Generative Adversarial Network

Deep Learning (ENGG*6600*07)

School of Engineering, University of Guelph, ON, Canada

Course Instructor: Benyamin Ghojogh Fall 2023 Introduction

Introduction

- Suppose we have a generative model which takes a random noise as input and generates a data point.
- We want the generated data point to be of good quality; hence, we should somehow judge its quality.
- One way to judge it is to observe the generated sample and assess its quality visually. In this case, the judge is a human. However, we cannot take derivative of human's judgment for optimization.
- Generative Adversarial Network (GAN), proposed in (2014) [1], has the same idea but it can take derivative of the judgment. For that, it uses a classifier as the judge rather than a human. Hence, we have a generator generating a sample and a binary classifier (or discriminator) to classify the generated sample as a real or generated sample.
- This classifier can be a <u>pre-trained network</u> which is already trained by some real and generated (fake) data points. However, GAN puts a step ahead and lets the classifier be <u>trained simultaneously</u> with training the generator. This is the core idea of adversarial learning where the classifier, also called the discriminator, and the generator compete with one another; hence, they make each other stronger gradually by this competition [2].

Introduction

- It is noteworthy that the term "adversarial" is used in two main streams of research in machine learning and they should not be confused. These two research areas are:
 - Adversarial attack, also called learning with adversarial examples or adversarial machine learning. This line of research inspects some examples which can be changed slightly but wisely to fool a trained learning model. For example, perturbation of some specific pixels in the input image may change the decision of learning model. The reason for this can be analyzed theoretically. Some example works in this area are [3, 4, 5, 6, 7].
 - Adversarial learning for generation. This line of research is categorized as generative models [8] and/or methods based on that. GAN is in this line of research. This paper focuses on this research area.

Adversarial Learning: The Adversarial Game

Adversarial Learning: The Adversarial Game

- The original GAN, also called the vanilla GAN, was proposed in (2014) [1].
- Consider a *d*-dimensional dataset with <u>*n*</u> data points, i.e., $\{x_i \in \mathbb{R}^d\}_{i=1}^n$. In GAN, we have a generator *G* which takes a *p*-dimensional random noise $z \in \mathbb{R}^p$ as input and outputs a *d*-dimensional generated point $x \in \mathbb{R}^d$. Hence, it is the mapping $\underline{G}: z \to x$ where:

$$\underbrace{G(z)}_{}=\underbrace{x.}$$
(1)

Let the distribution of random noise be denoted by z ~ p_z(z). We want the generated x̂ to be very similar to some original (or real) data point x in the dataset. We need a module to judge the quality of the generated point to see how similar it is to the real point. This module can be a human but we cannot take derivative of human's judgment for optimization! A good candidate for the judge is a classifier, also called the discriminator. The discriminator (also called the critic), denoted by D : x → [0, 1], is a binary classifier which classifies the generated point as a real or generated point:

$$D(\mathbf{x}) := \begin{cases} \frac{1}{0} & \text{if } \mathbf{x} \text{ is real.} \\ \frac{1}{0} & \text{if } \mathbf{x} \text{ is generated (fake).} \end{cases}$$
(2)

- The perfect discriminator outputs one for real points and zero for generated points.
- The discriminator's output is in the range [0, 1] where the output for real data is closer to one and the output for fake data is closer to zero.
- If the generated point is very good and closely similar to a real data point, the classifier may make a mistake and output a value close to one for it. Therefore, if the classifier makes a mistake for the generated point, the generator has done a good job in generating a data point.

Adversarial Learning: The Adversarial Game

- The discriminator can be a <u>pre-trained classifier</u>, but we can make the problem more sophisticated, as explained in the following.
- Let us train the discriminator simultaneously while we are training the generator. This makes the discriminator *D* and the generator *G* stronger gradually while they compete each other.
- On the one hand, the generator tries to generate realistic points to fool the discriminator and make it a hard time to distinguish the generated point from a real point.
- On the other hand, the discriminator tries to discriminate the fake (i.e., generated) point from a real point.
- When one of them gets stronger in training, the other one tries to become stronger to be able to compete. Therefore, there is an **adversarial game** between the generator and the discriminator. This game is **zero-sum** because whatever one of them loses, the other wins.



Optimization and Loss Function

Optimization and Loss Function

- We denote the probability distributions of dataset and noise by $p_{data}(x)$ and $p_z(z)$, respectively.
- As the figure shows, the discriminator is trained by real points from dataset as well as generated points from the generator.
- The discriminator and generator are trained simultaneously.
- The optimization loss function for both the discriminator and generator is:

$$\underbrace{\operatorname{mir}}_{\mathbb{G}} \underbrace{V(D,G)}_{\mathcal{D}} := \mathbb{E}_{\mathbf{x} \sim p_{\operatorname{data}}(\mathbf{x})} \Big[\log(D(\mathbf{x})) \Big] + \mathbb{E}_{\mathbf{z} \sim p_{z}(z)} \Big[\log(1 - D(G(z))) \Big],$$
(3)

where $\mathbb{E}[.]$ denotes the expectation operator and the loss function V(D, G) is also called the value function of the game.



• In practice, we can use the Monte Carlo approximation [9] of expectation where the expectations are replaced with averages over the mini-batch.

