

Generative Adversarial Networks for Data Modeling

Lecture Notes on Deep Learning

Avi Kak and Charles Bouman

Purdue University

Tuesday 28th March, 2023 09:33

©2023 A. C. Kak, Purdue University

Preamble

When you create a probabilistic model for your data, you acquire the power to generate new samples of the data from the model. Depending on how good a job you did of modeling the data, the new samples you generate from the model may look deceptively similar to those in your data without being exactly the same as any one of them.

In general, probabilistic modeling may involve fitting a parametric form to the data, the choice of the form based on your understanding of the phenomenon that produced the data. Obviously, you would want to choose the parameters that can account for all of the observed data in a maximum-likelihood sense.

It may also happen that you are really NOT interested in fitting a parametric model to your data, **but you are interested in generating new samples from the data nevertheless**. In such cases, it is possible you could get away with just constructing a multi-dimensional histogram from the data and using a generator of some sort that would spit out new samples according to that histogram.

Preamble

Regardless of whether you have an analytic model for the data or just a good-quality histogram, **generating new samples is not easy**. It has been the subject of much research by probability theorists and statisticians the last several decades. The best techniques fall under the label Markov-Chain Monte-Carlo (MCMC) sampling and the most commonly used algorithm for MCMC sampling is the Metropolis-Hastings algorithm.

The basic intuition in these algorithms is based on conducting a random walk through the space in which the model is defined and subjecting each successive randomly generated sample to an acceptance test that is based on the model probability distribution. As you generate a candidate for the next sample at your current point on the walk, you subject the acceptance of the candidate to the ratio of the probabilities at the candidate point and the current point. In this manner, you bias the acceptance of a candidate sample in such a way that you end up with more samples in those portions of the model space where the probabilities are relatively high. The generation of the new samples is according to what is known as a proposal distribution. Since the acceptance of each sample is predicated on just the previous sample that was already accepted, we obviously have a Markov Chain. Hence the name MCMC for such algorithms.

Preamble (contd.)

The following link is to a Perl module I created several years ago for helping generate positive and negative training samples for a machine learning algorithm using the Metropolis-Hastings algorithm for sample selection:

<https://metacpan.org/pod/Algorithm::RandomPointGenerator>

The machine learning program in this case was for classifying land-cover data obtained from wide-area satellite imagery as described in

https://engineering.purdue.edu/RVL/Publications/CVIU_2016_Chang_Comandur_Park_Kak.pdf

Fast forward to deep learning: Just as it has demolished so many of our previous approaches to solving data engineering problems, probabilistic modeling of data has suffered the same fate. The deep learning based approaches to data modeling produce stunning results that nobody could have even dared dream just a few years back. I am sure you have heard about what media refers to as “deep fakes”. That’s what I am talking about. My goal in this lecture is to introduce you to deep learning based approaches to probabilistic data modeling with neural networks.

Preamble (contd.)

The modern excitement in adversarial learning for data modeling began with the 2014 publication "Generative Adversarial Nets" by Goodfellow, Pouget-Abadie, Mirza, Xu, Warde-Farley, Ozair, Courville, and Bengio:

<https://arxiv.org/pdf/1406.2661.pdf>

Such learning involves two networks, a Discriminator network and a Generator network:

- We can think of the Discriminator network as a function $D(x, \theta_d)$ whose output is the probability that a sample x comes from the training data. The notation θ_d represents the learnable parameters in the Discriminator network.
- Similarly, we can think of the Generator network as a function $G(z, \theta_g)$ that maps noise vectors z to samples that we want to look like the samples in our training data. The vector θ_g represents the learnable parameters in the Generator network.

Preamble (contd.)

If p_{data} represents the probability distribution that describes the training data and p_g represents the probability distribution that the Generator network has learned so far, the goal of deep learning for probabilistic data modeling would be to estimate the best values for the parameters θ_d and θ_g so that some measure of the distance between distributions p_{data} and p_g is minimized.

What's interesting is that the deep learning framework that was actually implemented by Goodfellow et al. did not directly minimize a distance between p_{data} and p_g . Nevertheless, the authors argued that if the Discriminator was trained to an optimum level, it was guaranteed to yield a solution for p_g that would be a minimum Jensen-Shannon divergence approximation to p_{data} .

The above paragraph points to the following fact: **In order to understand the algorithms for probabilistic data modeling, you must first understand how to measure the “distance” between two probability distributions.**

Preamble (contd.)

For the reason stated at the bottom of the previous slide, I'll start this lecture with a brief survey of the more popular distances and divergences between two given distributions.

For any such distance to be useful in a deep learning context, you would want to treat it as a loss for the backpropagation needed for updating the parameters θ_d and θ_g that I defined previously. That places an important constraint on what kinds of distances can actually be used a deep learning algorithm: **the distance must be differentiable so that we can calculate the gradients of the loss with respect to the network parameters.**

Over the last couple of years, the Wasserstein distance has emerged as a strong candidate for such a differentiable distance function. And that has led to a Generative Adversarial Network named WGAN that was presented by Arjovsky, Chintala, and Bottou in the following 2017 publication:

<https://arxiv.org/pdf/1701.07875.pdf>

Preamble (contd.)

As I mentioned on the previous slide, I'll start this lecture with a review of the distance functions for probability distributions. That will get us ready to talk about my implementation of DCGAN and WassersteinGAN in DLStudio:

<https://engineering.purdue.edu/kak/distDLS/>

If you are already familiar with the module and for whatever reason you just need to “pip install” the latest version of the code, here is a link to its PyPi repository:

<https://pypi.org/project/DLStudio/>

The DCGAN that I mentioned above was first presented by Radford, Metz, and Chintala in the following 2016 publication:

<https://arxiv.org/pdf/1511.06434.pdf>

It was the first fully convolutional implementation of a GAN.

Outline

1	Distance Between Two Probability Distributions	10
2	Total Variation (TV) Distance	12
3	Kullback-Leibler Divergence	16
4	Jensen-Shannon Divergence and Distance	23
5	Earth Mover's Distance	27
6	Wasserstein Distance	37
7	A Random Experiment for Studying Differentiability	44
8	Differentiability of Distance Functions	47
9	PurdueShapes5GAN Dataset for Adversarial Learning	55
10	DCGAN Implementation in DLStudio	62
11	Making Small Changes to the DCGAN Architecture	84
12	Wasserstein GAN Implementation in DLStudio	90
13	Improving Wasserstein GAN with Gradient Penalty	106

Outline

1	Distance Between Two Probability Distributions	10
2	Total Variation (TV) Distance	12
3	Kullback-Leibler Divergence	16
4	Jensen-Shannon Divergence and Distance	23
5	Earth Mover's Distance	27
6	Wasserstein Distance	37
7	A Random Experiment for Studying Differentiability	44
8	Differentiability of Distance Functions	47
9	PurdueShapes5GAN Dataset for Adversarial Learning	55
10	DCGAN Implementation in DLStudio	62
11	Making Small Changes to the DCGAN Architecture	84
12	Wasserstein GAN Implementation in DLStudio	90
13	Improving Wasserstein GAN with Gradient Penalty	106

Estimating the Distance Between Two Distributions

- Given two probability distributions, p_{data} and p_g , the former representing the training data and the latter an approximation to the former as learned by some ML framework, the question is: As a measure of the dissimilarity of the two distributions, what is the distance between the two?
- Along the lines of a review of such distances that was presented in

<https://arxiv.org/pdf/1701.07875.pdf>

let's briefly review the following popular distances and divergences between a pair of probability distributions:

- Total Variation Distance**
- Kullback-Liebler Divergence**
- Jensen-Shannon Divergence**
- Earth Mover's Distance**
- Wasserstein Distance**

Outline

1	Distance Between Two Probability Distributions	10
2	Total Variation (TV) Distance	12
3	Kullback-Leibler Divergence	16
4	Jensen-Shannon Divergence and Distance	23
5	Earth Mover's Distance	27
6	Wasserstein Distance	37
7	A Random Experiment for Studying Differentiability	44
8	Differentiability of Distance Functions	47
9	PurdueShapes5GAN Dataset for Adversarial Learning	55
10	DCGAN Implementation in DLStudio	62
11	Making Small Changes to the DCGAN Architecture	84
12	Wasserstein GAN Implementation in DLStudio	90
13	Improving Wasserstein GAN with Gradient Penalty	106
	Purdue University	12

Total Variation (TV) Distance

- We start with a continuous random variable $\{X \mid x \in R^n\}$ and consider two different probability distributions (densities, really), denoted f and g , over X . The Total Variation (TV) distance between f and g is given by

$$d_{TV}(f, g) = \sup_A \left[\left| \int_A f(x)dx - \int_A g(x)dx \right| \quad : \quad A \subset R^n \right] \quad (1)$$

- What that says is that we check every subset A of the domain R^n and find the difference between the probability masses over that subset for the f and g densities. The largest value for this difference is the TV distance between the two.
- The important thing here is that the TV distance is a *metric*, in the sense that it satisfies all the conditions for a distance measure to be a *metric*: Must never be negative; must be symmetric; and must obey the triangle inequality.

TV for the Discrete Case

- Let's now consider the case when the random variable X is discretized. That is, the observed values for X are confined to the set shown below:

$$X = \{x_1, x_2, \dots, x_N\}$$

- We are now interested in the distance between two discrete probability distributions, to be denoted P and Q , over a countable set. These distributions must obviously satisfy the unit summation condition:

$$\sum_{i=1}^N P(x_i) = 1 \quad \sum_{i=1}^N Q(x_i) = 1 \quad (2)$$

- In this case, the Total Variation distance is given by:

$$d_{TV}(P, Q) = \sup_A \left[\left| \sum_{x_i \in A} P(x_i) - \sum_{x_i \in A} Q(x_i) \right| : A \subset X \right] \quad (3)$$

TV for the Discrete Case (contd.)

- Let's now consider the following two subsets of the set X :

$$A_1 = \{x_i \in X \mid P(x_i) \geq Q(x_i)\}$$

$$A_2 = \{x_i \in X \mid Q(x_i) < P(x_i)\} \quad (4)$$

- On account of the absolute value operator in Eq. (3), for the optimizing set A , it must either be the case that $P(x_i) \geq Q(x_i)$ or that $Q(x_i) \geq P(x_i)$. What that implies that both A_1 and A_2 are a part of the optimizing set A . However, since $A_1 \cup A_2 = X$, we can write for the discretized case:

$$\begin{aligned} d_{TV}(P, Q) &= \frac{1}{2} \sum_{x_i \in X} |P(x_i) - Q(x_i)| \\ &= \frac{1}{2} L_1(P, Q) \end{aligned} \quad (5)$$

where the L_1 norm is the Minkowski norm L_p with $p = 1$.

