

Generative Adversarial Networks (GANs)

Lapland University of Applied Sciences

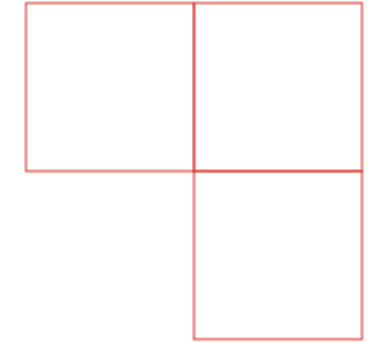
Dr. Manuel García Fernández

manuel.garcia2@universidadeuropea.es

Ve más allá

Index

- Fundamentals of Representation Learning **01**
- Introduction to GANs: the Vanilla-GAN **02**
- Flavours of GANs: f-GAN, WGAN and KL-WGAN **03**
- The Conditional GAN (CGAN): pix2pix **04**



A decorative arrangement of squares in red and white. A large red square is the central element. To its left, two white squares are stacked vertically. Above the red square, one white square is centered. To the right of the red square, two white squares are stacked vertically. Below the red square, one white square is centered. Further to the right, there is a single white square. In the bottom right corner, there is a solid red square.

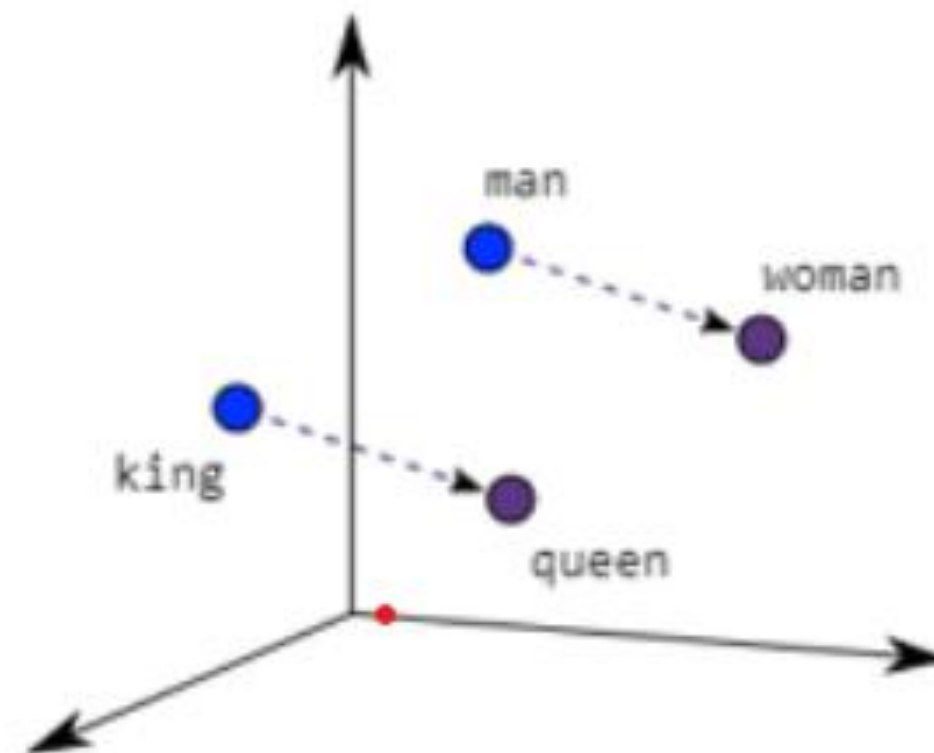
01

Fundamentals of Representation Learning

Fundamentals of Representation Learning

Representation learning is a foundational field within machine learning that focuses on automatically discovering and encoding meaningful features or abstractions from raw data to facilitate downstream tasks such as classification, prediction, or generation.

At the heart of this process lies the concept of a **latent space, which is a lower-dimensional, often continuous, embedding of the original high-dimensional data**. This latent space serves as the compressed form that captures the essential properties, underlying structure, and semantics of the data, making it central to how representation learning functions.



E.g.: King-queen are represented along the same axis. Man/King are represented at the same side as they refer both to the same gender

Fundamentals of Representation Learning

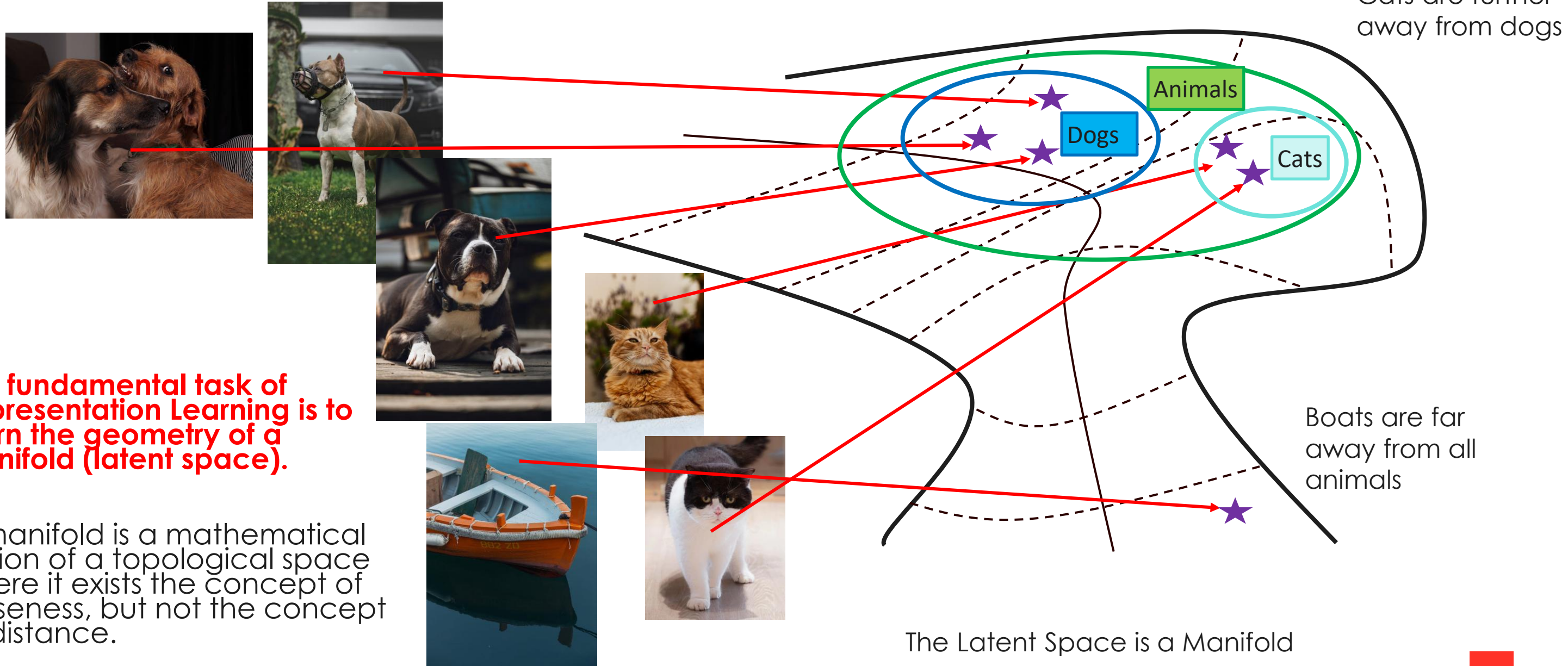
The primary objective of representation learning is to transform complex, high-dimensional data into a more compact, robust, and semantically meaningful representation. This transformation significantly reduces the input space's redundancy, benefiting various pattern recognition tasks. The latent space acts as this reduced-dimensionality vector space where the complex data is embedded, often capturing intricate semantic properties.

Latent spaces inherently perform dimensionality reduction by mapping high-dimensional input data into a lower-dimensional manifold. This process is crucial for managing computational complexity and focusing on the most salient features of the data. For instance, network representation learning embeds network vertices into a low-dimensional vector space, preserving network topology and other side information, thereby simplifying analysis. In complex scenarios like analyzing unprocessed images with thousands of pixels, latent variable models can effectively capture underlying structure. The representations learned in the latent space are designed to be more compact and robust, making them ideal for subsequent analysis.

A well-constructed latent space organizes data in a way that reflects meaningful attributes and preserves intrinsic geometric properties. It allows for the discovery of interpretable representations where different latent dimensions might correspond to distinct semantic features or generative factors



Fundamentals of Representation Learning



Fundamentals of Representation Learning

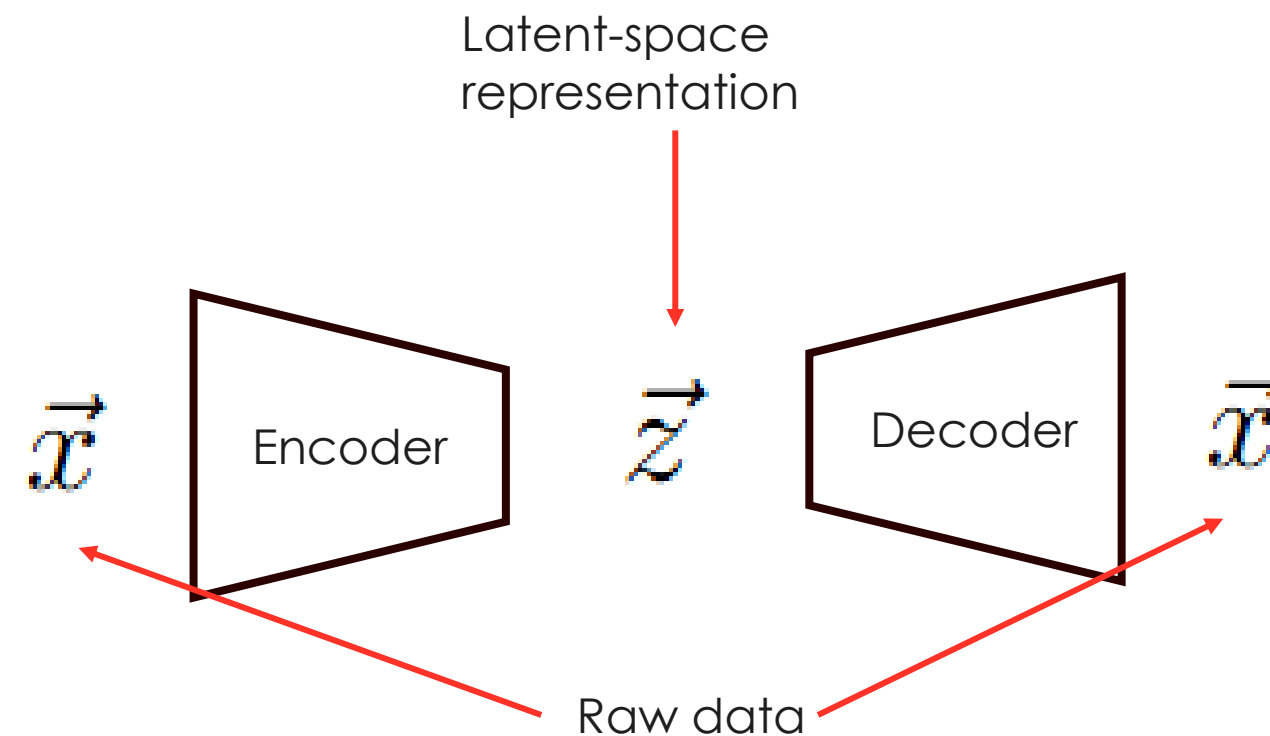
For leveraging Representation Learning and allowing a proper definition of the geometry of manifold, one has to learn the mapping between the raw high-dimensional space and the compressed latent-space.