Optimization and Loss Function • We had: $\min_{G} D = V(D,G) := \mathbb{E}_{x \sim p_{data}(x)} \left[\log(D(x)) \right] + \mathbb{E}_{z \sim p_z(z)} \left[\log\left(1 - D(G(z))\right) \right].$

- The first term in Eq. (3) is expectation over the real data. This term is only used for the discriminator while it is a constant for the generator. According to Eq. (2), D(x) outputs one (the larger label) for the real data; therefore, the discriminator maximizes this term because it assigns the larger label to the real data.
- The second term in Eq. (3) is expectation over noise. It inputs the noise z to the generator to have G(z). The output of generator, which is the generated point, is fed as input to the discriminator to have D(G(z)). The discriminator wants to minimize D(G(z)) because the smaller label is assigned to the generated data, according to Eq. (2). In other words, the discriminator wants to maximize 1 − D(G(z)).
- As logarithm is a monotonic function, we can say that the discriminator wants to maximize E_{z~p_z(z)}[log(1 D(G(z)))] which is the second term in Eq. (3). As opposed to the discriminator, the generator minimizes E_{z~p_z(z)}[log(1 D(G(z)))] which is the second term in Eq. (3). This is because the generator wants to fool the discriminator to label the generated data as real data.



Optimization and Loss Function

We had:

$$\underbrace{\min_{G} \max_{D}}_{W(D,G)} V(D,G) := \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})} \Big[\log(D(\mathbf{x})) \Big] + \mathbb{E}_{\mathbf{z} \sim p_{z}(\mathbf{z})} \Big[\log\Big(1 - D(G(\mathbf{z}))\Big) \Big].$$

- The Eq. (3) is a <u>minimax optimization problem</u> [10] and can be solved using <u>alternating</u> <u>optimization</u> [11] where we optimize over \overline{D} and over G iteratively <u>until convergence</u> (i.e., Nash equilibrium).
- The original GAN [1] uses a step of stochastic gradient descent [11] for updates of each variable in the alternating optimization. If we denote the loss function in Eq. (3) by V(D, G), the alternating optimization is done as:

$$\longrightarrow \begin{cases} & \underbrace{D^{(k+1)}} := \underbrace{D^{(k)}}_{O} \eta^{(k)} \frac{\partial}{\partial D} \left(V(D, G^{(k)}) \right), \end{cases}$$
(4)

$$\longrightarrow \int \mathbf{K} \ G^{(k+1)} := \underline{G^{(k)}} - \eta^{(k)} \frac{\partial}{\partial G} \Big(V(\underline{\widehat{D^{(k+1)}}}, \underline{G}) \Big), \tag{5}$$

where k is the index of iteration and $\eta^{(k)}$ is the learning rate at iteration k.

- Throughout this slide deck, **derivatives w.r.t.** *D* and *G* mean the derivatives w.r.t. the parameters (weights) of *D* and *G* networks, respectively.
- Eqs. (4) and (5) are one step of gradient ascent and gradient descent, respectively. Note that the gradients here are the average of gradients in the mini-batch. Every mini-batch includes both real and generated data.
- The <u>paper [1]</u> suggests that Eq. (4) can be <u>performed for several times</u> before performing Eq. (5); however, the experiments of that paper perform Eq. (4) for only one time before performing Eq. (5).

Optimization and Loss Function

- Also note that another way to solve the optimization problem in GAN is <u>simultaneous</u> optimization [12] in which Eqs. (4) and (5) are performed at the same time and not one after the other.
- Minimax versus maximin in GAN [13, Section 5]: We saw in Eq. (3) that the optimization of GAN is a minimax problem:

In fact, under some conditions, Eqs. (6) and (7) are equivalent [10].

Network Structure of GAN

- In practice, the discriminator and generator are two (deep) neural networks. The structure of GAN is depicted in this figure.
- The first layer of discriminator network is d-dimensional and its last layer is one dimensional with scalar output. In the original GAN, maxout activation function [14] is used for all layers except the last layer which has the sigmoid activation function to output a probability to model Eq. (2). The closer the output of D to one, the more probable its input is to be real.
- The generator network has a *p*-dimensional input layer for noise and a *d*-dimensional output layer for generating data. In the generator, a combination of <u>ReLU</u> [15] and sigmoid activation functions are used.
- The space of noise as the input to the generator is called the **latent space** or the **latent factor**.
- Each of the Eqs. (4) and (5) are performed using **backpropagation** in the neural networks.



• Theorem: [1, Proposition 1]: For a fixed generator G, the optimal discriminator is:

where $p_{data}(x)$ is the probability distribution of the real dataset evaluated at point x and $p_g(x)$ is the probability distribution of output of generator evaluated at point x.