Outline

1	Distance Between Two Probability Distributions	10
2	Total Variation (TV) Distance	12
3	Kullback-Leibler Divergence	16
4	Jensen-Shannon Divergence and Distance	23
5	Earth Mover's Distance	27
6	Wasserstein Distance	37
7	A Random Experiment for Studying Differentiability	44
8	Differentiability of Distance Functions	47
9	PurdueShapes5GAN Dataset for Adversarial Learning	55
10	DCGAN Implementation in DLStudio	62
11	Making Small Changes to the DCGAN Architecture	84
12	Wasserstein GAN Implementation in DLStudio	90
13	Improving Wasserstein GAN with Gradient Penalty	106

Kullback-Liebler Divergence

- Popularly known as **KL-Divergence**.
- In this case, let's start directly with the discrete case of a random variable X as stated in the first two bullets on Slide 14. The KL-Divergence between a true distribution P and its approximating distribution Q is given by

$$d_{KL}(P, Q) = \sum_{i=1}^N P(x_i) \log \frac{P(x_i)}{Q(x_i)} \quad (6)$$

- $d_{KL}(P, Q)$ is obviously the expectation of the ratios $\log \frac{P(x_i)}{Q(x_i)}$ with respect to the P distribution. For the ratios to be defined you must have $Q(x_i) > 0$ when $P(x_i) > 0$. $Q(x_i)$ is allowed to be zero when $P(x_i)$ is zero since $x \log x \rightarrow 0$ as $x \rightarrow 0+$.
- The logarithm shown above is taken to base 2 if the value of the divergence is required in bits. For natural logarithms, the value returned by KL Divergence is in nats.

KL-Divergence (contd.)

- Since, in general, $\log x$ can return negative and positive values as x increases from 0 to $+\infty$, and since a negative value for KL-divergence makes no sense, how can we be sure that the value of $d_{KL}(P, Q)$ is always non-negative?
- To see that the formula for $d_{KL}(P, Q)$ always returns a non-negative value, we first subject that formula to the following rewrites:

$$\begin{aligned}
 d_{KL}(P, Q) &= \sum_{i=1}^N P(x_i) \log \frac{P(x_i)}{Q(x_i)} \\
 &= - \sum_{i=1}^N P(x_i) \log \frac{Q(x_i)}{P(x_i)} \\
 &= - \sum_{i=1}^N P(x_i) \log \frac{P(x_i) + Q(x_i) - P(x_i)}{P(x_i)} \\
 &= - \sum_{i=1}^N P(x_i) \log \left[1 + \frac{Q(x_i) - P(x_i)}{P(x_i)} \right] \\
 &= - \sum_{i=1}^N P(x_i) \log(1 + a) \tag{7}
 \end{aligned}$$

KL-Divergence (contd.)

- In the last equation on the previous slide, $a = \frac{Q(x_i) - P(x_i)}{P(x_i)}$. The factor a is lower bounded by -1 , which happens when $P(x_i)$ takes on the largest possible value of 1 and $Q(x_i)$ takes on the smallest possible value of 0.
- Using Jensen's inequality to take advantage of the concavity of $\log x$ over the interval $(0, \infty)$, one can show that for all $a > -1$, $\log(1 + a) \leq a$. The derivation on the previous slide can therefore be extended as follows:

$$\begin{aligned}
 d_{KL}(P, Q) &\geq - \sum_{i=1}^N P(x_i) \frac{Q(x_i) - P(x_i)}{P(x_i)} \\
 &= - \sum_{i=1}^N [Q(x_i) - P(x_i)] \\
 &= 0
 \end{aligned} \tag{8}$$

which implies that we are guaranteed that $d_{KL}(P, Q) \geq 0$.

KL-Divergence (contd.)

- KL-Divergence **CANNOT** be a *metric distance*, not the least because what it calculates is asymmetric with respect to its two args.
- Given its limitations — requiring $Q(x) > 0$ when $P(x) > 0$ and not being a metric distance — **students frequently want to know as to why KL-Divergence is as “famous” as it is in the estimation-theoretic literature.** One reason for that is its interpretation as relative entropy:

$$d_{KL}(P, Q) = H_{P(Q)} - H(P) \quad (9)$$

which follows straightforwardly from the definition in Eq. (6). $H(P)$ is the entropy associated with the probability distribution P and $H_{P(Q)}$ the cross-entropy of an approximating distribution Q vis-a-vis the true distribution P . [See the definitions for $H(P)$ and $H_{P(Q)}$ on the next slide.]

- Since $d_{KL}(P, Q) \geq 0$, it must be the case that $H_{P(Q)} \geq H(P)$, **which constitutes a proof of the assertion made on Slide 17 of my Week 7 lecture that the smallest possible value for $H_{P(Q)}$ is $H(P)$.**

KL-Divergence (contd.)

- Whereas the entropy associated with a distribution P is defined as $H(P) = -\sum_{i=1}^N P(x_i) \log P(x_i)$, the cross-entropy of an approximate distribution Q with respect to a true distribution P is given by $H_P(Q) = -\sum_{i=1}^N P(x_i) \log Q(x_i)$. [Entropy based interpretations of uncertainty are valuable for developing powerful algorithms for data engineering. See Sections 2 through 4 of my Decision Trees tutorial at the clickable link <https://engineering.purdue.edu/kak/Tutorials/DecisionTreeClassifiers.pdf>.]
- The entropy based definition of KL-Divergence in Eq. (9) on the previous slide implies that the divergence is a measure of the uncertainty in the estimated distribution Q over and above what it is in the original distribution P . [See Slide 15 of my Week 7 lecture for why the entropy is a measure of uncertainty.]
- Perhaps the most important role that KL-Divergence plays in our discussion on Adversarial Networks is that it is a stepping stone to learning the Jensen-Shannon divergence (and the closely related Jensen-Shannon distance) that I present starting with the next section.

KL-Divergence (contd.)

- In Python, a call like:

```
import scipy.stats
scipy.stats.entropy(P,Q)
```

with P and Q standing for two normalized (or unnormalized) histograms, returns the KL-Divergence of Q vis-a-vis P . If $Q(x)$ is zero where $P(x)$ is not, it will throw an exception.

- In the calls shown above, the two histogram arrays must be of equal length. You can also specify the base of the logarithm with an optional third argument. The default for the base is e for the natural algorithm.

Outline

1	Distance Between Two Probability Distributions	10
2	Total Variation (TV) Distance	12
3	Kullback-Leibler Divergence	16
4	Jensen-Shannon Divergence and Distance	23
5	Earth Mover's Distance	27
6	Wasserstein Distance	37
7	A Random Experiment for Studying Differentiability	44
8	Differentiability of Distance Functions	47
9	PurdueShapes5GAN Dataset for Adversarial Learning	55
10	DCGAN Implementation in DLStudio	62
11	Making Small Changes to the DCGAN Architecture	84
12	Wasserstein GAN Implementation in DLStudio	90
13	Improving Wasserstein GAN with Gradient Penalty	106

Jensen-Shannon Divergence and Distance

- We again have a random variable X whose observed samples belong to the set:

$$X = \{x_1, x_2, \dots, x_N\} \quad (10)$$

- And, as for the case of KL-Divergence, we consider a true probability distribution P and its approximation Q over the values taken on by the random variable. The Jensen-Shannon divergence, defined below, is a symmetrized version of the KL-Divergence presented earlier in Eq. (6):

$$d_{JS}(P, Q) = d_{KL}(P, M) + d_{KL}(Q, M) \quad (11)$$

where M is the mean distribution for P and Q , as given by

$$M = \frac{P + Q}{2} \quad (12)$$

- We can also talk about Jensen-Shannon *distance*, which is given by the square-root of the Jensen-Shannon Divergence:

$$\text{dist}_{JS}(P, Q) = \sqrt{d_{JS}(P, Q)} \quad (13)$$

JS Divergence and Distance (contd.)

- Both the divergence $d_{JS}(P, Q)$ and the distance $dist_{JS}(P, Q)$ are symmetric with respect to the arguments P and Q . Additionally, they do away with the “ $Q(x) > 0$ when $P(x) > 0$ ” requirement of KL-Divergence.
- Since, as established earlier in these slides, the KL Divergence is always non-negative, the JS-Divergence is also non-negative.
- The value of $d_{JS}(P, Q)$ is always a real number in the closed interval $[0, 1]$. When the value is 0, the two distributions P and Q are identical. And when the value is 1, the two distributions are as different as they can possibly be.
- **Most significantly, $dist_{JS}(P, Q)$ is a valid metric distance.**

JS Divergence and Distance (contd.)

- Given two histogram arrays P and Q of equal length, normalized or unnormalized, a call like the following in Python

```
from scipy.spatial import distance
distance.jensenshannon(P,Q)
```

directly returns the Jensen-Shannon *distance* between the two histograms. If you wanted the Jensen-Shannon *divergence*, you would need to square the answer returned. The function call implicitly normalizes the histogram arrays if you supply them otherwise.

- With regard to the role of the Jensen-Shannon divergence (and, therefore, also of the KL-Divergence) in the context of this lecture, **the authors Goodfellow et al. of “Generative Adversarial Nets” have argued that if the Discriminator in a GAN is trained to its optimum, the distribution learned by the Generator is guaranteed to be the one whose Jensen-Shannon divergence from the training-data distribution is minimized.**

Outline

1	Distance Between Two Probability Distributions	10
2	Total Variation (TV) Distance	12
3	Kullback-Leibler Divergence	16
4	Jensen-Shannon Divergence and Distance	23
5	Earth Mover's Distance	27
6	Wasserstein Distance	37
7	A Random Experiment for Studying Differentiability	44
8	Differentiability of Distance Functions	47
9	PurdueShapes5GAN Dataset for Adversarial Learning	55
10	DCGAN Implementation in DLStudio	62
11	Making Small Changes to the DCGAN Architecture	84
12	Wasserstein GAN Implementation in DLStudio	90
13	Improving Wasserstein GAN with Gradient Penalty	106

Earth Mover's Distance

- **The distance function that the DL community is all excited about at the moment is the Wasserstein Distance.** The reason has to do with the fact this is the only differentiable distance function and, because it is differentiable, a loss based on this distance function can be backpropagated directly for updating the weights in a network.
- However, in order to fully appreciate what exactly is measured by the Wasserstein Distance, you first have to understand what is known as the Earth Mover's Distance (EMD). **Note that many researchers use the two names interchangeably. I personally think of the Wasserstein Distance as a stochastic version of EMD.**
- My goal in this section is to introduce you to EMD. My intro to EMD is based on the following classic paper by Rubner, Tomasi, and Guibas:

<http://robotics.stanford.edu/~rubner/papers/rubnerIjcv00.pdf>

Earth Mover's Distance (contd.)

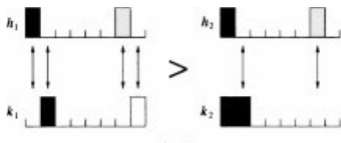
- To appreciate EMD, consider establishing similarity between two images on the basis of the histograms of their graylevels.
- Given two N -bin histograms f and g for the two images, you would not be too far off the mark if the first idea that pops up in your head would be to carry out a bin-by-bin comparison using a distance like:

$$d_{L_r}(f, g) = \left(\sum_{i=1}^N |g_i - h_i|^r \right)^{\frac{1}{r}} \quad (14)$$

- With $r = 1$, you'd be computing the L_1 distance between the two histograms, and with $r = 2$ the Euclidean distance. You will see both being used rather commonly, but you have to be careful as you will soon see. As mentioned on Slide 15, the general form of the distance shown above is known as the Minkowski distance.

Earth Mover's Distance (contd.)