Different approaches can be used. The simplest approach that uses Neural Networks is the Autoencoder.

An autoencoder is a set of two Neural Networks:

- An Encoder: which compresses data from an N-dimensional space to a M-dimensional space ($N < M$)
- A Decoder: which decompresses back from the M-dimensional space.

Typically, target goal of an Autoencoder seeks to minimize the reconstruction error, so Mean Squared Error (MSE) is used as loss function to be minimized.



$$\mathcal{L}_{AE} = \frac{1}{N} \sum_{i=1}^N \|\vec{x} - D(E(\vec{x}))\|^2$$

Fundamentals of Representation Learning

On the Autoencoder, it was learned a invertible map from the data-space to the latent space.

- Encoder is a direct map from the real-space to the latent-space.
- Decoder is an invers map from the latent-space to the real-space.

Autoencoders are very useful for compressing data as they have are an invertible relationship. In addition, as the reconstruction error can be determined, this can be used for other applications, such as anomaly detection.

Autoencoders can be used as synthetic data generators by just by samling on the latent space with random noise and decode it with the Decoder.

Nevertheless, synthetic data generated by Autoencoders is not good, as the induced geometry of the manifold is not very accurate, as we are limited by the training data instances and interpolation does not make any sense, as close points not necessarily are similar.

As we only measure the geometry of the manifold at real training data instances, the sampling of the latent-space is very sparse, being bad at non-seen data instances.



Fundamentals of Representation Learning

Encoder and Decoder are Feedforward neural networks. They are symmetric and read MNIST images flattened ($28 \times 28 = 784$)

The output of the Encoder is the input of the Decoder. The number of neurons of the output layer is the dimension number of the latent space

```

11 class Autoencoder(nn.Module):
12     def __init__(self, latent_dim=128):
13         super(Autoencoder, self).__init__()
14
15         # Encoder
16         self.encoder = nn.Sequential(
17             nn.Linear(784, 512),
18             nn.ReLU(),
19             nn.Linear(512, 256),
20             nn.ReLU(),
21             nn.Linear(256, 128),
22             nn.ReLU(),
23             nn.Linear(128, latent_dim)
24         )
25
26         # Decoder
27         self.decoder = nn.Sequential(
28             nn.Linear(latent_dim, 128),
29             nn.ReLU(),
30             nn.Linear(128, 256),
31             nn.ReLU(),
32             nn.Linear(256, 512),
33             nn.ReLU(),
34             nn.Linear(512, 784),
35             nn.Sigmoid()
36         )
37
38     def forward(self, x):
39         encoded = self.encoder(x)
40         decoded = self.decoder(encoded)
41         return decoded
42
43     def encode(self, x):
44         return self.encoder(x)
45
46     def train_autoencoder():
47         # Data loading
48         transform = transforms.Compose([
49             transforms.ToTensor(),
50             transforms.Normalize((0.5,), (0.5,))
51         ])
52
53         train_dataset = datasets.MNIST(root='./data', train=True, download=True, transform=transform)
54         train_loader = DataLoader(train_dataset, batch_size=128, shuffle=True)
55
56         # Model setup
57         device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
58         model = Autoencoder(latent_dim=16).to(device)
59         criterion = nn.MSELoss()
60         optimizer = optim.Adam(model.parameters(), lr=1e-4)
61
62         # Training
63         model.train()
64         loss_history = []
65
66         for epoch in range(300):
67             total_loss = 0
68             for data, _ in train_loader:
69                 data = data.view(-1, 784).to(device)
70
71                 optimizer.zero_grad()
72                 reconstructed = model(data)
73                 loss = criterion(reconstructed, data)
74                 loss.backward()
75                 optimizer.step()
76
77                 total_loss += loss.item()
78
79             avg_loss = total_loss / len(train_loader)
80             loss_history.append(avg_loss)
81             print(f'Epoch {epoch+1}/300, Loss: {avg_loss:.4f}')
82
83     return model, loss_history
84

```

Loss function is Mean Square Error (MSE)

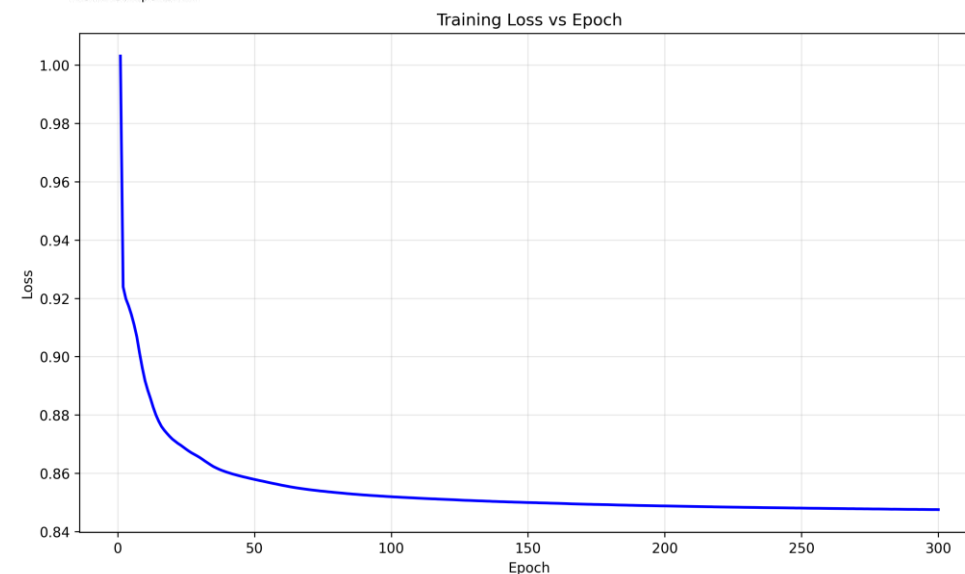
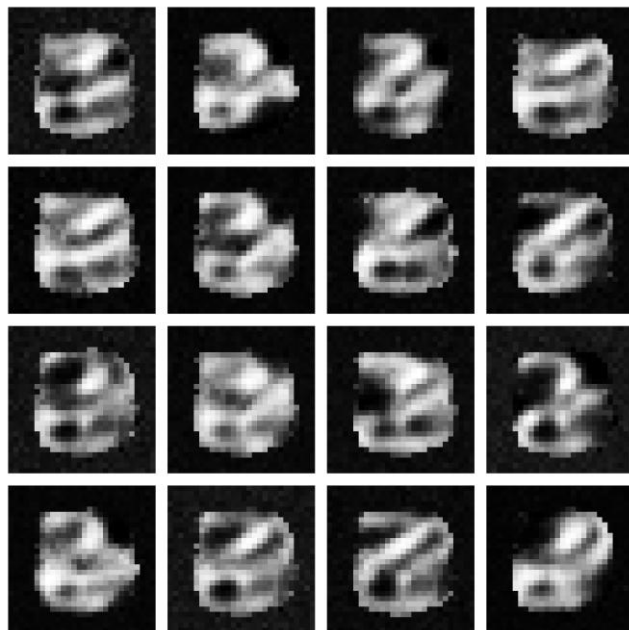
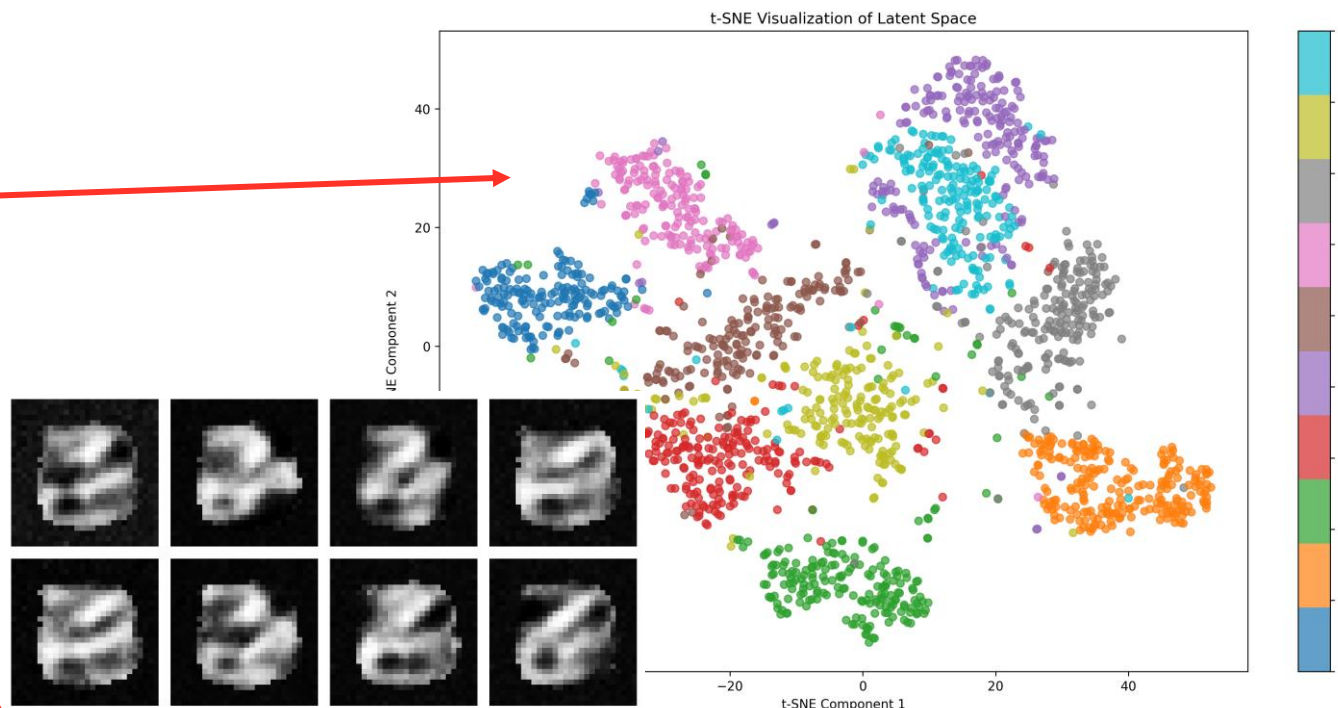
Fundamentals of Representation Learning

Although using a t-SNE visualization of the latent-space displays a fairly reasonable separation of digits for the training dataset.

The synthetic dataset displays poor images

The learned latent space by an Autoencoder has the following properties:

- Latent vectors are all deterministic (1-to-1 relationship).
- No guarantee of continuity or smoothnes.
- No guarantee simmilar points lie closer.



