for epoch in range(30):
16     curve_G.append(0)
17     curve_D.append(0)
18     for batch in tqdm(train_dataloader):
19
20         """ train discriminator """
21         opt_D.zero_grad()
22
23         real_imgs = torch.stack(batch["path"], dim=0).to(device)
24         batch_size = real_imgs.shape[0]
25         fake_imgs = model.generate_fake(batch_size, device)
26
27         real_loss = loss(model.discriminator(real_imgs), \
28                         torch.ones(batch_size, 1).to(device))
29         fake_loss = loss(model.discriminator(fake_imgs), \
30                         torch.zeros(batch_size, 1).to(device))
31         loss_D = (real_loss + fake_loss) / 2
32
33         loss_D.backward()
34         opt_D.step()
35
36         """ train generator """
37         opt_G.zero_grad()
```

10.4 Chest X-ray Image Generation with VAE, GAN, and Diffusion Model

217

```

38     loss_G = loss(model.discriminator(fake_imgs), \
39                   torch.ones(batch_size, 1).to(device))
40
41     loss_G.backward()
42     opt_G.step()
43
44     curve_G[-1] += loss_G.item()
45     curve_D[-1] += loss_D.item()
46
47     print (f"epoch: {epoch} --- loss of G: {curve_G[-1]}, \
48           loss of D: {curve_D[-1]}")
49
50     """
51     epoch: 0 --- loss of G: 223.14766198396683, loss of D: 14.372084847651422]
52     epoch: 1 --- loss of G: 466.95815896987915, loss of D: 0.5429049517260864
53     epoch: 2 --- loss of G: 541.9039497375488, loss of D: 0.154949375835713
54     epoch: 3 --- loss of G: 585.4883279800415, loss of D: 0.08368747876374982
55     ...
56     epoch: 26 --- loss of G: 83.02981567382812, loss of D: 35.7115815281868
57     epoch: 27 --- loss of G: 79.00199449062347, loss of D: 44.660759538412094
58     epoch: 28 --- loss of G: 76.95186269283295, loss of D: 39.98225525021553
59     epoch: 29 --- loss of G: 68.42279094457626, loss of D: 42.02642643451691
60     """

```

Synthetic image from GAN

Similarly, we let GAN to generate synthetic chest X-ray image, shown in Figure 10.17.

```

1  import torch
2  import numpy as np
3  import matplotlib.pyplot as plt
4
5  model.eval()
6  with torch.no_grad():
7      fake_imgs = model.generate_fake(1, device).detach().cpu()
8      plt.imshow(fake_imgs[0][0], cmap="gray")
9      plt.show()

```

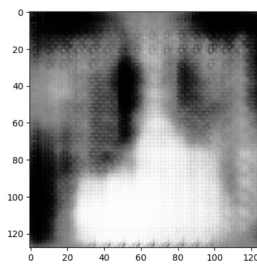


Figure 10.17 A synthesized chest xray image generated by GAN

10.4.8 Use the latent diffusion model for generation

We will be using the diffusion model to create synthetic chest X-ray images. However, simple diffusion models can be difficult to train when working with raw images that have high dimensions. To address this, we use a latent diffusion model (LDM) that operates on a latent space. Here’s how it works:

- It first trains a VAE model on the training data, and then the VAE model is completely fixed at later steps.
- Then, on the same training data, we encode the training images into latent space by the VAE encoder. We train a diffusion model on the latent embeddings.
- During data generation, a random vector is drawn from Gaussian distribution, and the vector runs through the reverse process of the diffusion model to generate a latent embedding. The latent embedding is then fed into the VAE decoder to generate a synthetic chest X-ray image.

Remarks

The LDM combines the strengths of both VAE and the simple diffusion model. VAE is great at encoding data into a low-dimensional latent space, but its generation power is often limited due to the distribution of the latent embeddings not following an independent multi-variate Gaussian. On the other hand, the diffusion model excels at generating synthetic data that follows the distribution in low-dimensional space, but can fail or be hard to train on high-dimensional spaces. The LDM overcomes these issues by leveraging the diffusion model to transform random noise into a low-dimensional distribution that aligns with VAE latent space. Then, it uses the VAE decoder to generate a high-dimensional real image. Below, you can find the code.

Train a VAE model as before

We copy the previous codes below to train a VAE model first.