- That a distance function of the sort shown on the previous slide might give nonsensical answers for image similarity is made beautifully clear by the following example from the Rubner et al. paper:



Comparing histograms

- In the figure shown above, first focus on the (h_1, k_1) histograms shown in the left column. The h_1 image has half its pixels very dark and the other half of the pixels very white. Perceptually, the k_1 image is going to look very similar to the h_1 image since the two dominant gray levels are merely shifted to the right by one unit. If the number of bins is, say, greater than 64, you will not even notice the shift.

Earth Mover's Distance (contd.)

- Next, focus on the (h_2, k_2) histograms in the figure on the previous slide. While the h_2 image has half its pixels very dark and the other half very white, the k_2 image contains only dark pixels.
- Therefore, to a human observer, the two images in the (h_1, k_1) pair will look very similar, while the two images in the (h_2, k_2) pair will look very different. **However, the d_{L_r} distance in Eq. (14) will give you exactly the opposite answer.**
- **Since distances like d_{L_r} in Eq. (14) cannot be trusted to yield meaningful results when comparing histograms for image similarity, EMD has emerged as a powerful alternative.**
- EMD is based on associating a cost with moving pixels from one bin to another in a hypothetical attempt that tries to make the two histograms as similar looking as possible, constructing an overall cost with all such pixel transfers, and then minimizing the overall cost. **31**

Earth Mover's Distance (contd.)

- Consider the following as an example of the cost associated with moving a pixel from one bin to another in a one-dimensional grayscale histogram whose bins are one-unit wide:

$$c_{ij} = 1 - e^{-\alpha|i-j|} \quad (15)$$

where you can think of $\alpha > 0$ as a heuristic parameter that is approximately proportional to the overall variability in the bin populations. It was shown by Rubner et al. that such a cost function is a metric. What it says is that cost of moving pixels from a bin to another close-by bins is close to zero. However, the costs go up if the transfer is between more widely separated bins.

- The problem of comparing two histograms can now be stated as an instance of the classic “transportation simplex” problem in optimal transport theory for resource distribution, as explained on the next slide.**

Earth Mover's Distance (contd.)

- You have M providers of some resource who possess different quantities ($\{g_i | i = 1, \dots, M\}$) of the resource and N consumers of the same resource whose needs vary according to ($\{h_j | j = 1, \dots, N\}$).
- And you also have a cost estimate c_{ij} that is the cost of transporting a unit of the resource from the i^{th} provider to the j^{th} consumer.
- Our goal is to come up with with an optimum flow matrix F , whose f_{ij} element tells us how much of the resource to transport from the i^{th} provider to the j^{th} consumer. We must obviously solve the following minimization problem for F :

$$\min_F \sum_{i=1}^M \sum_{j=1}^N c_{ij} f_{ij} \quad (16)$$

with the minimization subject to the constraints shown on the next slide.

Earth Mover's Distance (contd.)

- The minimization problem on the previous slide must be solved subject to the constraints:

$$f_{ij} \geq 0 \quad i = 1, \dots, M, \quad j = 1, \dots, N \quad (17)$$

$$\sum_{j=1}^N f_{ij} \leq h_i \quad i = 1, \dots, M \quad (18)$$

$$\sum_{i=1}^M f_{ij} \leq g_j \quad j = 1, \dots, N \quad (19)$$

$$\sum_{i=1}^M \sum_{j=1}^N f_{ij} = \min \left\{ \sum_{i=1}^M g_i, \sum_{j=1}^N h_j \right\} \quad (20)$$

- All four constraints are straightforward because they are so intuitive. [The constraints in Eqs. (17) and (18) are straightforward: The flow can never be negative and the total outgoing flow from a provider cannot exceed what the provider has in stock. The constraint in Eq. (19) also makes sense since the accumulated in-flows for the j^{th} consumer should not exceed to total demand for that consumer. The constraint in Eq. (20) is important only when the total supply provided by all the providers is not equal to the total demand at all the consumers. Should there be such a disparity between total supply and total demand, summing all of elements of the flow matrix should not exceed the smaller of the total supply and the total demand.]

Earth Mover's Distance (contd.)

- Having calculated the optimal transport by solving the minimization problem described on the previous two slides, we use the following formula to compute the EMD between the suppliers distribution for the resource and the consumers distribution:

$$EMD(g, h) = \frac{\sum_{i=1}^M \sum_{j=1}^N c_{ij} f_{ij}}{\sum_{i=1}^M \sum_{j=1}^N f_{ij}} \quad (21)$$

where we normalize the cost of the optimal transport of the goods by the total amount of the goods transported.

- Such optimization problems have received much attention by the OR (Operations Research) folks over the last several decades. We now have polynomial-time solutions for the problem that fall under the general category of “simplex algorithms for linear programming”. Rubner et al. used such a solution in their work on retrieval from image databases and showed impressive results.

Earth Mover's Distance (contd.)

- It was shown by Rubner et al. that EMD is a metric when the supplier and the consumer distributions are normalized. For the case of comparing image histograms, we can say that EMD between two histograms is a metric for the case of normalized histograms.
- With that as an intro to EMD, the issue that should come up next would be whether it is possible to create a loss function directly from EMD for adversarial learning. I'll address this question later when I get into the differentiability of the different distance functions.
- For now, let's move on to the Wasserstein distance. **As mentioned earlier, I consider the Wasserstein distance to be a stochastic version of EMD.**

Outline

1	Distance Between Two Probability Distributions	10
2	Total Variation (TV) Distance	12
3	Kullback-Leibler Divergence	16
4	Jensen-Shannon Divergence and Distance	23
5	Earth Mover's Distance	27
6	Wasserstein Distance	37
7	A Random Experiment for Studying Differentiability	44
8	Differentiability of Distance Functions	47
9	PurdueShapes5GAN Dataset for Adversarial Learning	55
10	DCGAN Implementation in DLStudio	62
11	Making Small Changes to the DCGAN Architecture	84
12	Wasserstein GAN Implementation in DLStudio	90
13	Improving Wasserstein GAN with Gradient Penalty	106

Wasserstein Distance

- Using $d_W(P, Q)$ to denote the Wasserstein distance between the distributions P and Q , here is its definition:

$$d_W(P, Q) = \inf_{\gamma(X, Y) \in \Gamma(P, Q)} E_{(X, Y) \sim \gamma} [\|x - y\|] \quad (22)$$

- In the above definition, $\Gamma(P, Q)$ is the set of all possible joint distributions $\gamma(X, Y)$ over two random variables X and Y such that the marginal of $\gamma(X, Y)$ with respect to X is P and the marginal of $\gamma(X, Y)$ with respect to Y is Q .
- Since the marginal of $\gamma(X, Y)$ with respect to X is $P(x)$ and the marginal of the same with respect to Y is $Q(x)$, $\gamma(X, Y)$ encodes in it the probability mass that must be shifted from the distribution P to the distribution Q if for whatever reason we wanted them to become identical. [If $\gamma(X, Y)$ encodes in it the probability mass that must be shifted from the distribution P to the distribution Q , is there any way to construct a "cost" — a single number — associated with this transfer of mass? The cost itself is proportional to the absolute difference between the value x for the random variable X and the value y for the random variable Y if the joint distribution $\gamma(X, Y)$ indicates there is a non-zero probability associated with mass transfer from x to y . For vector random variables, this would be the same as the norm $\|x - y\|$. In order to get a single-number cost, we would need to average the norm $\|x - y\|$ as indicated in Eq. (22) above.]

Wasserstein Distance (contd.)

- The $d_W(P, Q)$ distance is a metric as it obeys the constraints on metrics: its values are guaranteed to be non-negative, it is symmetric with respect to its args, and it obeys the triangle inequality. Let's now focus on what it might take to compute the Wasserstein distance.
- The infimum required on the right side of Eq. (22) says that from the set $\Gamma(P, Q)$ of all joint distributions defined in the second bullet on the previous slide, we need to zero in on the joint distribution $\gamma(X, Y)$ that minimizes the mean value of the normed difference $\|x - y\|$ with the sample pair (x, y) drawn from the joint distribution.
- In a computation based on a literal interpretation of the definition in Eq. (22), we are required to carry out a random experiment in which we sample the (infinite) set $\Gamma(P, Q)$ of the joint distributions for the two random variables X and Y for a candidate distribution $\gamma(X, Y)$.

Wasserstein Distance (contd.)

- Subsequently, in another random experiment, we sample the distribution $\gamma(X, Y)$ for specific values x and y for the random variables X and Y . We carry out the second random experiment repeatedly in order to form a good estimate for the average value for $\|x - y\|$. Subsequently, we go back to the first random experiment and choose a second candidate for $\gamma(X, Y)$, and so on. **Such a computation is obviously not feasible.**
- Fortunately, the infimum in the theoretical definition of Wasserstein Distance in Eq. (22) can be converted into a computationally tractable supremum calculated separately over the component distributions P and Q as shown below

$$d_W(P, Q) = \sup_{\|f\|_L \leq 1} \left[E_{x \sim P}\{f(x)\} - E_{y \sim Q}\{f(y)\} \right] \quad (23)$$

for ALL 1-Lipschitz functions $f : X \rightarrow R$ where X is the domain from which the elements x and y mentioned above are drawn and R is the set of all reals.

Wasserstein Distance (contd.)

- The result shown in Eq. (23) is from a famous book in Optimal Transport Theory by Cédric Villani:

https://cedricvillani.org/sites/dev/files/old_images/2012/08/preprint-1.pdf

- Despite the use of "ALL" for the family of 1-Lipschitz functions $f()$ in Eq. (23), **a better way to state the same thing would be that there exists a 1-Lipschitz function $f()$ for which the maximization shown on the right in Eq. (23) yields the value for the Wasserstein distance.**
- **But what is a k-Lipschitz Function?** A function $f : X \rightarrow R$ is a k-Lipschitz function if $|f(x_1) - f(x_2)| \leq k \cdot d(x_1, x_2)$ for every $x_1, x_2 \in X$. Note that X is the domain of the function. In this definition, $d(.,.)$ is the metric distance defined on the domain of f . So $d(x_1, x_2)$ is the distance between the points x_1 and x_2 .

Wasserstein Distance (contd.)

- In general, the Lipschitz functions allow us to prescribe functions with “levels” of continuity properties. The larger the value of the integer k , the more rapidly the function would be allowed to change when you go from a point x_1 to another point x_2 in its domain.
- In general, at all x in the domain X of f :

$$f(x) = \inf_{y \in X} [f(y) + k \cdot d(x, y)] = \sup_{y \in X} [f(y) - k \cdot d(x, y)] \quad (24)$$

- Note that the definition $|f(x) - f(y)| \leq k \cdot d(x, y)$ implies $f(y) - k \cdot d(x, y) \leq f(x) \leq f(y) + k \cdot d(x, y)$. When you apply the definitions of infimum and supremum to these inequalities, you get the form shown in Eq. (24).

Wasserstein Distance (contd.)

- We are faced with the following questions if we want to use the form in Eq. (23) for computing the Wasserstein Loss in adversarial learning:
 - How do we find the function $f()$ that would solve the maximization problem in Eq. (23)?
 - The expectation operator $E()$ in Eq. (23) is meant to be applied over the entire domain of the distributions P and Q . How do we do that in a practical setting?
- I'll address each of these issues separately in Section 12 on how to use the Wasserstein distance for adversarial learning. That material begins on Slide 90.