Fundamentals of Representation Learning

An improvement can be made at standard Autoencoder with the introduction of the Variational Autoencoder (VAE).

A Variational Autoencoder does not map raw point to latent-space points, but as probability distributions at the latent-space.

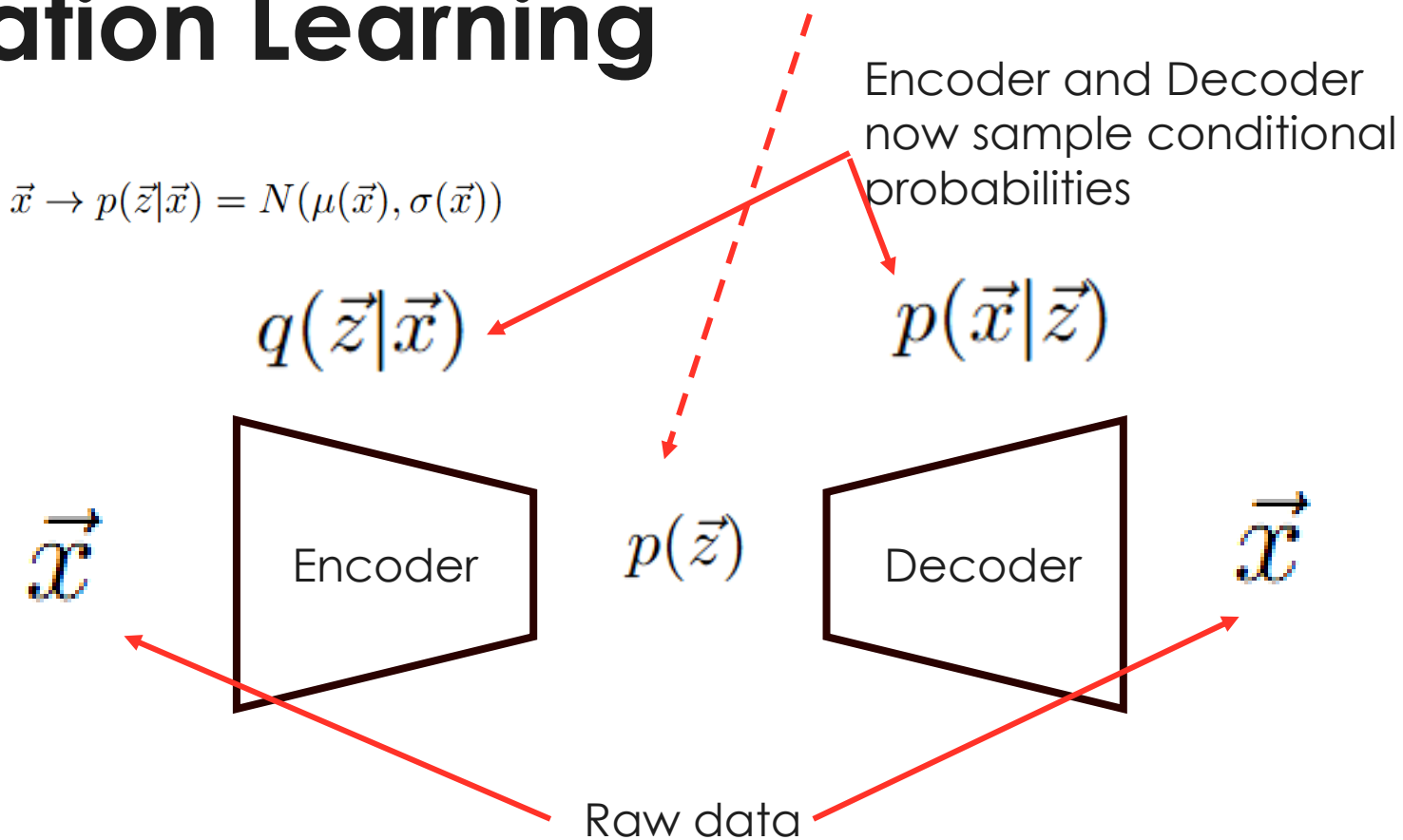
Probability distribution can be anything, but normally is assumed to be a Gaussian. Thus, we aim to learn the mean and variance of that Gaussian.

To enforce the learning of the probability distribution, we demand to minimize the reconstruction loss, but also that probability distributions resemble each other (with KL-divergence).

$$\mathcal{L}_{VAE} = \mathcal{L}_{AE} - D_{KL}(q(\vec{z}|\vec{x}), p(\vec{z})) = \frac{1}{N} \sum_{i=1}^N \|\vec{x} - D(E(\vec{x}))\|^2 - \frac{1}{2} \sum_{i=1}^N (\sigma_i^2 + \mu_i^2 - 1 - \log \sigma_i^2)$$

Latent-space representation now induces a probability distribution.

$$\vec{x} \rightarrow p(\vec{z}|\vec{x}) = N(\mu(\vec{x}), \sigma(\vec{x}))$$

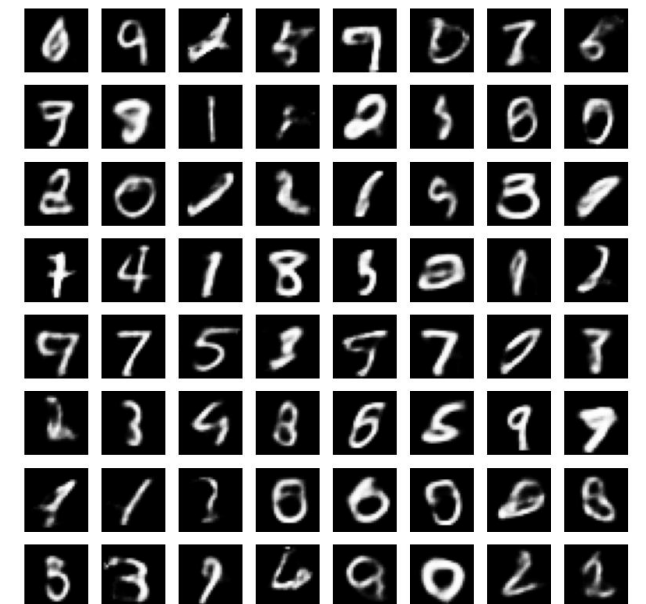
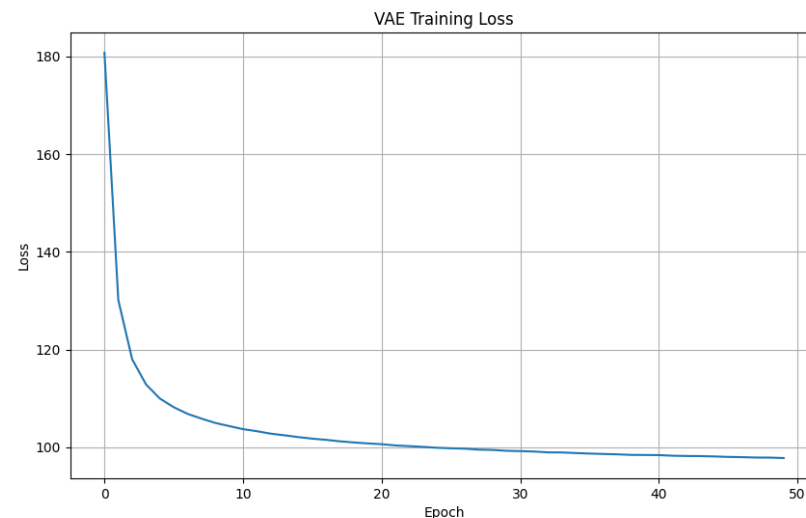
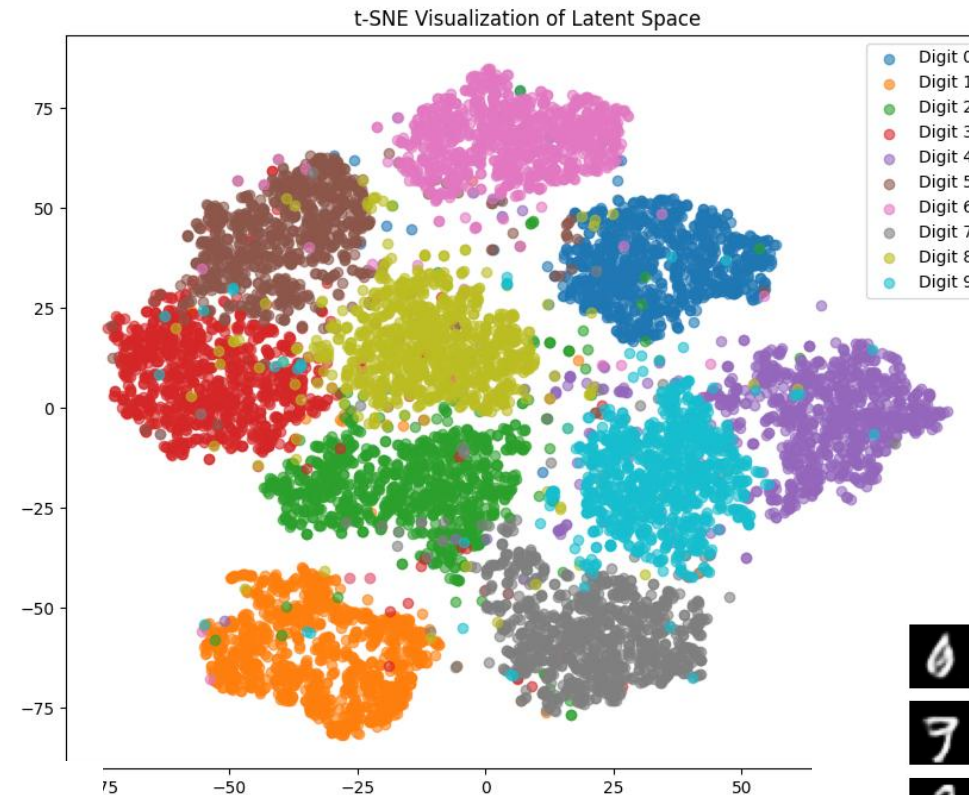


KL-divergence of two gaussians

Fundamentals of Representation Learning

The learned manifold contrary to the Autoencoder has:

- Smoothnes. The transition between different datapoints is not abrupt as in AE.
- Closeness as a meaning (close points mean similar things).
- It is a continuous manifold (no holes).
- The mapping between real data and latent space is probabilistic.
- **Smoothnes ensures that interpolation has a meaning: this implies that unseen regions of the manifold mean something. Thus it is a truly generative model.**