```
1 from pyhealth.models import VAE
2
3 vae = VAE(
4     dataset=sample_dataset,
5     input_channel=3,
6     input_size=128,
7     feature_keys=["path"],
8     label_key="path",
9     mode="regression",
10    hidden_dim = 256,
11 )
```

```
1 from pyhealth.trainer import Trainer
2
3 trainer = Trainer(model=vae, device=device, metrics=["kl_divergence", "mse", "
4                                                       mae"])
5
6 trainer.train(
7     train_dataloader=train_dataloader,
8     val_dataloader=val_dataloader,
```

10.4 Chest X-ray Image Generation with VAE, GAN, and Diffusion Model

219

```

7     epochs=20,
8     monitor = "kl_divergence",
9     monitor_criterion = 'min',
10    optimizer_params={"lr": 1e-3},
11 )

```

We generate random vectors and directly use the VAE decoder for generation, shown in Figure 10.18. As we mentioned in the Remark, the latent space of VAE usually does not align with independent multi-variate Gaussian distribution, and thus the generation from a Gaussian random noise does not look good.

```

1  import torch
2  import numpy as np
3
4  N = 4
5
6  vae = trainer.model
7  vae.eval()
8  # we generate a random 64-dimensional multi-Gaussian vector
9  x = np.random.normal(0, 1, (N, 256))
10 x = x.astype(np.float32)
11 x = torch.from_numpy(x).to(trainer.device).unsqueeze(2).unsqueeze(3)
12
13 with torch.no_grad():
14     rec = vae.decoder(x).detach().cpu()
15
16 plt.figure(figsize=(7, 5))
17 for i in range(N):
18     plt.subplot(1, N, i+1)
19     plt.imshow(rec[i].reshape(3, 128, 128).permute(1,2,0))
20     plt.axis('off')
21 plt.show()

```

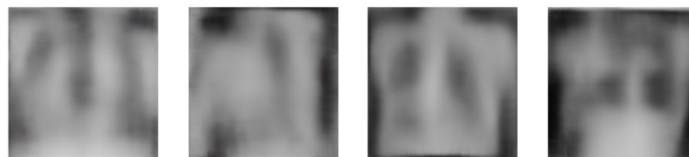


Figure 10.18 Several synthesized chest x-ray images generated by VAE

Build the LDM model

Now, we will train a diffusion model to improve the generation quality. We build the diffusion model, which is almost the same as the one we used for MNIST that takes the image `x`, the diffusion steps `step`, and the image label `label` as input. However, the input data dimension becomes the VAE bottleneck layer dimension 256.

```

1  import torch
2  import torch.nn as nn
3  import torch.nn.functional as F

```



```

4
5 class Diffuser(nn.Module):
6     def __init__(self, diff_steps, n_labels):
7         super(Diffuser, self).__init__()
8
9         self.diffusion = nn.Sequential(
10             nn.Linear(256 + diff_steps + n_labels, 256),
11             nn.ReLU(),
12             nn.Linear(256, 256),
13         )
14
15         self.diff_steps = diff_steps
16         self.n_labels = n_labels
17
18     def forward(self, x, step, label):
19
20         # one hot for diffusion step
21         step_encoding = F.one_hot(
22             torch.tensor([step]*x.shape[0]), num_classes = self.diff_steps
23         ).to(device)
24
25         # one hot for label
26         labels_encoding = F.one_hot(label, num_classes = self.n_labels).to(
27             device)
28
29         # concat features and run diffusion process
30         emb = torch.cat([x, step_encoding, labels_encoding], 1)
31         out = self.diffusion(emb)
32
33         # normalize the output
34         MEAN = out.mean(dim=1, keepdim=True)
35         STD = out.std(dim=1, keepdim=True) + 1e-8
36         out = (out - MEAN) / STD * sigma
37         return out
38
39 diffusion_steps = 25
40 ldm = Diffuser(25, 4)
41 ldm.to(device) # move model to GPU

```

We use the same random noise sampler as defined below.