Outline

1	Distance Between Two Probability Distributions	10
2	Total Variation (TV) Distance	12
3	Kullback-Leibler Divergence	16
4	Jensen-Shannon Divergence and Distance	23
5	Earth Mover's Distance	27
6	Wasserstein Distance	37
7	A Random Experiment for Studying Differentiability	44
8	Differentiability of Distance Functions	47
9	PurdueShapes5GAN Dataset for Adversarial Learning	55
10	DCGAN Implementation in DLStudio	62
11	Making Small Changes to the DCGAN Architecture	84
12	Wasserstein GAN Implementation in DLStudio	90
13	Improving Wasserstein GAN with Gradient Penalty	106

A Random Experiment for Studying Differentiability

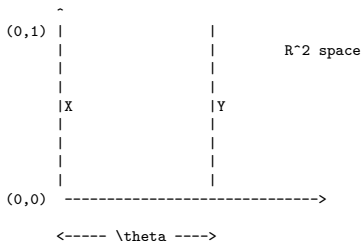
- The discussion in this section is an elaboration of the “learning parallel lines” example in the paper

<https://arxiv.org/pdf/1701.07875.pdf>

- We start with a random variable Z whose values, z , are uniformly distributed over the unit interval $[0, 1]$.
- We assume that the ground-truth consists of Z -values on the y -axis in \mathbb{R}^2 — this would presumably be our “training” data (to make an analogy with GAN training). Now imagine a GAN Generator that is also capable of producing the same kind of points in \mathbb{R}^2 but the points produced by the Generator are offset horizontally by a learnable parameter θ . The true value of θ is obviously 0, but the Generator has to learn that during training.
- We use X as the random variable to denote the points on the ground-truth line in \mathbb{R}^2 and Y to denote the points being produced by the Generator on another vertical line that is horizontally offset by 0.5

Studying Differentiability (contd.)

- Let P denote the distribution for the ground-truth points X and Q the distribution for the GAN-generated points Y .
- Note again that the ground-truth points X are the set of all points $\{(0, z) \in \mathbb{R}^2 | z \sim U[0, 1]\}$ and the GAN-generated points Y form the set $\{(\theta, z) \in \mathbb{R}^2 | z \sim U[0, 1]\}$.
- The following figure illustrates the relationship between X , Y , and the sole learnable parameter θ .



Outline

1	Distance Between Two Probability Distributions	10
2	Total Variation (TV) Distance	12
3	Kullback-Leibler Divergence	16
4	Jensen-Shannon Divergence and Distance	23
5	Earth Mover's Distance	27
6	Wasserstein Distance	37
7	A Random Experiment for Studying Differentiability	44
8	Differentiability of Distance Functions	47
9	PurdueShapes5GAN Dataset for Adversarial Learning	55
10	DCGAN Implementation in DLStudio	62
11	Making Small Changes to the DCGAN Architecture	84
12	Wasserstein GAN Implementation in DLStudio	90
13	Improving Wasserstein GAN with Gradient Penalty	106

Differentiability of Distance Functions

- Given the sets X and Y as defined on Slides 45 and 46, **we start with examining the differentiability of the Wasserstein Distance.**
- Given the definition that X is set of all points $\{x = (0, z) \in \mathbb{R}^2 | z \sim U[0, 1]\}$ and Y is the set of all points $\{y = (\theta, z) \in \mathbb{R}^2 | z \sim U[0, 1]\}$, we can say that the difference $\|x - y\|$ needed for calculating the Wasserstein distance using Eq. (22) on Slide 38 will always be equal to the value of the parameter θ .
- The same would be the case if we used the supremum based estimate of the Wasserstein distance using Eq. (23). Therefore, for the random experiment under consideration, we can claim:

$$d_W(P, Q) = \theta \quad (25)$$

- So we see that the Wasserstein distance is continuous and differentiable with respect to the learnable parameter θ . That makes it a good candidate as a loss function in a neural network.

Differentiability of Distance Functions (contd.)

- What is interesting is that the closely related EMD distance does not possess the property of differentiability with respect to the learnable parameters. That is because it involves comparing histograms directly. Since a histogram is a discretization of continuous values, it is not possible to backpropagate any partial derivatives through such a step.
- **Let's now consider the differentiability of KL-Divergence.**
- The definition of KL-Divergence provided earlier in Eq. (6) is for the case of random variables that take discrete values. But the “parallel lines” example involves two continuous random variables X and Y . Here is the definition of KL-Divergence for the continuous case:

$$d_{KL}(P, Q) = \int P(x) \log \frac{P(x)}{Q(x)} dx \quad (26)$$

- The scope of the variable x of integration is the space of all random outcomes over which both the distributions P and Q are defined.

Differentiability of Distance Functions (contd.)

- The last bullet on the previous implies that x must span both the lines X and Y for this integration. However, the sets X and Y are disjoint except when the Generator parameter θ equals zero.
- When X and Y are disjoint, we run headlong into the condition $Q(x) = 0$ when $P(x) > 0$ that makes the divergence d_{KL} become infinity. Hence we can write:

$$\begin{aligned} d_{KL}(P, Q) &= 0 && \theta = 0 \\ &= +\infty && \theta \neq 0 \end{aligned} \tag{27}$$

- Obviously, KL-Divergence is not differentiable with respect to the learnable parameter θ .
- **Next we take up the case of differentiability of JS-Divergence.**

Differentiability of Distance Functions (contd.)

- The formula for JS-Divergence was presented in Eq. (11) on Slide 24. Given two distributions P and Q , the formula in that equation requires that we first calculate the mean distribution M as defined in Eq. (12).
- For what follows, recall the fact that JS-Divergence is a symmetrization of KL-Divergence that is meant to get around the main shortcoming of the latter in those regions of the probability space where $Q(x) = 0$ whereas $P(x) > 0$.
- Note that M in Eq. (12) is a **mixture distribution**. By definition, given two separate distributions P and Q defined over the same set of random outcomes, a mixture means merely that the next sample will be drawn randomly either from P or from Q . Since the two component distributions P and Q in the mixture M are weighted equally (by a factor $\frac{1}{2}$), the individual distributions will be selected with equal probability for the realizations of M .
- On the next slide, we will consider the first term in the summation in Eq. (11). The result for the second term would be the same.

Differentiability of Distance Functions (contd.)

- Focusing on the case when the learnable parameter θ is nonzero, that is, when we are going to encounter the condition $Q(x) = 0$ when $P(x) > 0$ (which will happen on line x as explained previously for the case of differentiability of KL-Divergence), let's focus on the first term on the RHS in Eq. (11) on Slide 24:

$$\begin{aligned}
 d_{KL}(P, M) &= \int P(x) \log \frac{P(x)}{M(x)} dx \\
 &= \int P(x) \left[\log P(x) - \log \frac{P(x) + Q(x)}{2} \right] dx \\
 &= \int P(x) \left[\log P(x) - \log(P(x) + Q(x)) + \log 2 \right] dx \\
 &= \int P(x) \log 2 dx \\
 &= \log 2
 \end{aligned} \tag{28}$$

- As expected, the expressions on the RHS of Eq. (11) are now inoculated against going to infinity under the condition $Q(x) = 0$ when $P(x) > 0$.

Differentiability of Distance Functions (contd.)

- Since both the component expressions on the RHS of Eq. (11) lead to exactly the same result that is shown above, we can say that $d_{JS}(P, Q) = \log 2$ for the case $\theta \neq 0$.

- Therefore, we can write:

$$\begin{aligned} d_{JS}(P, Q) &= 0 & \theta = 0 \\ &= \log 2 & \theta \neq 0 \end{aligned} \tag{29}$$

which is again not differentiable with respect to the parameter θ .

- We next take up the differentiability of the Total Variation Distance**
- The Total Variation (TV) distance for the continuous case was defined in Eq. (1).
- That definition calls for identifying a subset A of the probability space defined by all possible outcomes that maximizes the difference between P 's probability mass over A and Q 's probability mass over A .

Differentiability of Distance Functions (contd.)

- When $\theta \neq 0$, we could choose for such an A the set X itself. Since the probability mass of P over this set equals 1 whereas the probability mass of Q over the same set equals 0. The difference of the two integrals in Eq. (1) on Slide 13 for such an A is the largest it can be — equal to 1.
- On the other hand, when the Generator's parameter θ equals 0, the sets X and Y become congruent. In this case, the difference of the two integrals in Eq. (1) would be zero.
- So we can write:

$$\begin{aligned}
 d_{TV}(P, Q) &= 0 & \theta = 0 \\
 &= 1 & \theta \neq 0
 \end{aligned}
 \tag{30}$$

- TV is obviously not a differentiable distance function.

Outline

1	Distance Between Two Probability Distributions	10
2	Total Variation (TV) Distance	12
3	Kullback-Leibler Divergence	16
4	Jensen-Shannon Divergence and Distance	23
5	Earth Mover's Distance	27
6	Wasserstein Distance	37
7	A Random Experiment for Studying Differentiability	44
8	Differentiability of Distance Functions	47
9	PurdueShapes5GAN Dataset for Adversarial Learning	55
10	DCGAN Implementation in DLStudio	62
11	Making Small Changes to the DCGAN Architecture	84
12	Wasserstein GAN Implementation in DLStudio	90
13	Improving Wasserstein GAN with Gradient Penalty	106

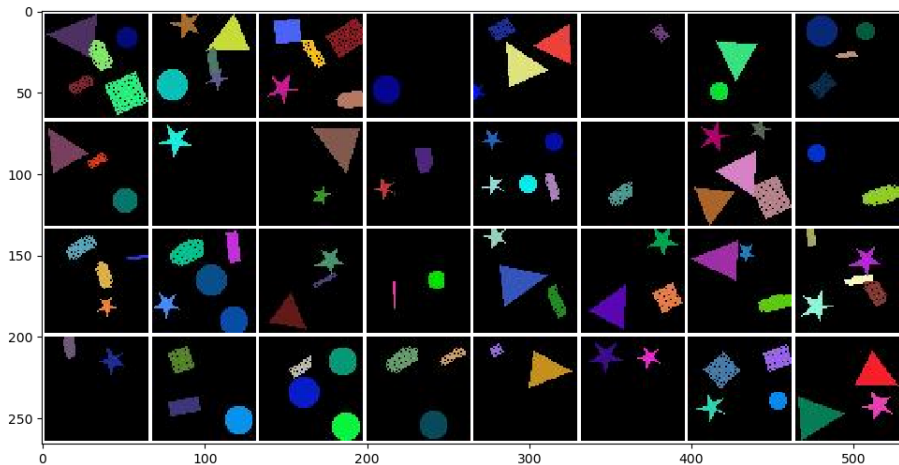
PurdueShapes5GAN Dataset of Images

- I have created a dataset, PurdueShapes5GAN, for experimenting with the three GANs in version 2.0.3 (or higher) of the DLStudio module. **Each image in the dataset is of size 64×64 .** The dataset consists of 20,000 images.
- **This dataset of rather small-sized images was created to make it easier to give classroom demonstrations of the training code and also for the students to be able to run the code on their laptops (at least those that come equipped with a GPU for graphics rendering, as many of them do these days).**
- The program that generates the PurdueShapes5GAN dataset is a modification of the script I used for the PurdueShapes5MultiObject dataset that I used previously in the lecture on semantic segmentation.