02

Introduction to GANs

Introduction to GANs

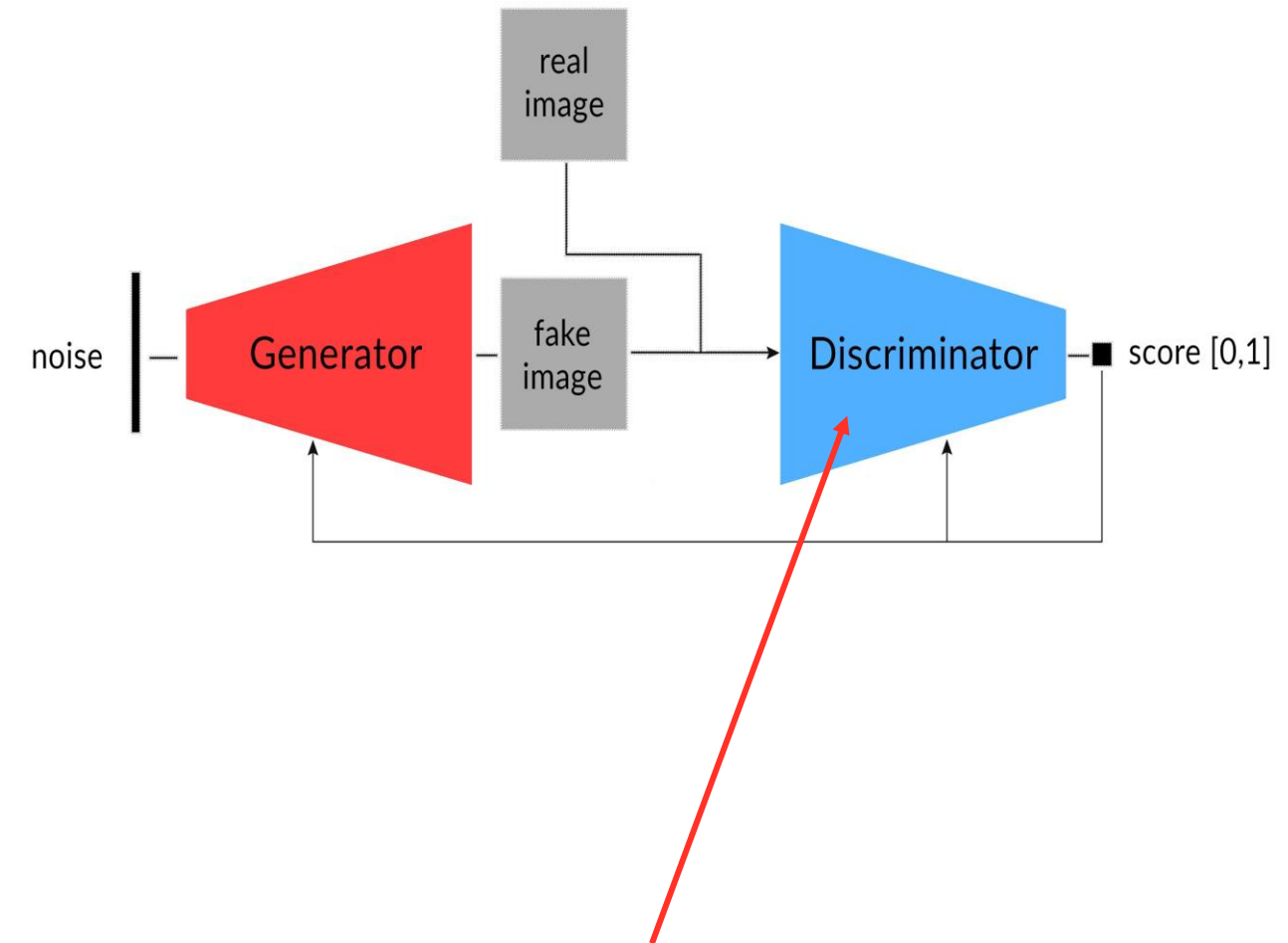
Generative Adversarial Networks (GANs) are set of two networks competing one with the other one.

It is composed by a Generator network and Discriminator network.

- Generator: it is a network that from pure random noise produces synthetic data instances.
- Discriminator: It is a network that tries to guess if a given data is real (from the training dataset) or is synthetic (generated by the Generator network).

Training process is such as Generator and Discriminator compete with each other in a competitive game, where we aim Generator win and generate so Good data instances that it fools the Discriminator.

Training process is such that only the Discriminator sees the real training data. Generator learns as the Discriminator “teaches” that network.



Discriminator sees synthetic and real data.
Discriminator can learn to differentiate both, as we at the training process know which are the real data and which from the Generator

Introduction to GANs

On the GAN, only the inverse map from the latent-space to the real-space is learned (which is the Generator).

The induced geometry of a GAN on the manifold, is much more accurate due to the sampling process. At the sampling of the latent-space, we generate random noise, each of which will correspond to a point of the manifold. Thus, if we generate a large enough sample of random noise, we can reach virtually any desired accuracy on the geometry of the manifold.

Nevertheless, practical issues happen at training: as is a competitive game between two neural networks, convergence of the solution is not always stable and is very complex.



Introduction to GANs

In practice GANs produce better synthetic data than VAEs because:

- Generator is not rewarded to cover the full data distribution, rather for producing images that lie on the manifold of realistic images.
- Producing accurate data comes at expenses of ignoring very unlikely events.
- Contrary VAEs are trained to explain all possible variations of data. Thus, when multimodality happens, the safest bet is to average modes (e.g. a person with eyes looking at the left or looking at right it appear with blurry eyes, because VAE is averaging left-looking and right-looking).
- VAEs are rewarded for covering all modes of the distribution, GANs are rewarded for covering very well any of the modes (even if not covering some of them). **This is the problem of mode-collapse of GANs.**

Currently, state-of-the art of generative models for general image generation include Diffusion models, which combine advantages of GANs and VAEs: training stability, full probability density coverage while producing realistic data.

Introduction to GANs

The simplest GAN we can make is the Vanilla-GAN.

The Vanilla-GAN constitutes two feedforward networks: one for the Generator and another for the Discriminator.

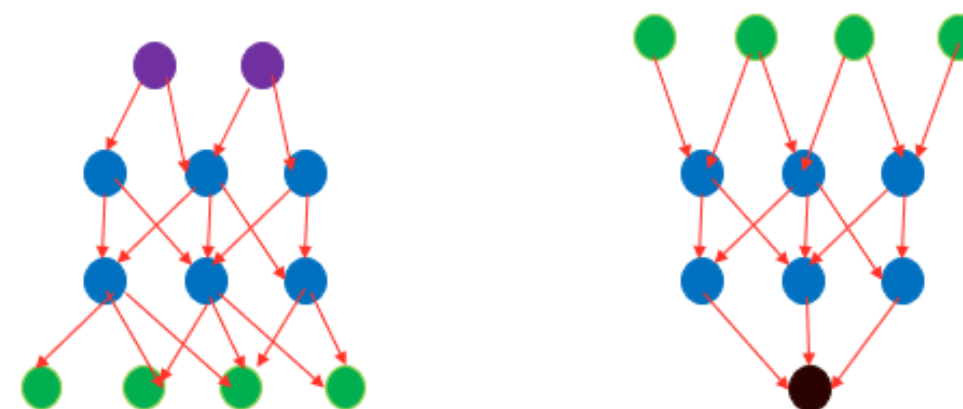
Each network Will have its own loss function.

Generator

$$\mathcal{L}_G = \frac{1}{m} \sum_{i=1}^m \log(1 - D(G(z_i)))$$

Discriminator

$$\mathcal{L}_D = \frac{1}{m} \sum_{i=1}^m [\log D(x_i) + \log(1 - D(G(z_i)))]$$



Generator

Discriminator

- Synthetic/Real data
- Random noise
- Probability synthetic/real

Introduction to GANs

Latent dimension is the number of neurons of the input layer of the Generator

Generator is a feedforward neural network with ReLU activation layers.

Discriminator is a feedforward neural network with LeakyReLU activation layers. LeakyReLU are needed because Generator learns through the Discriminator. Negative part of ReLU has no gradient, complicating the learning of the Generator

```
mnist_vanillagan.py > ...
1  import torch
2  import torchvision
3  from torch.utils.data import DataLoader
4
5  import torch.nn as nn
6  import torch.optim as optim
7  import torchvision.transforms as transforms
8  import matplotlib.pyplot as plt
9
10 # Hyperparameters
11 latent_dim = 100
12 hidden_dim = 256
13 image_dim = 28 * 28
14 batch_size = 64
15 learning_rate = 0.0002
16 epochs = 50
17
18 # Generator
19 class Generator(nn.Module):
20     def __init__(self, latent_dim, hidden_dim, image_dim):
21         super(Generator, self).__init__()
22         self.model = nn.Sequential(
23             nn.Linear(latent_dim, hidden_dim),
24             nn.ReLU(),
25             nn.Linear(hidden_dim, hidden_dim * 2),
26             nn.ReLU(),
27             nn.Linear(hidden_dim * 2, hidden_dim * 4),
28             nn.ReLU(),
29             nn.Linear(hidden_dim * 4, image_dim),
30             nn.Tanh()
31         )
32
33     def forward(self, x):
34         return self.model(x)
35
36 # Discriminator
37 class Discriminator(nn.Module):
38     def __init__(self, image_dim, hidden_dim):
39         super(Discriminator, self).__init__()
40         self.model = nn.Sequential(
41             nn.Linear(image_dim, hidden_dim * 4),
42             nn.LeakyReLU(0.2),
43             nn.Linear(hidden_dim * 4, hidden_dim * 2),
44             nn.LeakyReLU(0.2),
45             nn.Linear(hidden_dim * 2, hidden_dim),
46             nn.LeakyReLU(0.2),
47             nn.Linear(hidden_dim, 1),
48             nn.Sigmoid()
49         )
50
51     def forward(self, x):
52         return self.model(x)
53
```