• **Proof**: According to the definition of expectation, the loss function in Eq. (3) can be stated as:

$$\bigstar V(D,G) = \underbrace{\int_{x}}_{x} p_{data}(x) \log(D(x)) dx + \underbrace{\int_{z}}_{z} p_{z}(x) \log(1 - D(G(z))) dz.$$

• According to Eq. (1), we have: $\underbrace{G(z) = x}_{G(z) = x} \Longrightarrow \underbrace{(z = G^{-1}(x))}_{Z = G^{-1}(x)} \Longrightarrow \underbrace{dz = (G^{-1})'(x)dx}_{G(z), X}$

where $(G^{-1})'(x)$ is the derivative of $(G^{-1})(x)$ with respect to (w.r.t.) x. Hence:

$$V(D,G) = \int_{\mathbf{x}} p_{data}(\mathbf{x}) \log(D(\mathbf{x})) d\mathbf{x} + \int_{\mathbf{z}} p_z(\underline{G^{-1}(\mathbf{x})}) \log(1 - D(\underline{\mathbf{x}}))(\underline{G^{-1}})'(\mathbf{x}) d\mathbf{x}.$$

• We found: • We found: • $V(D, G) = \int_{x} p_{data}(x) \log(D(x)) dx + \int_{z} o_{z}(G^{-1}(x)) og(1 - D(x)) (G^{-1})'(x) dx.$ • The relation of distributions of input and output of generator is: • $p_{g}(x) = p_{z}(z) \times G^{-1}(x) = p_{z}(G^{-1}(x)) G^{-1}(x),$ (9) where $G^{-1}(x)$ is the Jacobian of distribution at point x. Hence: • $V(D, G) = \int_{x} p_{data}(x) \log(D(x)) dx + \int_{Q} p_{g}(x) \log(1 - D(x)) dx$

$$= \int_{\mathbf{x}} \left(p_{data}(\mathbf{x}) \log(D(\mathbf{x})) + p_g(\mathbf{x}) \log(1 - D(\mathbf{x})) \right) d\mathbf{x}.$$
(10)

• For optimization in Eq. (3), taking derivative w.r.t. D(x) gives:

$$\frac{\partial V(D,G)}{\partial D(\mathbf{x})} \stackrel{(a)}{=} \frac{\partial}{\partial D(\mathbf{x})} \left(p_{data}(\mathbf{x}) \log(D(\mathbf{x})) + p_g(\mathbf{x}) \log(1 - D(\mathbf{x})) \right) \\ = \underbrace{p_{data}(\mathbf{x})}_{D(\mathbf{x})} - \underbrace{p_g(\mathbf{x})}_{1 - D(\mathbf{x})} \stackrel{=}{=} \underbrace{p_{data}(\mathbf{x})(1 - D(\mathbf{x})) - p_g(\mathbf{x})D(\mathbf{x})}_{D(\mathbf{x})(1 - D(\mathbf{x}))} \stackrel{\text{set }}{=} \underbrace{0} \\ \implies p_{data}(\mathbf{x}) - p_{data}(\mathbf{x})D(\mathbf{x}) - p_g(\mathbf{x})D(\mathbf{x}) = 0 \implies \left(D(\mathbf{x}) = \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})}, \right)$$

where (a) is because taking derivative w.r.t. D(x) considers a specific x and hence it removes the integral (summation). Q.E.D.

• **Theorem**: [1, Theorem 1] The optimal solution of GAN is when the distribution of generated data becomes equal to the distribution of data:

$$\oint p_{g^*}(\mathbf{x}) = p_{data}(\mathbf{x}). \tag{11}$$

• **Proof**: Putting the optimum $D^*(x)$, i.e. Eq. (8), in Eq. (10) gives:

$$\begin{array}{l} \checkmark \quad V(\widehat{D^*}, G) = \int_{x} \left(p_{\text{data}}(x) \log(\widehat{D^*(x)}) + p_g(x) \log(1 - D^*(x)) \right) dx \\ \stackrel{(8)}{=} \int_{x} \left[p_{\text{data}}(x) \log\left(\frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)}\right) + p_g(x) \log\left(\frac{p_g(x)}{p_{\text{data}}(x) + p_g(x)}\right) \right] dx \\ = \int_{x} \left[p_{\text{data}}(x) \log\left(\frac{p_{\text{data}}(x)}{2}\right) + p_g(x) \log\left(\frac{p_g(x)}{2}\right) \right] dx \\ = \int_{x} \left[p_{\text{data}}(x) \log\left(\frac{p_{\text{data}}(x)}{2}\right) + p_g(x) \log\left(\frac{p_g(x)}{2}\right) \right] dx \\ = \int_{x} \left[p_{\text{data}}(x) \log\left(\frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)}\right) + p_g(x) \log\left(\frac{p_g(x)}{p_{\text{data}}(x) + p_g(x)}\right) \right] dx \\ = \int_{x} \left[p_{\text{data}}(x) \log\left(\frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)}\right) + p_g(x) \log\left(\frac{p_g(x)}{p_{\text{data}}(x) + p_g(x)}\right) \right] dx \\ = \int_{x} \left[p_{\text{data}}(x) \log\left(\frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)}\right) + p_g(x) \log\left(\frac{p_g(x)}{p_{\text{data}}(x) + p_g(x)}\right) \right] dx \\ = \int_{x} \left[p_{\text{data}}(x) \log\left(\frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)}\right) + p_g(x) \log\left(\frac{p_g(x)}{p_{\text{data}}(x) + p_g(x)}\right) \right] dx \\ = \int_{x} \left[p_{\text{data}}(x) \log\left(\frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)}\right) + p_g(x) \log\left(\frac{p_g(x)}{p_{\text{data}}(x) + p_g(x)}\right) \right] dx \\ = \int_{x} \left[p_{\text{data}}(x) \log\left(\frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)}\right) + VL\left(p_g(x)\right) \right] \frac{p_{\text{data}}(x) + p_g(x)}{2} \right) - \log(4), \quad (12) \end{aligned} \right]$$

where (a) is because of the definition of KL divergence.