PurdueShapes5GAN Dataset (contd.)

- Compared to its predecessor semantic-segmentation dataset, the annotations that were needed for the semantic segmentation dataset (the bounding boxes and masks) are no longer necessary for adversarial learning of a probabilistic data model for a set of images. That makes a GAN dataset much simpler compared to a semantic-segmentation dataset.
- Each image in the PurdueShapes5GAN dataset contains a random number of up to five shapes: rectangle, triangle, disk, oval, and star. Each shape is located randomly in the image, oriented randomly, and assigned a random color. Since the orientation transformation is carried out without bilinear interpolation, it is possible for a shape to acquire holes in it. Shown in the next slide is a batchful of images that is processed in each iteration of the training loop. The batch size is 32.

PurdueShapes5GAN Dataset (contd.)



A batch of images from the PurdueShapes5GAN dataset

About the “Complexity” of the Dataset Images

- I would not be surprised if your first reaction to the dataset images is that they couldn't possibly present a great challenge to a data modeler.
- Shown in the next slide are enlarged views of two of the images on the previous slide. In addition to the sharp shape boundaries, you can also small holes inside the shapes.
- The holes that you see inside the shapes were caused by intentionally suppressing bilinear interpolation as the shapes were randomly reoriented.
- So the challenge for the data modeler would be its ability to not only reproduce the shapes while preserving the sharp edges, but also to incorporate the tiny holes inside the shapes, **and do so with the probabilities that reflect the training data.**

About the “Complexity” of the Images (contd.)



PurdueShapes5GAN Dataset (contd.)

You can download the dataset archive

```
datasets_for_AdversarialNetworks.tar.gz
```

through the link "[Download the image dataset for AdversarialNetworks](#)" provided at the top of the HTML version of the main webpage for the DLStudio module (version 2.0.3 or higher). You would need to store it in the [ExamplesAdversarialLearning](#) directory of the distribution. Subsequently, you would need to execute the following command in that directory:

```
tar zxvf datasets_for_AdversarialNetworks.tar.gz
```

This command will create a [dataGAN](#) subdirectory and deposit the following dataset archive in that subdirectory:

```
PurdueShapes5GAN-20000.tar.gz
```

Now execute the following in the [dataGAN](#) directory:

```
tar zxvf PurdueShapes5GAN-20000.tar.gz
```

With that, you should be able to execute the adversarial learning based scripts in the [ExamplesAdversarialLearning](#) directory.

Outline

1	Distance Between Two Probability Distributions	10
2	Total Variation (TV) Distance	12
3	Kullback-Leibler Divergence	16
4	Jensen-Shannon Divergence and Distance	23
5	Earth Mover's Distance	27
6	Wasserstein Distance	37
7	A Random Experiment for Studying Differentiability	44
8	Differentiability of Distance Functions	47
9	PurdueShapes5GAN Dataset for Adversarial Learning	55
10	DCGAN Implementation in DLStudio	62
11	Making Small Changes to the DCGAN Architecture	84
12	Wasserstein GAN Implementation in DLStudio	90
13	Improving Wasserstein GAN with Gradient Penalty	106

DCGAN Implementation in DLStudio

- The main goal of this section is to tell you about the implementation of DCGAN in DLStudio's co-class [AdversarialLearning](#).
- DCGAN, short for "**Deep Convolutional Generative Adversarial Network**", was presented in a paper that I cited in the Preamble to this lecture.
- However, before actually getting into the DCGAN architecture, I need to take you back to the first paper that started the modern excitement in adversarial learning. I am talking about the 2014 publication "[Generative Adversarial Nets](#)" by Goodfellow, Pouget-Abadie, Mirza, Xu, Warde-Farley, Ozair, Courville, and Bengio that was also cited in the Preamble.
- The reason I need to take you back to this paper is because the basic training logic in DCGAN is the same as that proposed in the above cited publication by Goodfellow et al.

Adversarial Learning Requires Generator and Discriminator

- Adversarial learning as described in the Goodfellow et al. paper involves two networks, a Discriminator and a Generator. We can think of the Discriminator as a function $D(x, \theta_d)$ where x is the image and θ_d the weights in the Discriminator network. The $D(x, \theta_d)$ function **returns the probability that the input x is from the probability distribution that describes the training data.**
- Similarly, we can think of the Generator as a function $G(z, \theta_g)$ that maps noise vectors to images that we want to look like the images in our training data. The vector θ_g represents the learnable parameters in the Generator network.
- We assume that the training images are described by some probability distribution that we denote p_{data} . The goal of the Generator is to transform a noise vector, denoted z , into an image that should look like a training image.

Discriminator and Generator (contd.)

- Regarding z , we also assume that the noise vectors z are generated with a probability distribution $p_Z(z)$. Obviously, z is a realization of a vector random variable Z .
- The output of the Generator consists of images that correspond to some probability distribution that we will denote p_G . So you can think of the Generator as a function that transforms the probability distribution p_Z into the distribution p_G .
- The question now is how do we train the Discriminator and the Generator networks.
- The Discriminator is trained to maximize the probability of assigning the correct label to an input image that looks like it came from the same distribution as the training data.

Discriminator Training vs. Generator Training

- That is, for Discriminator training, we want the parameters θ_d to maximize the following expectation:

$$\max_{\theta_d} E_{x \sim p_{data}} [\log D(x)] \quad (31)$$

- The expression $x \sim p_{data}$ means that x was pulled from the distribution p_{data} . In other words, x is one of the training images.
- While we are training D to exhibit the above behavior, we train the Generator for the following minimization:

$$\min_{\theta_g} E_{z \sim p_Z} [\log(1 - D(G(z)))] \quad (32)$$

- Combining the two expressions shown above, we can express the combined optimization as:

$$\min_{\theta_g} \max_{\theta_d} \left[E_{x \sim p_{data}} [\log D(x)] + E_{z \sim p_Z} [\log(1 - D(G(z)))] \right] \quad (33)$$

Discriminator Training vs. Generator Training (contd.)

- We'll translate the min-max form in Eq. (33) into a “protocol” for training the two networks.
- For each training batch of images, we will **first** update the parameters in the Discriminator network and **then** we'll do the same in the Generator network.
- If we use `nn.BCELoss` as the loss criterion for training the Discriminator, that will automatically take care of the logarithms in the expression shown on the previous slide.
- We first train the Discriminator by subjecting it to a maximization that involves the three steps listed on the next slide.
- Subsequently, we train the Generator by a minimization to be described on the slide that follows.

The Two Targets for Discriminator Training

- The maximization steps required for the Discriminator training:
 - ① The maximization of the first term in the expression on the previous slide requires that we use the target "1" for the network output $D(x)$.
 - ② The maximization of the second term in the same expression is a bit more involved since it requires applying the Discriminator network to the output of the Generator for noise input. The second term also requires that we now use "-1" as the target for the Discriminator.

The phrase "we now use -1 as the target for the Discriminator" is to be taken figuratively. Since the Discriminator is a binary classifier (that's what you get with `nn.BCELoss`), its targets can only be 1 and 0. We use 1 as the target in Step 1 and 0 as the target in Step 2.

- ③ After we have calculated the two losses for the Discriminator, we can sum the losses and call `backwards()` on the sum for calculating the gradients of the loss with respect to its weights. A subsequent call to the `step()` of the optimizer would update the weights in the Discriminator network.

The Target for Generator Training

- For the training required for the Generator, only the second term inside the square brackets in Eq. (33) matters. We proceed through the following 4 steps:
 - ① We note that the logarithm is a monotonically increasing function and also because the output $D(G(z))$ in the second term will always be between 0 and 1.
 - ② Therefore, the needed minimization translates into maximizing $D(G(z))$ with respect to a target value of 1.
 - ③ With 1 as the target, we again find the `nn.BCELoss` associated with $D(G(z))$. We call `backwards()` on this loss — **making sure that we have turned off `requires_grad()` on the Discriminator parameters as we are updating the Generator parameters.**
 - ④ A subsequent call to the `step()` for the optimizer would update the weights in the Generator network.

How the GAN Code is Organized in AdversarialLearning

- Now that you have become familiar with the basic idea of Adversarial Learning for data modeling, it's time to get to know better the `AdversarialLearning` co-class in the DLStudio platform.
- All of the GAN related code is in the inner class `DataModeling` of the `AdversarialLearning` class.
- The code in the `DataModeling` class allows you to experiment with the following Discriminator-Generator pairs and Critic-Generator pairs [I'll be talking about "Critics" in the next section on Wasserstein GANs.]:

DG1: This is a Discriminator-Generator pair that corresponds to the original formulation of DCGAN.

DG2: This is a slight variant of the Discriminator-Generator pair in DG1.

CG1: This is a Critic-Generator pair for the Wasserstein GAN in Section 12.

CG2: This is another Critic-Generator pair for the Wasserstein GAN.

DG1: Discriminator and Generator Networks

- Slides 74 and 75 show the DCGAN networks for the DG1 Discriminator-Generator pair.
- **Regarding the Discriminator network on Slide 74**, I refer to the DCGAN network topology as the 4-2-1 network. Each layer of the Discriminator network carries out a strided convolution with a 4×4 kernel, a 2×2 stride, and a 1×1 padding for all but the final layer.
- The output of the final convolutional layer in the Discriminator is pushed through a sigmoid to yield a scalar value as the final output for each image in a batch.
- **Next, on Slide 75, is the implementation of the DCGAN Generator in the example DG1.** As was the case with the Discriminator network, you again see the 4-2-1 topology here.

DG1: Discriminator and Generator (contd.)

- Recall that a Generator's job is to transform a random noise vector into an image that is supposed to look like it came from the training dataset. (Most people refer to the images constructed from noise vectors in this manner as fakes.)
- As you will see in `run_gan_code()`, the starting noise vector is a 1×1 image with 100 channels. In order to output a 64×64 output image from the noise vector, the Generator code shown on Slide 75 uses the Transpose Convolution operator `nn.ConvTranspose2d` with a stride of 2.
- If (H_{in}, W_{in}) are the height and the width of the image at the input to a `nn.ConvTranspose2d` layer and (H_{out}, W_{out}) the same at the output, the input/output sizes are related by [See Slides 45 through 61 of my Week 9 Lecture on Semantic Segmentation]:

$$\begin{aligned}H_{out} &= (H_{in} - 1) * s + k - 2 * p \\W_{out} &= (W_{in} - 1) * s + k - 2 * p\end{aligned}$$

DG1: Discriminator and Generator (contd.)

- In the last bullet on the previous slide, s is the stride and k the size of the kernel. (I am assuming square strides, kernels, and padding).
- Therefore, each `nn.ConvTranspose2d` layer doubles the size of the input.