Introduction to GANs

Real and fake labels are generated as tensors with ones and zeros

We first show to the discriminator real data and compute the loss

Then we generate a fake image from random noise

We show the fake images to the discriminator and compute the loss

Apply the backpropagation and advance the optimizer

Now we generate a new fake image and compute the loss, but with “real” labels (fooling the discriminator).

```

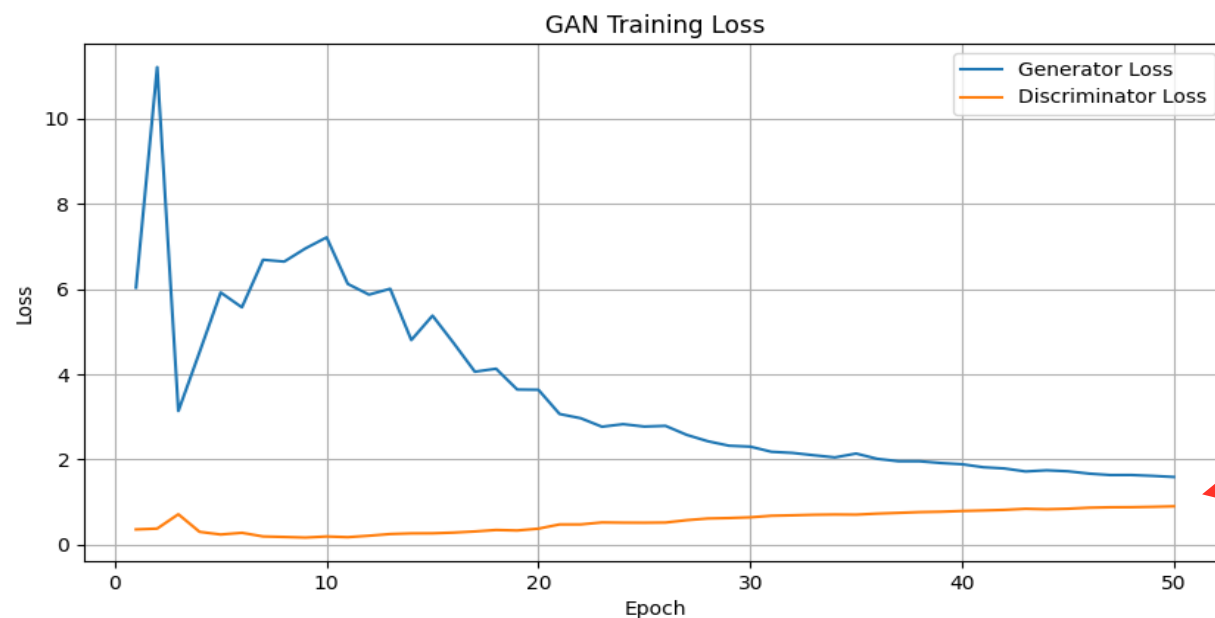
147
148 dataloader = DataLoader(dataset, batch_size=batch_size, shuffle=True)
149
150 # Initialize models
151 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
152 generator = Generator(latent_dim, hidden_dim, image_dim).to(device)
153 discriminator = Discriminator(image_dim, hidden_dim).to(device)
154
155 # Loss and optimizers
156 criterion = nn.BCELoss()
157 optimizer_G = optim.Adam(generator.parameters(), lr=learning_rate)
158 optimizer_D = optim.Adam(discriminator.parameters(), lr=learning_rate)
159
160 # Lists to store losses
161 g_losses = []
162 d_losses = []
163
164 # Training loop
165 for epoch in range(epochs):
166     epoch_g_loss = 0
167     epoch_d_loss = 0
168     num_batches = len(dataloader)
169
170     for batch_idx, (real_images, _) in enumerate(dataloader):
171         real_images = real_images.view(-1, image_dim).to(device)
172         batch_size = real_images.size(0)
173
174         # Train Discriminator
175         real_labels = torch.ones(batch_size, 1).to(device)
176         fake_labels = torch.zeros(batch_size, 1).to(device)
177
178         # Real images
179         outputs = discriminator(real_images)
180         d_loss_real = criterion(outputs, real_labels)
181
182         # Fake images
183         z = torch.randn(batch_size, latent_dim).to(device)
184         fake_images = generator(z)
185         outputs = discriminator(fake_images.detach())
186         d_loss_fake = criterion(outputs, fake_labels)
187
188         d_loss = d_loss_real + d_loss_fake
189         optimizer_D.zero_grad()
190         d_loss.backward()
191         optimizer_D.step()
192
193         # Train generator
194         z = torch.randn(batch_size, latent_dim).to(device)
195         fake_images = generator(z)
196         outputs = discriminator(fake_images)
197         g_loss = criterion(outputs, real_labels)
198
199         optimizer_G.zero_grad()
200         g_loss.backward()
201         optimizer_G.step()
202
203         epoch_g_loss += g_loss.item()
204         epoch_d_loss += d_loss.item()
205
206     # Calculate average losses for the epoch
207     avg_g_loss = epoch_g_loss / num_batches
208     avg_d_loss = epoch_d_loss / num_batches
209
210     g_losses.append(avg_g_loss)
211     d_losses.append(avg_d_loss)
212
213     print(f'Epoch [{epoch+1}/{epochs}], d_loss: {avg_d_loss:.4f}, g_loss: {avg_g_loss:.4f}')
214

```

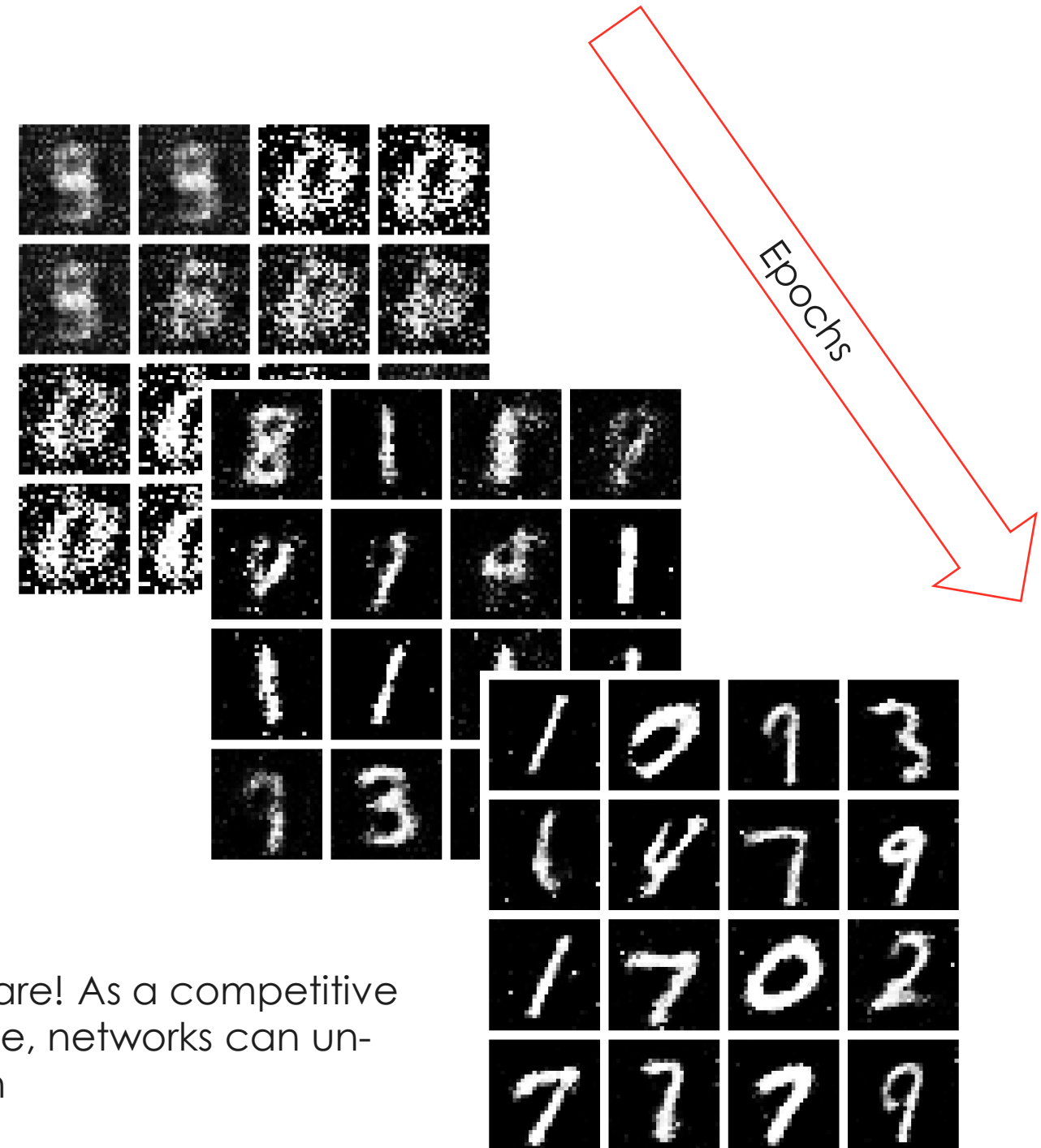
Introduction to GANs

Epoch by epoch generated images are more close to the real ones.

We do not have control over which number Will be generated



Beware! As a competitive game, networks can un-learn



A decorative arrangement of squares in red and white. A large red square is the central element. To its left, two white squares are stacked vertically. Above the red square, one white square is centered. To the right of the red square, two white squares are stacked vertically. Below the red square, one white square is centered. Further to the right, there is a single white square. In the bottom right corner, there is a solid red square.

03

Flavours of GANs

Flavours of GANs

But what is the target mathematical goal of a GAN?

We aim to follow the min-max $\min_{\theta_G} \max_{\theta_D} V(D, G)$ where V can be selected as the most general form as

$$V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_d(\mathbf{x})} [\mathbb{E}_{y \sim p_d(y)} [g_f(D(y|\mathbf{x}))]] + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [-f^*(g_f(D(G(\mathbf{z}|\mathbf{x})|\mathbf{x})))] ,$$

Where g and f are one the Frenchel conjugate of the other. They can be derived freely provided they are derived from a f -divergence. Thus, losses are:

$$\mathcal{L}_G(\theta_G) = \frac{-1}{N_{batch}} \sum_{i=1}^{N_{batch}} g_f(D(G(\mathbf{z}_i|\mathbf{x}_i; \theta_G)|\mathbf{x}_i; \theta_D)).$$

$$\mathcal{L}_D(\theta_D) = \frac{-1}{N_{batch}} \sum_{i=1}^{N_{batch}} [g_f(D(y_i|\mathbf{x}_i; \theta_D)) - f^*(g_f(D(G(\mathbf{z}_i|\mathbf{x}_i); \theta_G)|\mathbf{x}_i; \theta_D))] ,$$



Flavours of GANs

But what is a f-divergence?

An f-divergence is a concept of Information Theory and measures how different two probability distributions are. It actually measures how much difference of information there is.

The most remarkable f-divergence is the Kullback-Lieberman divergence.

$$D_{KL}(P, Q) = \int P(\mathbf{x}) \ln \left[\frac{P(\mathbf{x})}{Q(\mathbf{x})} \right] d\mathbf{x}.$$

Where P and Q are probability distributions. **BEWARE of non-commutativity!**

Giving result as

$$g_f(x) = x, \quad f^*(g_f(x)) = e^{x-1}.$$

Flavours of GANs

GANs can go beyond what has been stated before, just by changing the lossfunction.
Depending on the selected target function to be learned, a different property the GAN Will have

GAN type	Discriminator Loss	Generator Loss
Standard	$\mathcal{L}_D(\theta_D) = \frac{1}{N_{batch}} \sum_{i=1}^{N_{batch}} [\log D(\mathbf{x}_i, \theta_D) + \log(1 - D(G(\mathbf{z}_i, \theta_G), \theta_D))]$	$\mathcal{L}_G(\theta_G) = \frac{1}{N_{batch}} \sum_{i=1}^{N_{batch}} \log(1 - D(G(\mathbf{z}_i, \theta_G), \theta_D))$
KL-GAN	$\mathcal{L}(\theta_D) = \frac{-1}{N_{batch}} \sum_{i=1}^{N_{batch}} [D(\mathbf{x}_i, \theta_D) - e^{D(G(\mathbf{z}_i, \theta_G), \theta_D)} - 1]$	$\mathcal{L}(\theta_G) = \frac{-1}{N_{batch}} \sum_{i=1}^{N_{batch}} G(\mathbf{z}_i, \theta_G)$
WGAN	$\mathcal{L}_D = \frac{1}{N_{batch}} \left[\sum_{i=1}^{N_{batch}} D(G(\mathbf{z}_i, \theta_G), \theta_D) - \sum_{i=1}^{N_{batch}} D(\mathbf{x}_i, \theta_D) \right].$	$\mathcal{L}_G = -\frac{1}{N_{batch}} \sum_{i=1}^{N_{batch}} G(\mathbf{z}_i, \theta_G)$
KL-WGAN	$\mathcal{L}_D(\theta_D) = \frac{1}{N_{batch}} \sum_{i=1}^{N_{batch}} e^{r_i D(G(\mathbf{z}_i, \theta_G), \theta_D)} - \frac{1}{N_{batch}} \sum_{i=1}^{N_{batch}} e^{D(\mathbf{x}_i, \theta_D)}$	$\mathcal{L}_G(\theta_G) = \frac{1}{N_{batch}} \sum_{i=1}^{N_{batch}} r_i D(G(\mathbf{z}_i, \theta_G), \theta_D),$

Flavours of GANs

Based on the original formulation, the loss functions have evolved giving origin to different formulations.