• We found: $V(D^*, G) = \mathsf{KL}\left(p_{\mathsf{data}}(\mathbf{x}) \| \frac{p_{\mathsf{data}}(\mathbf{x}) + p_g(\mathbf{x})}{2}\right) + \mathsf{KL}\left(p_g(\mathbf{x}) \| \frac{p_{\mathsf{data}}(\mathbf{x}) + p_g(\mathbf{x})}{2}\right) - \log(4).$

• The Jensen-Shannon Divergence (JSD) is defined as [16]:

$$\underbrace{JSD(P\|Q)}_{=} = \underbrace{\frac{1}{2}}_{=} KL(P\|\frac{1}{2}(P+Q)) + \frac{1}{2}KL(Q\|\frac{1}{2}(P+Q)), \quad (13)$$

where P and Q denote the probability densities. In contrast to KL divergence, the JSD is symmetric.

- The obtained $V(D^*, G)$ can be restated as: $V(D^*, G) = 2 JSD(p_{data}(x) || p_g(x)) - \log(4).$ (14)
- According to Eq. (3), the generator minimizes $V(D^*, G)$. As the JSD is non-negative, the above loss function is minimized if we have:

$$\underbrace{JSD(p_{data}(\mathbf{x}) \parallel p_{g^*}(\mathbf{x})) = 0}_{JSD(p_{data}(\mathbf{x}) \parallel p_{g^*}(\mathbf{x})) = 0} \implies p_{data}(\mathbf{x}) = p_{g^*}(\mathbf{x}).$$

Q.E.D.

• Corollary: [1, Theorem 1]: From Eqs. (11) and (14):

$$\begin{cases} p_{g^*}(\mathbf{x}) = p_{data}(\mathbf{x}), \\ V(D^*, G) = 2 \operatorname{JSD}(p_{data}(\mathbf{x}) || p_g(\mathbf{x})) - \log(4), \end{cases}$$

we conclude that the optimal loss function in GAN is:

$$V(D^*, G^*) = -\log(4). \tag{15}$$

Convergence and Equilibrium Analysis of GAN

Convergence and Equilibrium Analysis of GAN

• **Theorem**: [1, Proposition 2] If the <u>discriminator</u> and <u>generator</u> have enough capacity and, at every iteration of the alternating optimization, the followings occur:

the discriminator is allowed to reach its optimum value as in Eq. (8),
$$\underbrace{D^*(x) = \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)}}_{p_{\text{data}}(x) + p_g(x)} \text{ and}$$
*p*_g(*x*) is updated to minimize *V*(*D**, *G*), as stated in Eq. (14),
$$V(D^*, G) = 2 \operatorname{JSD}(p_{\text{data}}(x) \parallel p_g(x)) - \log(4),$$
then, $p_g(x)$ converges to $p_{\text{data}}(x)$, as stated in Eq. (11), $p_g^*(x) = p_{\text{data}}(x)$.

$$D^{\dagger}(x) = \frac{1}{2}$$

Convergence and Equilibrium Analysis of GAN

• Proof: The KL divergences in Eq. (12):

$$\bigstar' V(D^*, G) = \mathsf{KL}\left(p_{\mathsf{data}}(\boldsymbol{x}) \| \frac{p_{\mathsf{data}}(\boldsymbol{x}) + p_g(\boldsymbol{x})}{2}\right) + \mathsf{KL}\left(p_g(\boldsymbol{x}) \| \frac{p_{\mathsf{data}}(\boldsymbol{x}) + p_g(\boldsymbol{x})}{2}\right) - \log(4).$$

are convex functions w.r.t. $p_g(\mathbf{x})$. Hence, with sufficiently small updates of $p_g(\mathbf{x})$, it converges to $p_{data}(\mathbf{x})$. Note that Eq. (12), which we used here, holds if Eq. (8) holds, i.e., the discriminator is allowed to reach its optimum value. Q.E.D.

• The GAN loss, i.e. Eq. (3), can be restated as [17]:

$$\underset{G}{\min \max} V(D,G) := \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})} \Big[f(D(\mathbf{x})) \Big] + \mathbb{E}_{\mathbf{z} \sim p_{z}(\mathbf{z})} \Big[f(-D(G(\mathbf{z}))) \Big],$$
(16)

where <u>f</u> is the negative logistic function, i.e., $f(x) := -\log(1 + \exp(-x))$. In fact, the function f(.) can be any concave function. This formulation is slightly different from the original GAN in the sense that, here, the discriminator D outputs a real-valued scalar (without any activation function) while the discriminator of Eq. (3) outputs values in the range (0, 1) after a sigmoid activation function.

If D outputs 0.5 and 0, it means that it is completely confused in Eqs. (3) and (16), respectively. The Eq. (16) is a concave-concave loss function in most of the domain of discriminator [17, Proposition 3.1].