The Discriminator Network (DG1)

```
##### Discriminator-Generator DG1 #####
class DiscriminatorDG1(nn.Module):

    def __init__(self):
        super(AdversarialLearning.DataModeling.DiscriminatorDG1, self).__init__()
        self.conv_in = nn.Conv2d( 3, 64, kernel_size=4, stride=2, padding=1)
        self.conv_in2 = nn.Conv2d( 64, 128, kernel_size=4, stride=2, padding=1)
        self.conv_in3 = nn.Conv2d( 128, 256, kernel_size=4, stride=2, padding=1)
        self.conv_in4 = nn.Conv2d( 256, 512, kernel_size=4, stride=2, padding=1)
        self.conv_in5 = nn.Conv2d( 512, 1, kernel_size=4, stride=1, padding=0)
        self.bn1 = nn.BatchNorm2d(128)
        self.bn2 = nn.BatchNorm2d(256)
        self.bn3 = nn.BatchNorm2d(512)
        self.sig = nn.Sigmoid()

    def forward(self, x):
        x = torch.nn.functional.leaky_relu(self.conv_in(x), negative_slope=0.2, inplace=True)
        x = self.bn1(self.conv_in2(x))
        x = torch.nn.functional.leaky_relu(x, negative_slope=0.2, inplace=True)
        x = self.bn2(self.conv_in3(x))
        x = torch.nn.functional.leaky_relu(x, negative_slope=0.2, inplace=True)
        x = self.bn3(self.conv_in4(x))
        x = torch.nn.functional.leaky_relu(x, negative_slope=0.2, inplace=True)
        x = self.conv_in5(x)
        x = self.sig(x)
        return x
```

The Generator Network (DG1)

```
class GeneratorDG1(nn.Module):
    def __init__(self):
        super(AdversarialLearning.DataModeling.GeneratorDG1, self).__init__()
        self.latent_to_image = nn.ConvTranspose2d(100, 512, kernel_size=4, stride=1, padding=0, bias=False)
        self.upsampler2 = nn.ConvTranspose2d( 512, 256, kernel_size=4, stride=2, padding=1, bias=False)
        self.upsampler3 = nn.ConvTranspose2d( 256, 128, kernel_size=4, stride=2, padding=1, bias=False)
        self.upsampler4 = nn.ConvTranspose2d( 128, 64,  kernel_size=4, stride=2, padding=1, bias=False)
        self.upsampler5 = nn.ConvTranspose2d( 64,  3,  kernel_size=4, stride=2, padding=1, bias=False)
        self.bn1 = nn.BatchNorm2d(512)
        self.bn2 = nn.BatchNorm2d(256)
        self.bn3 = nn.BatchNorm2d(128)
        self.bn4 = nn.BatchNorm2d(64)
        self.tanh = nn.Tanh()

    def forward(self, x):
        x = self.latent_to_image(x)
        x = torch.nn.functional.relu(self.bn1(x))
        x = self.upsampler2(x)
        x = torch.nn.functional.relu(self.bn2(x))
        x = self.upsampler3(x)
        x = torch.nn.functional.relu(self.bn3(x))
        x = self.upsampler4(x)
        x = torch.nn.functional.relu(self.bn4(x))
        x = self.upsampler5(x)
        x = self.tanh(x)
        return x
```

The Training Loop for DCGAN (DG1)

- The code shown on Slides 78 through 80 implements the training logic presented on Slides 67 through 69. It is meant for training a Discriminator-Generator based Adversarial Network. The implementation shown has borrowed several programming constructs from the "official" DCGAN implementation at GitHub.
- Sections of the training loop that begin in Lines (A) and (B) are for the Discriminator part of the training in Eq. (33). The statements in Part 1(a) implement the logic in the first bullet under Discriminator training on Slide 68. In these statements we use the target of "1" for the output of the Discriminator when it is invoked on a data image.
- The statements in Part 1(b) that begin at Line (B) implement the logic in the second bullet on the Slide 68. That is, now we subject the output of the Discriminator after it is applied to the Generator images to the target "-1".

The Training Loop for DCGAN (DG1) (contd.)

- The section of the code that begins in Line (C) is for Generator training through the steps outlined on Slide 69. The min part in Eq. (33) on Slide 66 requires that we minimize $1 - D(G(z))$ which, since D is constrained to lie in the interval $(0,1)$, requires that we maximize $D(G(z))$. We accomplish that by applying the Discriminator to the output of the Generator and use 1 as the target for each image, as mentioned in the second bullet on Slide 69.

The Training Loop for DCGAN (DG1) (contd.)

```

def run_gan_code(self, dlstudio, advers, discriminator, generator, results_dir):
    # Set the number of channels for the 1x1 input noise vectors for the Generator:
    nz = 100
    netD = discriminator.to(advers.device)
    netG = generator.to(advers.device)
    # Initialize the parameters of the Discriminator and the Generator networks according to the
    # definition of the "weights_init()" method:
    netD.apply(self.weights_init)
    netG.apply(self.weights_init)
    # We will use the same noise batch to periodically check on the progress made for the Generator:
    fixed_noise = torch.randn(self.dlstudio.batch_size, nz, 1, 1, device=advers.device)
    # Establish convention for real and fake labels during training
    real_label = 1
    fake_label = 0
    # Adam optimizers for the Discriminator and the Generator:
    optimizerD = optim.Adam(netD.parameters(), lr=dlstudio.learning_rate, betas=(advers.beta1, 0.999))
    optimizerG = optim.Adam(netG.parameters(), lr=dlstudio.learning_rate, betas=(advers.beta1, 0.999))
    # Establish the criterion for measuring the loss at the output of the Discriminator network:
    criterion = nn.BCELoss()
    # We will use these lists to store the results accumulated during training:
    img_list = []
    G_losses = []
    D_losses = []
    iters = 0
    print("\n\nStarting Training Loop...\n\n")
    start_time = time.perf_counter()
  
```

(Continued on the next slide

The Training Loop for DCGAN (DG1) (contd.)

(..... continued from the previous slide)

```

for epoch in range(dlstudio.epochs):
    g_losses_per_print_cycle = []
    d_losses_per_print_cycle = []
    for i, data in enumerate(self.train_data_loader, 0):
        ## Part 1(a) of Training (maximization of minmax objective for the Discriminator):    ## (A)
        netD.zero_grad()
        real_images_in_batch = data[0].to(advers.device)
        # Need to know how many images we pulled in since at the tailend of the dataset,
        # the number of images may not equal the user-specified batch size:
        b_size = real_images_in_batch.size(0)
        label = torch.full((b_size,), real_label, dtype=torch.float, device=advers.device)
        output = netD(real_images_in_batch).view(-1)
        errD_reals = criterion(output, label)
        errD_reals.backward()
        ## Part 1(b) of Training (maximization of the minmax object for the Discriminator
        ##                               when applied to fakes):    ## (B)
        noise = torch.randn(b_size, nz, 1, 1, device=advers.device)
        fakes = netG(noise)
        label.fill_(fake_label)
        ## The call to fakes.detach() in the next statement returns a copy of the 'fakes' tensor
        ## such that the copy that is returned does not exist in the computational graph. That is,
        ## the copy of the tensor is removed from the computational graph. However, the original
        ## 'fakes' tensor continues to remain in the computational graph. This play ensures that
        ## a subsequent call to backward() in the 3rd statement below would only result in a
        ## calculation of the gradients for the netD weights:
        output = netD(fakes.detach()).view(-1)
        errD_fakes = criterion(output, label)
        errD_fakes.backward()
        errD = errD_reals + errD_fakes
        d_losses_per_print_cycle.append(errD)
    optimizerD.step()    ## Only the Discriminator weights are incremented

```

The Training Loop for DCGAN (DG1) (contd.)

(..... continued from the previous slide)

```

## Part 2 of Training (minimization of the minmax objective for learning
##                                     the Generator):                ## (C)
##
## The min part requires that we MINIMIZE "1 - D(G(z))" which, since D is constrained to
## lie in the interval (0,1), requires that we maximize D(G(z)). We accomplish that by
## applying the Discriminator to the output of the Generator and use 1 as the target:
netG.zero_grad()
label.fill_(real_label)
output = netD(fakes).view(-1)
errG = criterion(output, label)
g_losses_per_print_cycle.append(errG)
errG.backward()
optimizerG.step()

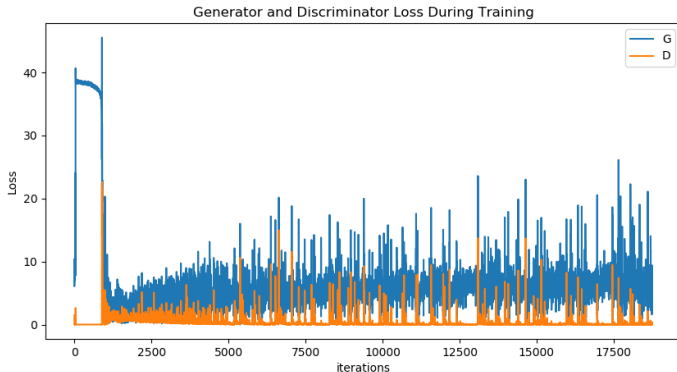
if i % 100 == 99:
    current_time = time.perf_counter()
    elapsed_time = current_time - start_time
    mean_D_loss = torch.mean(torch.FloatTensor(d_losses_per_print_cycle))
    mean_G_loss = torch.mean(torch.FloatTensor(g_losses_per_print_cycle))
    print("[epoch=%d/%d  iter=%4d  elapsed_time=%5d secs]      mean_D_loss=%7.4f
                                                mean_G_loss=%7.4f" %
          ((epoch+1), dlstudio.epochs, (i+1), elapsed_time, mean_D_loss, mean_G_loss))
    d_losses_per_print_cycle = []
    g_losses_per_print_cycle = []

```

NOTES:

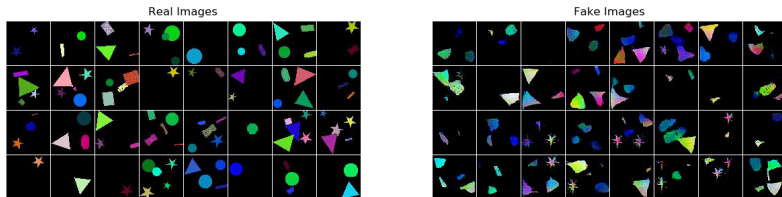
- A statement like `label = torch.full((b.size,), real_label)` means that we want to set `label` to a single-axis tensor of size `b.size` and we want all its elements to be set to the value given by `real_label`.
- A statement like `label.fill_(value)` means that the previously declared tensor `label` needs to be filled in-place with the specified value.