- **Standard:** Original formulation by Goodfellow. Bad convergence. (<https://arxiv.org/abs/1406.2661>)
- **KL-GAN:** A generalization of the original. Better convergence but still hard. It has an added advantage, the output of the Generator fully matches a probability distribution. (<https://arxiv.org/abs/1606.00709>)
- **WGAN:** Alternative formulation. Faster and easier convergence, specially for images. But the output of the Generator can not be interpreted as a probability distribution (<https://arxiv.org/abs/1701.07875>)
- **KL-WGAN:** Generalizes both KL-GAN and WGAN. Gets the best of both worlds: improved convergence of the WGAN and the interpretability of a KL-GAN. (<https://arxiv.org/pdf/1910.09779>)

A large red square is the central focus. It is surrounded by several smaller squares with red outlines. To the left of the red square is a vertical stack of two squares. Above the red square is one square. To the right of the red square is a horizontal row of two squares, with one square positioned above the first square of that row. Further to the right is a single square. Below the red square is one square. In the bottom right corner of the slide is another single square.

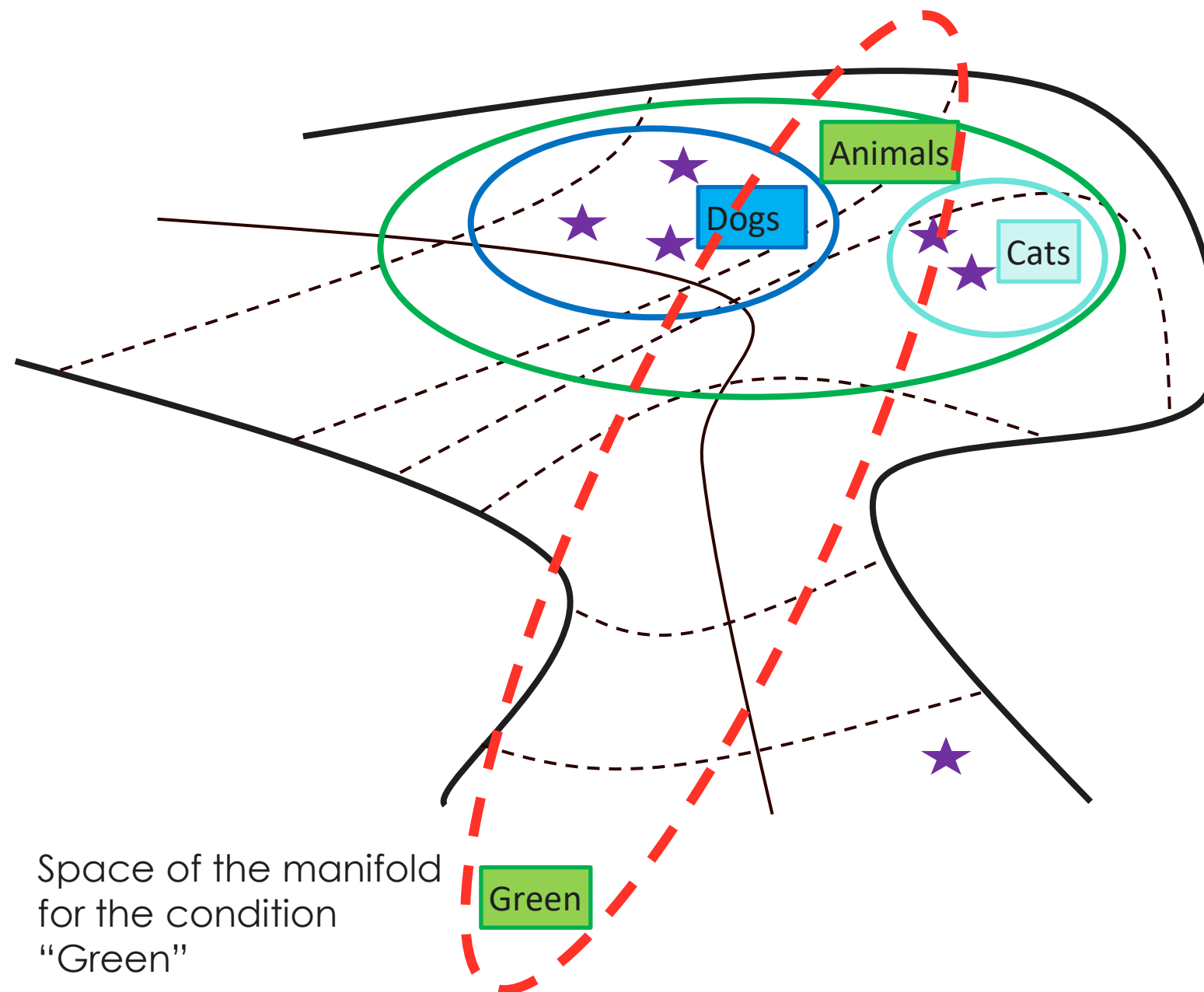
04

Conditional GANs

Conditional GANs

As seen previously, GANs will trace the underlying probability distribution of training data. The Conditional GANs (CGANs), trace the probability distribution given a condition. That is, they trace conditional probabilities.

At a conditional GAN, when sampling the latent-space, we will be only restricting ourselves to a specific region delimited by the condition.



Conditional GANs

CGANs have many different variants and applications in a wide range of different areas of application and different input type of data.

On one field CGANs have showcased extraordinary application is:

- Image synthesis and manipulation.
- Text-to-image generation.
- Image-to-image translation.
 - Image segmentation
 - Style transfer
 - Super-resolution and image enhancement.
- Speech enhancement.
- Music generation.

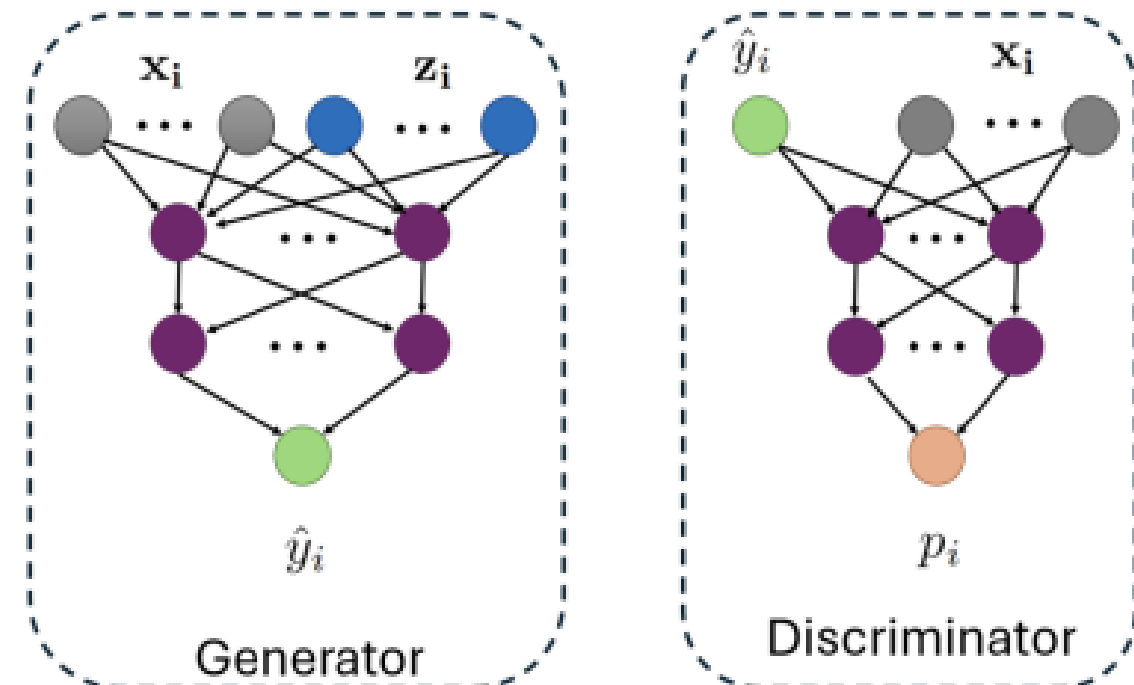


Conditional GANs

At CGANs, the typical losses are the same as in the case of normal GANs.

CGANs, also have a Discriminator and a Generator. Nevertheless, a slight modification shall be made to trace the conditional probability.

Now, Generator network gets both random noise and the input condition. Discriminator takes both the generated data and the condition.



Conditional GANs

Generator and discriminator are both feedforward neural networks.

The difference with Vanilla GAN is that input neurons take the latent-space dimension and the dimensions of the vector of conditions

Discriminator takes the image dimensions and the dimensions of the condition vector

```

64 # Generator
65 class Generator(nn.Module):
66     def __init__(self):
67         super(Generator, self).__init__()
68
69         self.model = nn.Sequential(
70             nn.Linear(z_dim + num_classes, 256),
71             nn.LeakyReLU(0.2),
72             nn.Linear(256, 512),
73             nn.LeakyReLU(0.2),
74             nn.Linear(512, 1024),
75             nn.LeakyReLU(0.2),
76             nn.Linear(1024, image_size),
77             nn.Tanh()
78         )
79
80     def forward(self, noise, labels):
81         label_onehot = torch.zeros(labels.size(0), num_classes).to(labels.device)
82         label_onehot.scatter_(1, labels.unsqueeze(1), 1)
83         input_tensor = torch.cat([noise, label_onehot], dim=1)
84         return self.model(input_tensor)
85
86 # Discriminator
87 class Discriminator(nn.Module):
88     def __init__(self):
89         super(Discriminator, self).__init__()
90
91         self.model = nn.Sequential(
92             nn.Linear(image_size + num_classes, 512),
93             nn.LeakyReLU(0.2),
94             nn.Dropout(0.3),
95             nn.Linear(512, 256),
96             nn.LeakyReLU(0.2),
97             nn.Dropout(0.3),
98             nn.Linear(256, 1),
99             nn.Sigmoid()
100         )
101
102     def forward(self, image, labels):
103         label_onehot = torch.zeros(labels.size(0), num_classes).to(labels.device)
104         label_onehot.scatter_(1, labels.unsqueeze(1), 1)
105         input_tensor = torch.cat([image, label_onehot], dim=1)
106         return self.model(input_tensor)