Convergence and Equilibrium Analysis of GAN

- **Theorem**: [17, Theorem 3.1]: After satisfying several reasonable assumptions (see [17] for details), a GAN with loss function of Eq. (16) is locally exponentially stable.
- Lemma: Nash equilibrium in GAN [18]: Nash equilibrium is the state of game where no player can improve its gain by choosing a different strategy. At the Nash equilibrium of GAN, we have:

$$\underbrace{V(D, G^*)}_{V(D, G^*)} \underbrace{\leq V(D^*, G^*)}_{V(D^*, G^*)} \underbrace{\leq V(D^*, G)}_{V(D^*, G)},$$
(17)

which is obvious because we are minimizing and maximizing V(G, D) by the generator and discriminator, respectively, in Eq. (3):

$$\min_{G} \max_{D} V(D,G).$$

 Empirical experiments have shown that GAN <u>may not reach its Nash equilibrium in</u> practice [18]. Regularization can help convergence of GAN to the Nash equilibrium [19]. It is shown in [19] that an effective regularization for GAN is <u>noise injection</u> [20] in which independent Gaussian noise is added to the training data points.

Conditional GAN

Conditional GAN

- As was explained before, in the original GAN, we randomly draw noise z from a prior distribution p_z(z) and feed it to the generator. The generator outputs a point x from the noise z. Assume that the dataset with which GAN is trained has c number of classes. The original GAN generates points from any class and we do not have control to generate a point from a specific class.
- Although, note that after GAN is trained, the latent space for *z* is meaningful in the sense that every part of the latent space results in generation of some specific points from some class. However, the user cannot choose specifically what class to generate points from.
- Conditional GAN (2014) [21], also called the conditional adversarial network, gives the user the opportunity to choose the class of generation of points. For the dataset {*x_i* ∈ ℝ^d}ⁿ_{i=1}, let the one-hot encoded class labels be {*y_i* ∈ ℝ^c}ⁿ_{i=1}. In conditional GAN, we use the following loss function instead of Eq. (3):

$$\mathbf{x} \quad \min_{G} \max_{D} \quad V_{C}(D,G) := \mathbb{E}_{\mathbf{x} \sim \rho_{data}(\mathbf{x})} \Big[\log \big(D(\mathbf{x}|\mathbf{y}) \big) \Big] + \mathbb{E}_{z \sim \rho_{z}(z)} \Big[\log \big(1 - D\big(G(z|\mathbf{y}) \big) \big) \Big],$$
(18)

where the discriminator and generator are both conditioned on the labels.

- In practice, for implementing the loss function (18), we concatenate the one-hot encoded label y to the point x for the input to discriminator. We also concatenate the one-hot encoded label y to the noise z for the input to generator. For these, the input layers of discriminator and generator are enlarged to accept the concatenated inputs.
- In the test phase, the user chooses the desired class label and the generator generates a new point from that class.

Deep Convolutional GAN (DCGAN)

Deep Convolutional GAN (DCGAN)

- Deep Convolutional GAN (DCGAN), proposed in (2016) [22], made GAN deeper and generated higher resolution images than GAN.
- It also showed that it is very hard to train a GAN.
- In DCGAN, we use an <u>all-convolutional network</u> [23] which <u>replaces the pooling</u> <u>functions with strided convolutions</u>. In this way, the network <u>learns its own spatial</u> <u>downsampling</u>. This network is used for both generator and discriminator.
- In DCGAN, we also have only convolutional layers in the input layer of generator and output layer of discriminator, without any fully-connected layer. This elimination of fully connected layers is inspired by [24].
- We also apply <u>batch normalization</u> [25] to all layers <u>except the last layer of generative</u> and the first layer of discriminator. This is because batch normalization makes the mean of input to each neuron zero and its variance one; however, we should learn the mean and variance of data in the first layer of discriminator and the mean and variance of data should be reproduced by the last layer of generator.
- Batch normalization reduces the problem of mode collapse, which will be introduced later, with the price of causing some fluctuations and instability [22].
- Virtual batch normalization [26, Section 3.5]: Batch normalization has a problem; it makes the effect of every input x on network dependent on other inputs in the mini-batch. To not have this problem, Virtual Batch Normalization (VBN) fixes the mini-batches initially once before start of training; these mini-batches are called the reference batches. Every reference batch is normalized by only its own statistics (i.e., mean and covariance). VBN has been found to be effective in training of generator [26].

Deep Convolutional GAN (DCGAN)

- In DCGAN, the last layer of generator has the hyperbolic tangent activation function and its other layers have the ReLU activation function [15].
- As in GAN, the one-to-last layer of discriminator is flattened and connected to one neuron with the sigmoid activation function.
- In contrast to GAN, which uses the <u>maxout activation function</u> [14] for discriminator layers, DCGAN uses the **leaky rectified action function** for discriminator.

Vector Arithmetic in Latent Space

- DCGAN showed that we can generate images from a specific domain if we train GAN on that domain. For example, bedroom images were generated by DCGAN after being trained on a dataset of bedroom images.
- DCGAN also showed that the learned latent space is meaningful and we can do vector arithmetic in the latent space. Vector arithmetic in the latent space was previously used for showing the ability of <u>Word2Vec</u> [27]. DCGAN made it possible to do vector arithmetic in the latent space for images.
- An example of vector arithmetic by DCGAN is shown in this figure. In the latent space, the latent variables corresponding to man with glasses, man without glasses, and woman without glasses are used. Using one latent point for each of these does not work very well. An average of three latent vectors for each has worked properly in practice. The vector arithmetic works because "man with glasses" minus "man without glasses" plus "woman without glasses" results in "woman with glasses".