Losses vs. Iterations for DG1



Discriminator and Generator losses over 30 epochs of training

Comparing Real and Fake Images for DG1



At the end of 30 epochs of training, shown at left is a batch of real images and, at right, the images produced by the Generator from noise vectors

An Animated GIF of the Generator Output for DG1

The following animated GIF shows how the Generator's output evolves over 30 epochs using the same set of noise vectors.

https://engineering.purdue.edu/DeepLearn/pdf-kak/DG1_generation_animation.gif

Outline

1	Distance Between Two Probability Distributions	10
2	Total Variation (TV) Distance	12
3	Kullback-Leibler Divergence	16
4	Jensen-Shannon Divergence and Distance	23
5	Earth Mover's Distance	27
6	Wasserstein Distance	37
7	A Random Experiment for Studying Differentiability	44
8	Differentiability of Distance Functions	47
9	PurdueShapes5GAN Dataset for Adversarial Learning	55
10	DCGAN Implementation in DLStudio	62
11	Making Small Changes to the DCGAN Architecture	84
12	Wasserstein GAN Implementation in DLStudio	90
13	Improving Wasserstein GAN with Gradient Penalty	106

Making Small Changes to the DCGAN Architecture (DG2)

- My personal experience with the DCGAN architecture is that when it works, it produces beautiful results. However, as you change the initializations for the parameters, or as you make minor tweaks to the Generator and/or the Discriminator network, **more often than not, what you get is what is known as mode collapse.** Mode collapse means that the different randomly chosen noise vectors for the input to the Generator will yield the same garbage output.
- To illustrate what I mean, The Discriminator network shown on the next slide is the same as the one you saw earlier for the DCGAN implementation, except for the additional layer `self.extra` that the incoming image is routed through at the beginning of the network in `forward()`
- I have also defined a batch normalization layer `self.bnX` for the output of the extra layer `self.extra`.

```

##### Discriminator-Generator DG2 #####
class DiscriminatorDG2(nn.Module):
    """
    This is essentially the same network as the DCGAN for DG1, except for the extra layer
    "self.extra" shown below. We also declare a batchnorm for this extra layer in the form
    of "self.bnX". In the implementation of "forward()", we invoke the extra layer at the
    beginning of the network.
    """
    def __init__(self, skip_connections=True, depth=16):
        super(AdversarialLearning.DataModeling.DiscriminatorDG2, self).__init__()
        self.conv_in = nn.Conv2d( 3, 64, kernel_size=4, stride=2, padding=1)
        self.extra = nn.Conv2d( 64, 64, kernel_size=4, stride=1, padding=2)
        self.conv_in2 = nn.Conv2d( 64, 128, kernel_size=4, stride=2, padding=1)
        self.conv_in3 = nn.Conv2d( 128, 256, kernel_size=4, stride=2, padding=1)
        self.conv_in4 = nn.Conv2d( 256, 512, kernel_size=4, stride=2, padding=1)
        self.conv_in5 = nn.Conv2d( 512, 1, kernel_size=4, stride=1, padding=0)
        self.bn1 = nn.BatchNorm2d(128)
        self.bn2 = nn.BatchNorm2d(256)
        self.bn3 = nn.BatchNorm2d(512)
        self.bnX = nn.BatchNorm2d(64)
        self.sig = nn.Sigmoid()

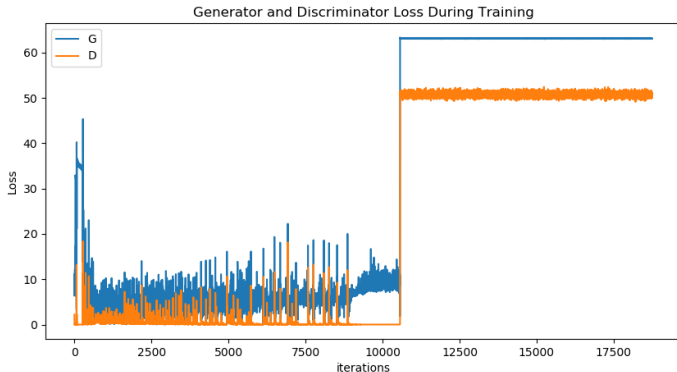
    def forward(self, x):
        x = torch.nn.functional.leaky_relu(self.conv_in(x), negative_slope=0.2, inplace=True)
        x = self.bnX(self.extra(x))
        x = torch.nn.functional.leaky_relu(x, negative_slope=0.2, inplace=True)
        x = self.bn1(self.conv_in2(x))
        x = torch.nn.functional.leaky_relu(x, negative_slope=0.2, inplace=True)
        x = self.bn2(self.conv_in3(x))
        x = torch.nn.functional.leaky_relu(x, negative_slope=0.2, inplace=True)
        x = self.bn3(self.conv_in4(x))
        x = torch.nn.functional.leaky_relu(x, negative_slope=0.2, inplace=True)
        x = self.conv_in5(x)
        x = self.sig(x)
        return x

class GeneratorDG2(nn.Module):
    """
    The Generator for DG2 is exactly the same as for the DG1. So please the comment block for that
    Generator.
    """
    def __init__(self):
        super(AdversarialLearning.DataModeling.GeneratorDG2, self).__init__()
        self.latent_to_image = nn.ConvTranspose2d( 100, 512, kernel_size=4, stride=1, padding=0, bias=False)
        self.upsampler2 = nn.ConvTranspose2d( 512, 256, kernel_size=4, stride=2, padding=1, bias=False)
        self.upsampler3 = nn.ConvTranspose2d( 256, 128, kernel_size=4, stride=2, padding=1, bias=False)
        self.upsampler4 = nn.ConvTranspose2d( 128, 64, kernel_size=4, stride=2, padding=1, bias=False)
        self.upsampler5 = nn.ConvTranspose2d( 64, 3, kernel_size=4, stride=2, padding=1, bias=False)
        self.bn1 = nn.BatchNorm2d(512)
        self.bn2 = nn.BatchNorm2d(256)
        self.bn3 = nn.BatchNorm2d(128)
        self.bn4 = nn.BatchNorm2d(64)
        self.tanh = nn.Tanh()

    def forward(self, x):
        x = self.latent_to_image(x)
        x = torch.nn.functional.relu(self.bn1(x))
        x = self.upsampler2(x)
        x = torch.nn.functional.relu(self.bn2(x))
        x = self.upsampler3(x)
        x = torch.nn.functional.relu(self.bn3(x))
        x = self.upsampler4(x)
        x = torch.nn.functional.relu(self.bn4(x))
        x = self.upsampler5(x)
        x = self.tanh(x)
        return x

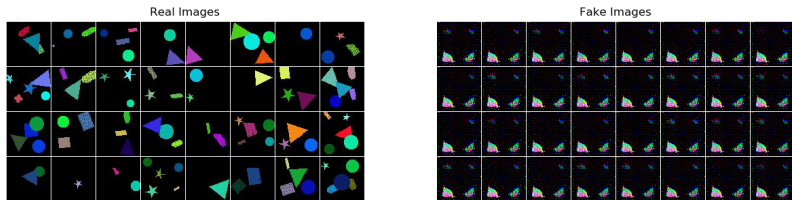
```

Losses vs. Iterations for DG2



Discriminator and Generator losses over 30 epochs of training

Comparing Real and Fake Images for DG2



At the end of 30 epochs of training, shown at left is a batch of real images and, at right, the images produced by the Generator from noise vectors

An Animated GIF of the Generator Output for DG2

The following animated GIF shows how the Generator's output evolves over 30 epochs using the same set of noise vectors for the case of a DCGAN with relatively minor alterations.

https://engineering.purdue.edu/DeepLearn/pdf-kak/DG2_generation_animation.gif

Outline

1	Distance Between Two Probability Distributions	10
2	Total Variation (TV) Distance	12
3	Kullback-Leibler Divergence	16
4	Jensen-Shannon Divergence and Distance	23
5	Earth Mover's Distance	27
6	Wasserstein Distance	37
7	A Random Experiment for Studying Differentiability	44
8	Differentiability of Distance Functions	47
9	PurdueShapes5GAN Dataset for Adversarial Learning	55
10	DCGAN Implementation in DLStudio	62
11	Making Small Changes to the DCGAN Architecture	84
12	Wasserstein GAN Implementation in DLStudio	90
13	Improving Wasserstein GAN with Gradient Penalty	106

Wasserstein GAN Implementation in DLStudio

- This implementation is based on the paper "Wasserstein GAN" by Arjovsky, Chintala, and Bottou that I cited previously in the Preamble.
- You will find my implementation of Wasserstein GAN (WGAN) in DLStudio's co-class [AdversarialLearning](#).
- As you would expect, WGAN is based on estimating the Wasserstein distance between the distribution that corresponds to the training images and the distribution that has been learned so far by the Generator. This distance was defined in Eq. (23) on Slide 40.
- The 1-Lipschitz function $f()$ that is required by the definition in Eq. (23) is implemented as a Critic — because, unlike what was the case for the Discriminator, the job of the Critic is NOT to accept or reject what is produced by the Generator, but to do what's mentioned on the next slide.

WGAN Implementation in DLStudio (contd.)

- In a WGAN, a Critic's job is to become adept at estimating the Wasserstein distance between the distribution that corresponds to the training dataset and the distribution that has been learned by the Generator so far.
- Since the Wasserstein distance is known to be differentiable with respect to the learnable weights in the Critic network, one can backprop the distance and update the weights in an iterative training loop. This is roughly the idea of the Wasserstein GAN that is incorporated as a Critic-Generator pair CG1 in the Adversarial Networks class.
- For the purpose of implementation, here is a rewrite of the Wasserstein distance presented earlier in Eq. (23) on Slide 40:

$$d_W(P_r, P_\theta) = \sup_{\|f\|_L \leq 1} \left[E_{x \sim P_r} \{f_w(x)\} - E_{z \sim P_z} \{f_w(g_\theta(z))\} \right] \quad (34)$$

WGAN Implementation in DLStudio (contd.)

- In the formula for Wasserstein distance shown on the previous slide, P_r is the “real” distribution that describes the training data and P_z describes the distribution of the noise vectors that are fed into the Generator for the production of the fake images. The Generator parameters are denoted θ and $g_\theta()$ stands for the function that describes the behavior of the Generator.
- Now that we have interpreted the role of the function $f_w()$ as a Critic — the Critic’s job being to learn the function $f_w()$ — the question is how does the Critic make sure that the function being learned is 1-Lipschitz?
- A heuristic answer to the vexing question posed above was provided by the original authors the “Wasserstein GAN” paper. For lack of any available well-principled approach as a solution to this issue, they experimented with tightly clipping the values being learned for the weights in the Critic network.

WGAN Implementation in DLStudio (contd.)

- It stands to reason that the closer the clipping level is to zero from both the positive and the negative sides, the less likely that the gradient of the function being learned will exhibit large swings.
- The calculation of the Wasserstein distance using Eq. (34) also calls for averaging of the output of the Critic in order for the maximization to yield the desired distance. This can be taken care of by having the Critic go through multiple iterations of the update of its parameters for each iteration for the Generator.
- For implementation, the expression for the Wasserstein distance shown in Eq. (34) can be rewritten as:

$$d_W(P, Q) = \max_{\|f\|_L \leq 1} \left[E_{x \sim P} \{f(x)\} - E_{y \sim Q} \{f(y)\} \right] \quad (35)$$

WGAN Implementation in DLStudio (contd.)

- Note that Eq. (35) can also be interpreted as: There is guaranteed to exist a 1-Lipschitz continuous function $f()$ that when applied to the samples drawn from the distributions P and Q will yield the Wasserstein distance between the two distributions.
- Let C denote a Critic network that can learn the function $f()$. Remember, our overarching goal remains that we need to also learn a Generator network G that is capable of converting noise into samples that look like those from the distribution P .
- We seek to create a GAN that can learn a G that MINIMIZES the Wasserstein distance between the true distribution P and its learned approximation Q . At the same time, the GAN must discover a C that seeks to maximize the same distance (in the sense that the Critic learns how to maximally distrust the Generator G).

WGAN Implementation in DLStudio (contd.)

- We thus end up with the following minimax objective for the learning framework:

$$\min_G \max_C \left[E_{x \sim P} [C(x)] - E_{z \sim p_Z} [C(G(z))] \right] \quad (36)$$

- In comparing this minimax objective with the one shown earlier in Eq. (33) of Section 10, note that the two components of the argument to the minimax in that equation were additive, whereas we subtract them in the objective shown above. In Eq. (33), we had a Discriminator in the GAN and our goal was to maximize its classification performance for images that look like they came from the true distribution P . On the other hand, the goal of the Critic here is to learn to maximize the Wasserstein distance between the true distribution P and its learned approximation Q . Note that the distribution Q is for the images that are constructed by the Generator from the white-noise samples z drawn from a distribution p_Z , as shown above.