```

1 ### sampler
2 sigma = 0.24
3 # normal distribution sampler (we could throw to cuda for once)
4 unit_normal = torch.distributions.Normal(0, sigma)
5 unit_normal.loc = unit_normal.loc.to(device)
6 unit_normal.scale = unit_normal.scale.to(device)

```

let us evaluate the performance of the `ldm` model before training. The generation process starts from a random noise, then it goes through the reverse process to be a latent vector, and next the latent vector is transformed into a Chest X-ray image by the fixed VAE decoder. We shown the initial performance in Figure 10.19. It is interesting that even before model training, the LDM generated images looks better than VAE generated images.

10.4 Chest X-ray Image Generation with VAE, GAN, and Diffusion Model

221

```

1 def evaluate(model, diffusion_steps):
2
3     # 0, 1, 2, 3
4     labels = torch.arange(4)
5
6     # random images
7     noisy_images = unit_normal.sample((len(labels), 256)) \
8                                     * np.sqrt(diffusion_steps)
9
10    plt.figure(figsize=(7, 5))
11
12    # we run diffusion process for multiple steps
13    # to denoise the noisy images into a digit image
14    for diffusion_step in range(diffusion_steps):
15        step = diffusion_steps - diffusion_step - 1
16        with torch.no_grad():
17            noise = ldm(noisy_images, step, labels)
18            noisy_images -= noise.view(len(labels), 256)
19
20    for i in range(4):
21        plt.subplot(1, 4, i+1)
22        # print the digit one by one
23        latent = noisy_images[i:i+1].unsqueeze(-1).unsqueeze(-1)
24        with torch.no_grad():
25            reconstruction = vae.decoder(latent)[0]
26        plt.imshow(reconstruction.cpu().permute(1, 2, 0))
27        plt.axis('off')
28    plt.show()
29
30    evaluate(model, diffusion_steps)

```

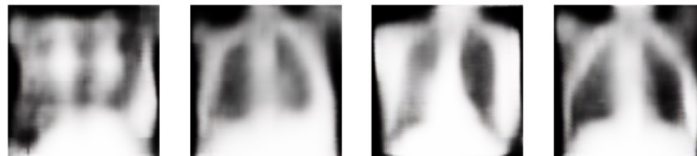


Figure 10.19 Several synthesized chest xray images generated by the LDM model (before training)

Train the LDM model

We train the LDM model before using similar pipeline as we used for training a diffusion model on MNIST, except that before feeding into the model, we transform the images into latent space by VAE encoder.

```

1 from tqdm import tqdm
2
3 optimizer = torch.optim.Adam(ldm.parameters(), 0.1)
4 criterion = nn.MSELoss()
5
6 mapping = {"Normal": 0, "COVID": 1, "Lung Opacity": 2, "Viral Pneumonia": 3}

```