king - mon + woman >> gueen

Mode Collapse Problem in GAN

Mode Collapse Problem in GAN

- We expect from a GAN to learn a <u>meaningful latent space of z</u> so that <u>every specific</u> value of z maps to a specific generated data point <u>x</u>. Also, nearby <u>z</u> values in the latent space should be mapped to similar but a little different generations.
- The <u>mode collapse problem</u> [28], also known as the <u>Helvetica scenario</u> [13], is a common problem in GAN models. It refers to when the <u>generator cannot learn a perfectly</u> <u>meaningful latent space</u> as was explained. Rather, it learns to <u>map several different z</u> values to the same generated data point.
- Mode collapse usually happens in GAN when the distribution of training data, <u>p_{data}(x)</u>, has <u>multiple modes</u>.
- An example of mode collapse is illustrated in this figure which shows training steps of a GAN model when the training data is a <u>mixture of Gaussians</u> [28]. In different training steps, GAN learns to map all z values to one of the modes of mixture. When the discriminator learns to reject generation of some mode, the generator learns to map all z values to another mode. However, it never learns to generate all modes of the mixture. We expect GAN to map some part, and not all parts, of the latent space to one of the modes so that all modes are covered by the whole latent space.



Mode Collapse Problem in GAN

- Another statement of the mode collapse is as follows [29, Fig. 1]. Assume $\underline{p_{data}(x)}$ is **<u>multi-modal</u>** while the latent space $p_z(z)$ has only **one mode**.
- Consider two points x_1 and x_2 from two modes of data whose corresponding latent noises are z_1 and z_2 , respectively. According to the **mean value theorem**, there is a latent noise with the **absolute gradient value** $||x_2 x_1||/||\overline{z_2 z_1}||$ where ||.|| is a norm. As this gradient is Lipschitz continuous, when the two modes are very far resulting in a large $||\overline{x_2 x_1}||$, we face a problem. In this case, the latent noises between z_1 and z_2 generate data points between x_1 and x_2 which are not in the modes of data and thus are not valid.
- There exist various methods which resolve the mode collapse problem in GAN and adversarial learning. Some of them make the latent <u>space a mixture distribution</u> to imitate generation of the multi-modal training data. Some of them, however, have other approaches.
- These methods are <u>Minibatch GAN</u> (2016) [26, Section 3.2], <u>Unrolled GAN</u> (2017) [28], <u>Bourgain GAN</u> (2018) [29], <u>Mixture GAN (MGAN)</u> (2018) [30], <u>Dual Discriminator GAN</u> (D2GAN) (2017) [31], and Wasserstein GAN (WGAN) (2017) [32].
- For more information on GAN and its variants, see our tutorial [33].

Some Applications of GAN



Some Applications of GAN



This small bird has a white breast, light grey head, and black wings and tail.

stubby beak.



Figure 10. Text-to-image translation by <u>StackGAN</u>. Images are from (Zhang et al., 2017).

credit: [36]

Some Applications of GAN



Figure 11. Mixing image characteristics using GAN: (a) Generating an image with background, shape, and color characteristics by FineGAN, (b) generating an image by borrowing its characteristics from three images using <u>MixNMatch</u>, and (c) generating an image by borrowing its characteristics from different domains using improved MixNMatch. Images are from (Singh et al., 2019), (<u>Li et al.</u>, 2020), and (Ojha et al., 2021b), respectively.

Images are from [37], [38], and [39], respectively.

- Previously, <u>variational Bayes</u> was used in an autoencoder setting to have <u>variational</u> autoencoder [40, 41].
- Likewise, adversarial learning can be used in an autoencoder setting [42].
- Several adversarial-based autoencoders exist:
 - Adversarial Autoencoder (AAE) (2015) [43]
 - PixelGAN autoencoder (2017) [44]
 - Implicit Autoencoder (IAE) (2018) [45]

- Adversarial Autoencoder (AAE) was proposed in (2015) [43].
- The structure of AAE is depicted in this figure. Each of the blocks *B*₁, *B*₂, and *B*₃ in this figure has several network layers with nonlinear activation functions.



- AAE has an encoder (i.e., block B₁) and a decoder (i.e., block B₂). The input of encoder is a real data point x ∈ ℝ^d and the output of decoder is the reconstructed data point x̂ ∈ ℝ^d. One of the low-dimensional middle layers is the latent (or code) layer, denoted by z ∈ ℝ^p, where p ≪ d.
- The encoder and decoder model **conditional distributions** p(z|x) and p(x|z), respectively. Let the distribution of the latent variable in the autoencoder be denoted by q(z). This is the **posterior distribution** of latent variable.
- The blocks B_1 and B_3 are the generator G and discriminator D of adversarial network, respectively. We also have a prior distribution, denoted by p(z), on the latent variable which is chosen by the user. This prior distribution can be a *p*-dimensional normal distribution $\mathcal{N}(\mathbf{0}, \mathbf{I})$, for example.



• The **encoder** of the autoencoder (i.e., block *B*₁) is the generator *G* which generates the latent variable from the posterior distribution:

$$G(\mathbf{x}) = \mathbf{z} \sim q(\mathbf{z}). \tag{19}$$

• The discriminator *D* (i.e., block *B*₃) has a single output neuron with **sigmoid** activation function. It classifies the latent variable *z* to be a **real latent variable from the prior distribution** *p*(*z*) or **a generated latent variable by the encoder** of autoencoder:

$$D(z) := \begin{cases} \frac{1}{0} & \text{if } \frac{z \text{ is real, i.e., } z \sim p(z)}{z \text{ is generated, i.e., } z \sim q(z).} \end{cases}$$

$$(20)$$

- As was explained, the block B₁ is shared between the autoencoder and the adversarial network.
- This adversarial learning makes both autoencoder and adversarial network stronger gradually because the **autoencoder tries to generate the latent variable which is very similar to the real latent variable from the prior distribution**. In this way, it tries to fool the discriminator. The discriminator, on the other hand, tries to become stronger not to be fooled by the encoder of autoencoder.