```

Conditional GANs

The only difference is that at CGAN, we now give the vector condition to the Discriminator when showing the real images

A random vector can be given as condition when showing the fake images. It is usually also given the real labels too.

```

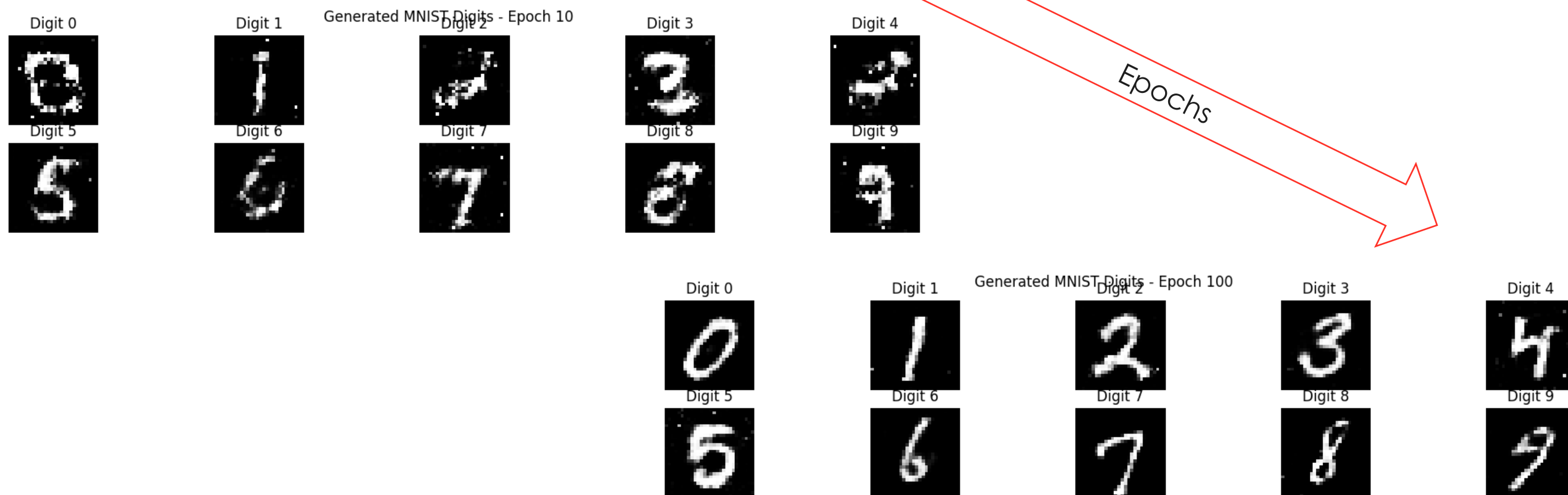
121 # Optimizers
122 optimizer_G = optim.Adam(generator.parameters(), lr=learning_rate, betas=(0.5, 0.999))
123 optimizer_D = optim.Adam(discriminator.parameters(), lr=learning_rate, betas=(0.5, 0.999))
124
125 # Loss function
126 criterion = nn.BCELoss()
127
128 # Lists to store losses
129 d_losses = []
130 g_losses = []
131
132 # Training
133 for epoch in range(num_epochs):
134     for i, (real_images, labels) in enumerate(dataloader):
135         real_images = real_images.view(-1, image_size).to(device)
136         labels = labels.to(device)
137
138         # Train Discriminator
139         optimizer_D.zero_grad()
140
141         # Real images
142         real_labels = torch.ones(real_images.size(0), 1).to(device)
143         real_output = discriminator(real_images, labels)
144         d_loss_real = criterion(real_output, real_labels)
145
146         # Fake images
147         noise = torch.randn(real_images.size(0), z_dim).to(device)
148         fake_labels = torch.randint(0, num_classes, (real_images.size(0),)).to(device)
149         fake_images = generator(noise, fake_labels)
150         fake_output = discriminator(fake_images.detach(), fake_labels)
151         fake_targets = torch.zeros(real_images.size(0), 1).to(device)
152         d_loss_fake = criterion(fake_output, fake_targets)
153
154         d_loss = d_loss_real + d_loss_fake
155         d_loss.backward()
156         optimizer_D.step()
157
158         # Train Generator
159         optimizer_G.zero_grad()
160
161         fake_output = discriminator(fake_images, fake_labels)
162         g_loss = criterion(fake_output, real_labels)
163         g_loss.backward()
164         optimizer_G.step()
165
166         # Store losses
167         d_losses.append(d_loss.item())
168         g_losses.append(g_loss.item())
169
170         if i % 100 == 0:
171             print(f'Epoch [{epoch+1}/{num_epochs}], Step [{i+1}/{len(dataloader)}], '
172                   f'D Loss: {d_loss.item():.4f}, G Loss: {g_loss.item():.4f}')
173

```

Conditional GANs

Example: MNIST image generation from input integer.

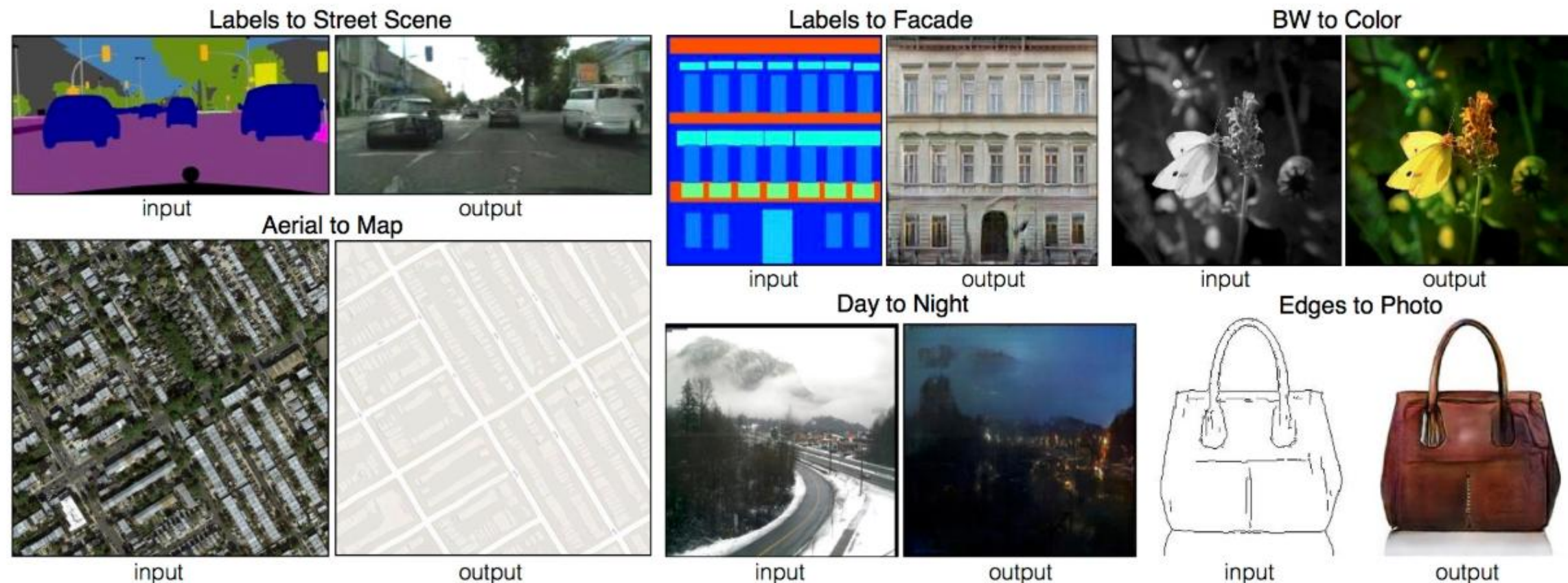
Now we have control over the generated number.



Conditional GANs

One interesting application is image-to-image translation, which have many interesting algorithms. One of which is pix2pix.

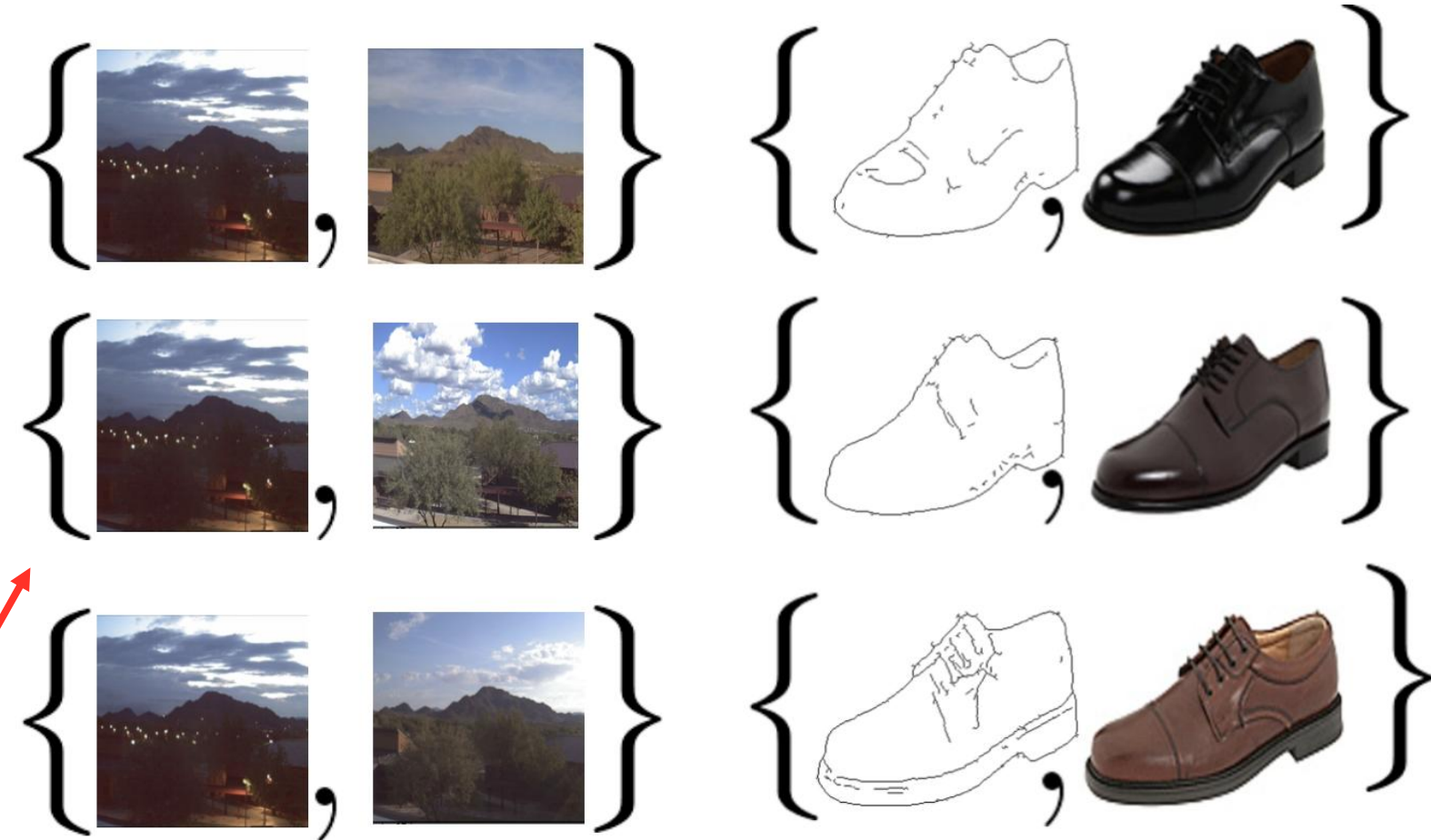
Pix2pix is a CGAN that enable image-to-image translation for pairs of images. That is, every image, will have its own counterpart.