```

7
8 for epoch in range(25):
9     epoch_loss = 0
10    for batch in train_dataloader:
11        images = torch.stack(batch["path"], dim=0).to(device)
12        labels = torch.LongTensor([mapping[item] \
13                                   for item in batch["label"]]).to(device)
14        with torch.no_grad():
15            latents = vae.encoder(images)[0]
16
17        for step in torch.randperm(diffusion_steps):
18            # get a sample at step: original image + sqrt(step) * normal sample
19            less_noisy = latents + \
20                unit_normal.sample(latents.shape) * np.sqrt(step)
21            # get a normal sample
22            one_step_noise = unit_normal.sample(latents.shape)
23            # get the sample at (step + 1)
24            more_noisy = less_noisy + one_step_noise
25
26            output = ldm(more_noisy, step, labels)
27            loss = criterion(output, one_step_noise)
28
29            optimizer.zero_grad()
30            loss.backward()
31            optimizer.step()
32
33            epoch_loss += loss.item()
34        print(f'Epoch: {epoch} \t Loss: {epoch_loss}')
35
36    if epoch % 5 == 0:
37        evaluate(ldm, diffusion_steps)
38
39    """
40    Epoch: 0    Loss: 144.3587410505861
41    Epoch: 1    Loss: 136.6373864300549
42    Epoch: 2    Loss: 135.48954015411437
43    Epoch: 3    Loss: 134.957813994959
44    Epoch: 4    Loss: 134.53044639341533
45    Epoch: 5    Loss: 134.18359554558992
46    ...
47    Epoch: 21    Loss: 132.3953073080629
48    Epoch: 22    Loss: 132.3351447628811
49    Epoch: 23    Loss: 132.276965896599
50    Epoch: 24    Loss: 132.23989922646433
51    """

```

We show the LDM generated images after training in Figure 10.20, which have higher qualities compared to the generated images from VAE or GAN.

The notebook resources could be found at the open repository ¹⁶, and readers could play with the latent diffusion model on their own laptop.

¹⁶ <https://github.com/sunlabuiuc/pyhealth-book/tree/main/chap8-GenAI/notebook>

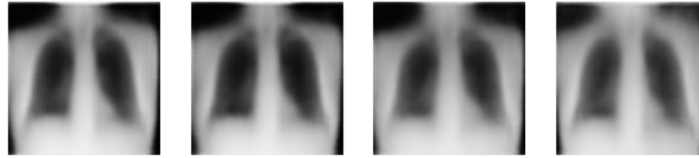


Figure 10.20 Several synthesized Chest X-ray images generated by the LDM model (after training)

10.5 Takeaways

- **Generative Models:** Different from previous supervised models that mainly predict some targeted values or classes based on complex data features, generative models are built mostly on unlabeled data and could synthesize new data samples.
- **Variational Auto-encoder (VAE):** VAE is one type of generative model that leverages the auto-encoder architecture and enforces the middle bottleneck layer to have standard distribution. It leverages MSE reconstruction loss and KL divergence as the objective.
- **Generative Adversarial Network (GAN):** GAN is another type of generative model that employs a discriminator and a generator where the generator strives to generate realistic samples from random noise and the discriminator tries to distinguish the real samples and the generated samples. The model is trained in an adversarial way by optimizing the generator and the discriminator alternatively.
- **Diffusion Models:** Diffusion model is a newly developed generative model that add noise to the real data progressively and then optimizes a model to predict the per-step noise which could gradually generate a realistic data from full noise.
- **Latent Diffusion Models (LDM):** LDM is practically a better choice than diffusion model. LDM needs a well-trained VAE model as predecessor. The trained VAE model first uses its encoder to map high-dimensional samples into latent low-dimensional space where the LDM is trained. During generation, the noise goes through the LDM model and then uses the VAE decoder to synthesize data.

Questions

- What are three generative models introduced in the chapter? When are they proposed?
- What is the main architecture of the AE model? How would you design the AE model for learning image representations from unlabeled image data? Could you specify the model architecture below?
- What are different ways to prevent overfitting in an AE model? Could you elaborate how they works?
- What is VAE? What are the main difference between VAE and common AE models?
- What is evidence lower bound (ELBO) in VAE? Could you derive the inequality?

- Elaborate the architecture of the GAN model and write down its objective functions.
- Discussion the difference between VAE and GAN.
- What are the major problems in training GAN models? Could you come up with some ideas to mitigate or improve the problem? Why they help? (you may search the google and show your answer with citations).
- Explain the diffusion models in a high level. Write down the diffusion forward and reverse process using equations or bullet points.
- What are two different ways to create intermediate noisy sample in diffusion models?
- In the reverse process of diffusion model, please use Bayesian rule to derive the quantity $q(\mathbf{X}_{t-1} | \mathbf{X}_t, \mathbf{X}_0)$ by three probability quantity from the diffusion process, $q(\mathbf{X}_t | \mathbf{X}_{t-1}, \mathbf{X}_0)$, $q(\mathbf{X}_{t-1} | \mathbf{X}_0)$, $q(\mathbf{X}_t | \mathbf{X}_0)$. Could you further simplify the results? HINT: by the Markov property, we could further decompose $q(\mathbf{X}_t | \mathbf{X}_0)$ into $q(\mathbf{X}_t | \mathbf{X}_{t-1})q(\mathbf{X}_{t-1} | \mathbf{X}_{t-2}) \cdots q(\mathbf{X}_1 | \mathbf{X}_0)$.
- **Programming Question:** Could you use either VAE or GAN or LDM to generate one synthetic image for bird and cat? Your training data is the entire CIFAR-10. You could build VAE or GAN or Diffusion model from scratch using PyTorch or directly use existing models from PyHealth. The CIFAR-10 dataset could be downloaded from ¹⁷.

¹⁷ <https://www.cs.toronto.edu/~kriz/cifar.html>