- In AAE, we have alternating optimization [11] where reconstruction and regularization phases are repeated iteratively.
- In the **reconstruction phase**, the mean squared error is minimized between the data x and the reconstructed data \hat{x} . In the **regularization phase**, the discriminator and generator are updated using the GAN approach.
- For each of these updates, we use stochastic gradient descent [11] with **backpropagation**. Overall, the two phases are performed as:

$$\longrightarrow \underbrace{ \begin{array}{c} \swarrow \\ B_1', B_2^{(k+1)} := \arg \min_{B_1, B_2} \|\widehat{\mathbf{x}} - \mathbf{x}\|_2^2, \\ \\ & & & \\ \end{array} }_{ \left\{ \begin{array}{c} B_3^{(k+1)} := B_3^{(k)} - \eta^{(k)} \frac{\partial}{\partial B_3} \left(V(B_3, B_1') \right), \\ B_1^{(k+1)} := B_1' - \eta^{(k)} \frac{\partial}{\partial B_1} \left(V(B_3^{(k+1)}, B_1) \right), \end{array} \right.$$
 (21)

where $B_1 = G$ and $B_3 = D$. Eq. (21) is the <u>reconstruction phase</u> and Eq. (22) is the regularization phase.

• For supervised and semi-supervised AAE, refer to our tutorial paper [33].

Acknowledgment

- Some slides are based on our tutorial paper: "Generative adversarial networks and adversarial autoencoders: Tutorial and survey" [33]
- For more information on GAN and its variants, see our tutorial [33]. It covers various variants of GAN:
 - vanilla GAN, conditional GAN, DCGAN, the mode collapse problem, minibatch GAN, unrolled GAN, BourGAN, mixture GAN, D2GAN, Wasserstein GAN, f-GAN, adversarial variational Bayes, Bayesian GAN, feature matching in GAN, InfoGAN, GRAN, LSGAN, energy-based GAN, CatGAN, MMD GAN, LapGAN, progressive GAN, triple GAN, LAG, GMAN, AdaGAN, CoGAN, inverse GAN, BiGAN, ALI, SAGAN, Few-shot GAN, SinGAN, interpolation and evaluation of GAN, image-to-image translation (including PatchGAN, CycleGAN, DeepFaceDrawing, simulated GAN, interactive GAN), text-to-image translation (including StackGAN), mixing image characteristics (including FineGAN and MixNMatch), adversarial autoencoder, PixelGAN, implicit autoencoder.
- Some slides of this slide deck are inspired by teachings of <u>Prof. Ali Ghodsi</u> at University of Waterloo, Department of Statistics.

References

- I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in neural information processing* systems, vol. 27, 2014.
- [2] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial networks," *Communications of the ACM*, vol. 63, no. 11, pp. 139–144, 2020.
- [3] L. Huang, A. D. Joseph, B. Nelson, B. I. Rubinstein, and J. D. Tygar, "Adversarial machine learning," in *Proceedings of the 4th ACM workshop on Security and artificial intelligence*, pp. 43–58, 2011.
- [4] S.-M. Moosavi-Dezfooli, A. Fawzi, and P. Frossard, "DeepFool: a simple and accurate method to fool deep neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2574–2582, 2016.
- [5] A. Kurakin, I. Goodfellow, and S. Bengio, "Adversarial machine learning at scale," in International Conference on Learning Representations, 2017.
- [6] A. Kurakin, I. Goodfellow, S. Bengio, *et al.*, "Adversarial examples in the physical world," in *International Conference on Learning Representations, Workshop Track*, 2017.
- [7] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," in *International Conference on Learning Representations*, 2018.

- [8] A. Y. Ng and M. I. Jordan, "On discriminative vs. generative classifiers: A comparison of logistic regression and naive Bayes," in *Advances in neural information processing systems*, pp. 841–848, 2002.
- [9] B. Ghojogh, H. Nekoei, A. Ghojogh, F. Karray, and M. Crowley, "Sampling algorithms, from survey sampling to Monte Carlo methods: Tutorial and literature review," arXiv preprint arXiv:2011.00901, 2020.
- [10] D.-Z. Du and P. M. Pardalos, *Minimax and applications*, vol. 4. Springer Science & Business Media, 2013.
- [11] B. Ghojogh, A. Ghodsi, F. Karray, and M. Crowley, "KKT conditions, first-order and second-order optimization, and distributed optimization: Tutorial and survey," arXiv preprint arXiv:2110.01858, 2021.
- [12] L. Mescheder, S. Nowozin, and A. Geiger, "The numerics of GANs," in Advances in neural information processing systems, 2017.
- [13] I. Goodfellow, "NIPS 2016 tutorial: Generative adversarial networks," in Advances in neural information processing systems, Tutorial rack, 2016.
- [14] I. Goodfellow, D. Warde-Farley, M. Mirza, A. Courville, and Y. Bengio, "Maxout networks," in *International conference on machine learning*, pp. 1319–1327, 2013.
- [15] V. Nair and G. E. Hinton, "Rectified linear units improve restricted Boltzmann machines," in *International Conference on Machine Learning*, 2010.