Conditional GANs

The most important consideration to take into account for training pix2pix is that images shall be placed in pairs.

It is important to remark that it is not strictly necessary that they are 1-to-1, but can be 1-to-many.



Beware! One night scene has many daylight counterparts.



Conditional GANs

Pix2pix is composed by:

- Generator: a Unet.
- Discriminator: a PatchGAN.

```

129
130 class Discriminator(nn.Module):
131     def __init__(self):
132         super().__init__()
133         self.model = nn.Sequential(
134             nn.Conv2d(6, 64, 4, 2, 1), # Input + Output channels
135             nn.LeakyReLU(0.2, inplace=True),
136
137             nn.Conv2d(64, 128, 4, 2, 1, bias=False),
138             nn.BatchNorm2d(128),
139             nn.LeakyReLU(0.2, inplace=True),
140
141             nn.Conv2d(128, 256, 4, 2, 1, bias=False),
142             nn.BatchNorm2d(256),
143             nn.LeakyReLU(0.2, inplace=True),
144
145             nn.Conv2d(256, 512, 4, 1, 1, bias=False),
146             nn.BatchNorm2d(512),
147             nn.LeakyReLU(0.2, inplace=True),
148
149             nn.Conv2d(512, 1, 4, 1, 1)
150         )
151
152     def forward(self, input_img, target_img):
153         return self.model(torch.cat([input_img, target_img], 1))
154

```

```

69 class UNetBlock(nn.Module):
70     def __init__(self, in_channels, out_channels, down=True, use_dropout=False):
71         super().__init__()
72         if down:
73             self.conv = nn.Conv2d(in_channels, out_channels, 4, 2, 1, bias=False)
74         else:
75             self.conv = nn.ConvTranspose2d(in_channels, out_channels, 4, 2, 1, bias=False)
76
77         self.norm = nn.BatchNorm2d(out_channels)
78         self.relu = nn.LeakyReLU(0.2, inplace=True) if down else nn.ReLU(inplace=True)
79         self.dropout = nn.Dropout(0.5) if use_dropout else None
80
81     def forward(self, x):
82         x = self.conv(x)
83         x = self.norm(x)
84         x = self.relu(x)
85         if self.dropout:
86             x = self.dropout(x)
87         return x
88
89 class Generator(nn.Module):
90     def __init__(self):
91         super().__init__()
92         # Encoder
93         self.down1 = nn.Conv2d(3, 64, 4, 2, 1) # 128x128 -> 64x64
94         self.down2 = UNetBlock(64, 128, down=True) # 32x32
95         self.down3 = UNetBlock(128, 256, down=True) # 16x16
96         self.down4 = UNetBlock(256, 512, down=True) # 8x8
97         self.down5 = UNetBlock(512, 512, down=True) # 4x4
98         self.down6 = UNetBlock(512, 512, down=True) # 2x2
99         self.down7 = UNetBlock(512, 512, down=True) # 1x1
100
101         # Decoder
102         self.up1 = UNetBlock(512, 512, down=False, use_dropout=True)
103         self.up2 = UNetBlock(1024, 512, down=False, use_dropout=True)
104         self.up3 = UNetBlock(1024, 512, down=False, use_dropout=True)
105         self.up4 = UNetBlock(1024, 256, down=False)
106         self.up5 = UNetBlock(512, 128, down=False)
107         self.up6 = UNetBlock(256, 64, down=False)
108         self.final = nn.ConvTranspose2d(128, 3, 4, 2, 1)
109
110     def forward(self, x):
111         # Encoder
112         d1 = self.down1(x)
113         d2 = self.down2(d1)
114         d3 = self.down3(d2)
115         d4 = self.down4(d3)
116         d5 = self.down5(d4)
117         d6 = self.down6(d5)
118         d7 = self.down7(d6)
119
120         # Decoder with skip connections
121         u1 = self.up1(d7)
122         u2 = self.up2(torch.cat([u1, d6], 1))
123         u3 = self.up3(torch.cat([u2, d5], 1))
124         u4 = self.up4(torch.cat([u3, d4], 1))
125         u5 = self.up5(torch.cat([u4, d3], 1))
126         u6 = self.up6(torch.cat([u5, d2], 1))
127
128         return torch.tanh(self.final(torch.cat([u6, d1], 1)))
129

```

Conditional GANs

Procedure is similar as in MNIST CGANT with integers as condition. Just now the condition is another image.

In this case, no random vector is used. UNET directly can map from original image to the other image.

PatchGAN can take the two images (real and generated) and see if they match

```

176 generator = Generator().to(device)
177 discriminator = Discriminator().to(device)
178
179 # Loss functions
180 criterion_GAN = nn.BCEWithLogitsLoss()
181 criterion_L1 = nn.L1Loss()
182
183 # Optimizers
184 optimizer_G = optim.Adam(generator.parameters(), lr=0.0002, betas=(0.5, 0.999))
185 optimizer_D = optim.Adam(discriminator.parameters(), lr=0.0002, betas=(0.5, 0.999))
186
187 # Training parameters
188 num_epochs = 100
189 lambda_L1 = 100
190
191 print("Starting training...")
192 print(f"Generator parameters: {sum(p.numel() for p in generator.parameters())}")
193 print(f"Discriminator parameters: {sum(p.numel() for p in discriminator.parameters())}")
194
195
196 loss_g = []
197 loss_d = []
198
199 # Training loop
200 for epoch in range(num_epochs):
201     for i, (input_imgs, target_imgs) in enumerate(train_dataloader):
202         input_imgs = input_imgs.to(device)
203         target_imgs = target_imgs.to(device)
204
205         # Train Generator
206         optimizer_G.zero_grad()
207         fake_imgs = generator(input_imgs)
208
209         # GAN loss
210         pred_fake = discriminator(input_imgs, fake_imgs)
211         real_labels = torch.ones_like(pred_fake)
212         loss_GAN = criterion_GAN(pred_fake, real_labels)
213
214         # L1 loss
215         loss_L1 = criterion_L1(fake_imgs, target_imgs)
216
217         # Total generator loss
218         loss_G = loss_GAN + lambda_L1 * loss_L1
219         loss_G.backward()
220         optimizer_G.step()
221
222         # Train Discriminator
223         optimizer_D.zero_grad()
224
225         # Real images
226         pred_real = discriminator(input_imgs, target_imgs)
227         real_labels = torch.ones_like(pred_real)
228         loss_real = criterion_GAN(pred_real, real_labels)
229
230         # Fake images
231         pred_fake = discriminator(input_imgs, fake_imgs.detach())
232         fake_labels = torch.zeros_like(pred_fake)
233         loss_fake = criterion_GAN(pred_fake, fake_labels)
234
235         # Total discriminator loss
236         loss_D = (loss_real + loss_fake) * 0.5
237         loss_D.backward()
238         optimizer_D.step()
239

```

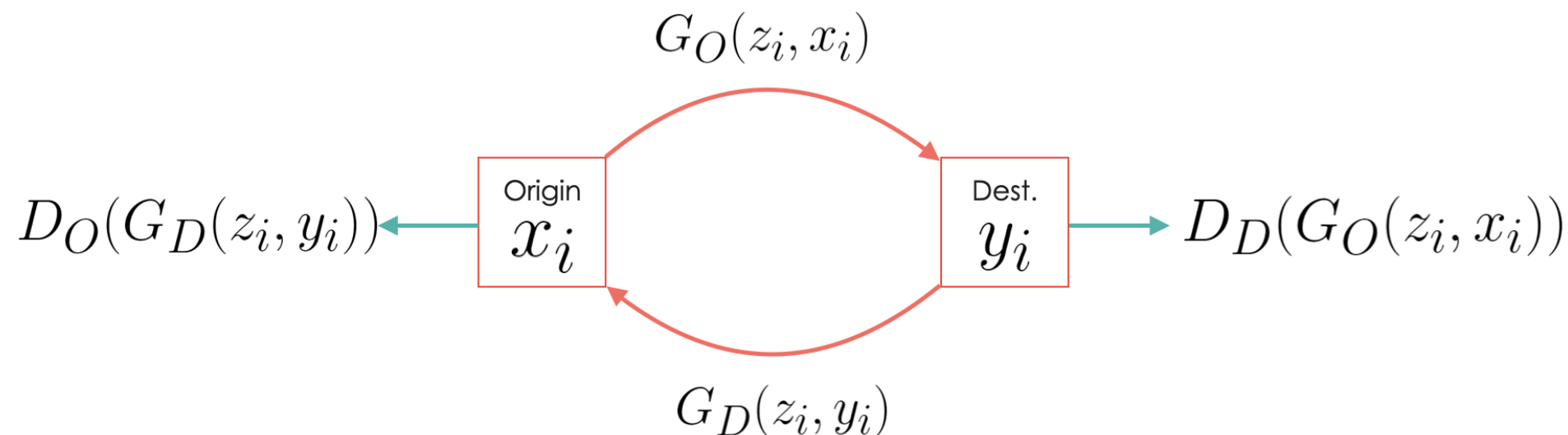

Conditional GANs

Nevertheless, there are cases where it is impossible to arrange images in image-pairs.

In this scenario, is where Cycle-GAN appears.

Cycle-GAN consist on 4 networks: 2 discriminators and 2 generators, all of them competing with each other.

In this case, images are arranged on just two sets: the origin set and the destination set.



Conditional GANs

While Diffusion Models have replaced largely GANs at general image generation because Diffusion models are much more stable and produce more realistic images, **CGAN still remain the state-of-the art for style transfer and image-to-image translation.**

This is because style-transfer is aligned better with overall CGAN objectives, because:

- Style-transfer is a conditional problem (like CGANs).
- PatchGAN stands for a top performer at Discriminator network because ensures both local and global realism.
- Most CGAN architectures were optimized for image-to-image translation, while diffusion architectures where optimized for text-to-image.
- Mode-collapse is what you aim when doing image-to-image translation (i.e. good realistic counterpart, rather the average between different modes).
- Image-to-image translation is just a guided transformation. Diffusion models are intended for image generation, they are an overkill for the task.
- Cycle consistency of Cycle-GAN ensures visual coherence across images and invertibility across transformed spaces (which Diffusion models do not have).





**Universidad
Europea**

Thanks

Dr. Manuel García Fernández

manuel.garcia2@universidadeuropea.es

Ve más allá