- [16] F. Nielsen, "A family of statistical symmetric divergences based on Jensen's inequality," arXiv preprint arXiv:1009.4004, 2010.
- [17] V. Nagarajan and J. Z. Kolter, "Gradient descent GAN optimization is locally stable," in Advances in neural information processing systems, 2017.
- [18] F. Farnia and A. Ozdaglar, "Do GANs always have Nash equilibria?," in International Conference on Machine Learning, pp. 3029–3039, 2020.
- [19] L. Mescheder, A. Geiger, and S. Nowozin, "Which training methods for GANs do actually converge?," in *International conference on machine learning*, pp. 3481–3490, PMLR, 2018.
- [20] B. Ghojogh and M. Crowley, "The theory behind overfitting, cross validation, regularization, bagging, and boosting: tutorial," arXiv preprint arXiv:1905.12787, 2019.
- [21] M. Mirza and S. Osindero, "Conditional generative adversarial nets," *arXiv preprint arXiv:1411.1784*, 2014.
- [22] A. Radford, L. Metz, and S. Chintala, "Unsupervised representation learning with deep convolutional generative adversarial networks," in *International Conference on Learning Representations*, 2016.
- [23] J. T. Springenberg, A. Dosovitskiy, T. Brox, and M. Riedmiller, "Striving for simplicity: The all convolutional net," in *International Conference on Learning Representations*, *Workshop Track*, 2015.

- [24] A. Mordvintsev, C. Olah, and M. Tyka, "Inceptionism: Going deeper into neural networks." Google Al Blog, 2015.
- [25] S. loffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *International conference on machine learning*, pp. 448–456, 2015.
- [26] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen, "Improved techniques for training GANs," *Advances in neural information processing systems*, vol. 29, pp. 2234–2242, 2016.
- [27] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," in *Advances in neural information processing systems*, pp. 3111–3119, 2013.
- [28] L. Metz, B. Poole, D. Pfau, and J. Sohl-Dickstein, "Unrolled generative adversarial networks," in *International Conference on Learning Representations*, 2017.
- [29] C. Xiao, P. Zhong, and C. Zheng, "BourGAN: Generative networks with metric embeddings," in Advances in neural information processing systems, 2018.
- [30] Q. Hoang, T. D. Nguyen, T. Le, and D. Phung, "MGAN: Training generative adversarial nets with multiple generators," in *International Conference on Learning Representations*, 2018.

- [31] T. D. Nguyen, T. Le, H. Vu, and D. Phung, "Dual discriminator generative adversarial nets," Advances in neural information processing systems, 2017.
- [32] M. Arjovsky, S. Chintala, and L. Bottou, "Wasserstein generative adversarial networks," in International conference on machine learning, pp. 214–223, 2017.
- [33] B. Ghojogh, A. Ghodsi, F. Karray, and M. Crowley, "Generative adversarial networks and adversarial autoencoders: Tutorial and survey," arXiv preprint arXiv:2111.13282, 2021.
- [34] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *Proceedings of the IEEE international conference* on computer vision, pp. 2223–2232, 2017.
- [35] S.-Y. Chen, W. Su, L. Gao, S. Xia, and H. Fu, "DeepFaceDrawing: Deep generation of face images from sketches," ACM Transactions on Graphics (TOG), vol. 39, no. 4, pp. 72–1, 2020.
- [36] H. Zhang, T. Xu, H. Li, S. Zhang, X. Wang, X. Huang, and D. N. Metaxas, "StackGAN: Text to photo-realistic image synthesis with stacked generative adversarial networks," in *Proceedings of the IEEE international conference on computer vision*, pp. 5907–5915, 2017.
- [37] K. K. Singh, U. Ojha, and Y. J. Lee, "FineGAN: Unsupervised hierarchical disentanglement for fine-grained object generation and discovery," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6490–6499, 2019.

- [38] Y. Li, K. K. Singh, U. Ojha, and Y. J. Lee, "MixNMatch: Multifactor disentanglement and encoding for conditional image generation," in *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition, pp. 8039–8048, 2020.
- [39] U. Ojha, K. K. Singh, and Y. J. Lee, "Generating furry cars: Disentangling object shape & appearance across multiple domains," in *International Conference on Learning Representations*, 2021.
- [40] D. P. Kingma and M. Welling, "Auto-encoding variational Bayes," in International Conference on Learning Representations, 2014.
- [41] B. Ghojogh, A. Ghodsi, F. Karray, and M. Crowley, "Factor analysis, probabilistic principal component analysis, variational inference, and variational autoencoder: Tutorial and survey," arXiv preprint arXiv:2101.00734, 2021.
- [42] A. Makhzani, Unsupervised representation learning with autoencoders. PhD thesis, University of Toronto, 2018.
- [43] A. Makhzani, J. Shlens, N. Jaitly, I. Goodfellow, and B. Frey, "Adversarial autoencoders," arXiv preprint arXiv:1511.05644, 2015.
- [44] A. Makhzani and B. Frey, "PixelGAN autoencoders," in Advances in neural information processing systems, 2017.
- [45] A. Makhzani, "Implicit autoencoders," arXiv preprint arXiv:1805.09804, 2018.