WGAN Implementation in DLStudio (contd.)

- As far as the Critic is concerned, the maximization needed in Eq. (36) can be achieved by using the following loss function:

$$\begin{aligned}\text{Critic Loss} &= E_{y \sim Q}[C(y)] - E_{x \sim P}[C(x)] \\ &= E_{z \sim p_z}[C(G(z))] - E_{x \sim P}[C(x)]\end{aligned}\tag{37}$$

- In the WGAN code shown in what follows, this is accomplished by using a "gradient target" of +1 for the mean of the output of the Critic when it sees the images produced by the Generator and the "gradient target" of -1 for the output of the Critic when it sees the training data directly.
- As to why we use the gradient targets of +1 and -1, it was shown by the original authors of WGAN that the optimal Critic C has unit gradient norm almost everywhere under P and Q . That is, the magnitude of the partial derivative of the output of the optimal C with respect to its input will almost always be 1.

The Critic and the Generator in DLStudio's WGAN

```
##### Critic-Generator CG1 #####
class CriticCG1(nn.Module):
    """
    I have used the SkipBlockDN as a building block for the Critic network. This I did with the hope
    that when time permits I may want to study the effect of skip connections on the behavior of the
    the critic vis-a-vis the Generator. The final layer of the network is the same as in the
    "official" GItHub implementation of Wasserstein GAN. And, as in WGAN, I have used the Leaky ReLU
    for activation.
    """
    def __init__(self):
        super(AdversarialLearning.DataModeling.CriticCG1, self).__init__()
        self.conv_in = AdversarialLearning.DataModeling.SkipBlockDN(3, 64, downsample=True, skip_connections=True)
        self.conv_in2 = AdversarialLearning.DataModeling.SkipBlockDN(64, 128, downsample=True, skip_connections=False)
        self.conv_in3 = AdversarialLearning.DataModeling.SkipBlockDN(128, 256, downsample=True, skip_connections=False)
        self.conv_in4 = AdversarialLearning.DataModeling.SkipBlockDN(256, 512, downsample=True, skip_connections=False)
        self.conv_in5 = AdversarialLearning.DataModeling.SkipBlockDN(512, 1, downsample=False, skip_connections=False)
        self.bn1 = nn.BatchNorm2d(128)
        self.bn2 = nn.BatchNorm2d(256)
        self.bn3 = nn.BatchNorm2d(512)
        self.final = nn.Linear(512, 1)

    def forward(self, x):
        x = torch.nn.functional.leaky_relu(self.conv_in(x), negative_slope=0.2, inplace=True)
        x = self.bn1(self.conv_in2(x))
        x = torch.nn.functional.leaky_relu(x, negative_slope=0.2, inplace=True)
        x = self.bn2(self.conv_in3(x))
        x = torch.nn.functional.leaky_relu(x, negative_slope=0.2, inplace=True)
        x = self.bn3(self.conv_in4(x))
        x = torch.nn.functional.leaky_relu(x, negative_slope=0.2, inplace=True)
        x = self.conv_in5(x)
        x = x.view(-1)
        x = self.final(x)
        x = x.mean(0)
        x = x.view(1)
        return x

class GeneratorCG1(nn.Module):
    """
    The Generator code remains the same as for the DCGAN shown earlier.
    """
    def __init__(self):
        super(AdversarialLearning.DataModeling.GeneratorCG1, self).__init__()
        self.latent_to_image = nn.ConvTranspose2d(100, 512, kernel_size=4, stride=1, padding=0, bias=False)
        self.upsampler2 = nn.ConvTranspose2d(512, 256, kernel_size=4, stride=2, padding=1, bias=False)
        self.upsampler3 = nn.ConvTranspose2d(256, 128, kernel_size=4, stride=2, padding=1, bias=False)
        self.upsampler4 = nn.ConvTranspose2d(128, 64, kernel_size=4, stride=2, padding=1, bias=False)
        self.upsampler5 = nn.ConvTranspose2d(64, 3, kernel_size=4, stride=2, padding=1, bias=False)
        self.bn1 = nn.BatchNorm2d(512)
        self.bn2 = nn.BatchNorm2d(256)
        self.bn3 = nn.BatchNorm2d(128)
        self.bn4 = nn.BatchNorm2d(64)
        self.tanh = nn.Tanh()

    def forward(self, x):
        x = self.latent_to_image(x)
        x = torch.nn.functional.relu(self.bn1(x))
        x = self.upsampler2(x)
        x = torch.nn.functional.relu(self.bn2(x))
        x = self.upsampler3(x)
        x = torch.nn.functional.relu(self.bn3(x))
        x = self.upsampler4(x)
        x = torch.nn.functional.relu(self.bn4(x))
        x = self.upsampler5(x)
        x = self.tanh(x)
        return x
```

Training the WGAN

- The code for training the Critic-Generator based WGAN shown next is based on the logic of a Wasserstein GAN as proposed by the original authors of WGAN. The implementation shown uses several programming constructs from the WGAN implementation at GitHub. I have also used several programming constructs from the DCGAN code at GitHub.
- The noise batch that is generated in Line (D) is used periodically check on the progress made by the Generator.
- The 'one' and 'minus_one' you see in Lines (E) and (F) are for training the Critic, 'minus_one' is for the part of the training with actual training images, and 'one' is for the part based on the images produced by the Generator.
- The inner 'while' loop in Line (G) is for updating the Critic in such a way that the discrimination function learned by the Critic satisfies the 1-Lipschitz condition.

Training the WGAN (contd.)

- The 1-Lipschitz condition is enforced by the clipping statements in Lines (H) and (I) along with the smoothing action of the inner 'while' loop.
- As mentioned previously, a minimization of the Wasserstein distance between the distribution that describes the training data and the distribution that has been learned so far by the Generator **can be translated into a *maximization* of the difference of the *average outputs* of a 1-Lipschitz function as applied to the training images and as applied to the output of the Generator.** Learning this 1-Lipschitz function is the job of the Critic.
- Training the Critic consists of two parts. In the first part that begins in Line (J), we apply the target 'one' to the training images and, in the second part that begins in Line (K), we use the target 'minus_one' for the output of the Critic when its input is the output of the Generator.

Training the WGAN (contd.)

- That brings us to the training of the Generator that begins in Line (L). We must start by turning off the `requires_grad` of the Critic parameters since the Critic and the Generator are meant to be updated independently.

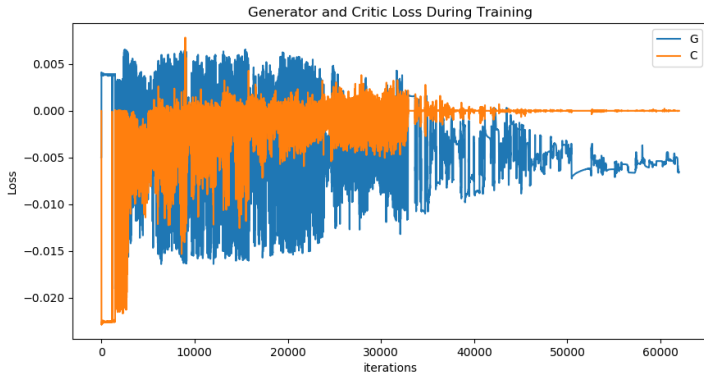
Training the WGAN

```

def run_wgan_code(self, dlstudio, adversarial, critic, generator, results_dir):
    nz = 100 # Set the number of channels for the 1x1 input noise vectors for the Generator ## (A)
    netC = critic.to(adversarial.device)
    netG = generator.to(adversarial.device)
    netC.apply(self.weights_init) # initialize Critic network parameters ## (B)
    netG.apply(self.weights_init) # initialize Generator network parameters ## (C)
    fixed_noise = torch.randn(self.dlstudio.batch_size, nz, 1, 1, device=adversarial.device) ## (D)
    one = torch.FloatTensor([1]).to(adversarial.device) ## (E)
    minus_one = torch.FloatTensor([-1]).to(adversarial.device) ## (F)
    # Adam optimizers for the Critic and the Generator:
    optimizerC = optim.Adam(netC.parameters(), lr=dlstudio.learning_rate, betas=(adversarial.beta1, 0.999))
    optimizerG = optim.Adam(netG.parameters(), lr=dlstudio.learning_rate, betas=(adversarial.beta1, 0.999))
    img_list = []
    Gen_losses = []
    Cri_losses = []
    iters = 0
    gen_iterations = 0
    start_time = time.perf_counter()
    dataloader = self.train_dataloader
    clipping_thresh = self.adversarial.clipping_threshold
    for epoch in range(dlstudio.epochs):
        data_iter = iter(dataloader)
        i = 0
        ncritic = 5
        while i < len(dataloader):
            for p in netC.parameters():
                p.requires_grad = True
            if gen_iterations < 25 or gen_iterations % 500 == 0: # the choices 25 and 500 are from WGAN
                ncritic = 100
            ic = 0
            ## The inner 'while' loop shown below calculates the expectations in Eq. (8) in the doc section
            ## at the beginning of this file:
            while ic < ncritic and i < len(dataloader):
                ic += 1 ## (G)
                for p in netC.parameters():
                    p.data.clamp_(-clipping_thresh, clipping_thresh) ## (H)
                ## Training the Critic with real images (Part 1): ## (I)
                netC.zero_grad()
                real_images_in_batch = data_iter.next()
                i += 1
                real_images_in_batch = real_images_in_batch[0].to(self.device)
                # Need to know how many images we pulled in since at the tailend of the dataset, the
                # number of images may not equal the user-specified batch size:
                b_size = real_images_in_batch.size(0)
                # Note that a single scalar is produced for all the data in a batch. This is probably
                # the reason why what the Generator learns is somewhat fuzzy.
                critic_for_reals_mean = netC(real_images_in_batch)
                ## 'minus_one' is the gradient target:
                critic_for_reals_mean.backward(minus_one)
                ## Training the Critic with fake images (Part 2): ## (K)
                noise = torch.randn(b_size, nz, 1, 1, device=self.device)
                fakes = netG(noise)
                # Again, a single number is produced for the whole batch:
                critic_for_fakes_mean = netC(fakes)
                ## 'one' is the gradient target:
                critic_for_fakes_mean.backward(one)
                wasser_dist = critic_for_reals_mean - critic_for_fakes_mean
                loss_critic = critic_for_fakes_mean - critic_for_reals_mean
                # Update the Critic
                optimizerC.step()
            ## Training the Generator: ## (L)
            for p in netG.parameters():
                p.requires_grad = False
            netG.zero_grad()
            # This is again a single scalar based characterization of the whole batch of the Generator images:
            noise = torch.randn(b_size, nz, 1, 1, device=self.device)
            fakes = netG(noise)
            critic_for_fakes_mean = netC(fakes)
            loss_gen = critic_for_fakes_mean
            critic_for_fakes_mean.backward(minus_one)
            # Update the Generator
            optimizerG.step()
            gen_iterations += 1

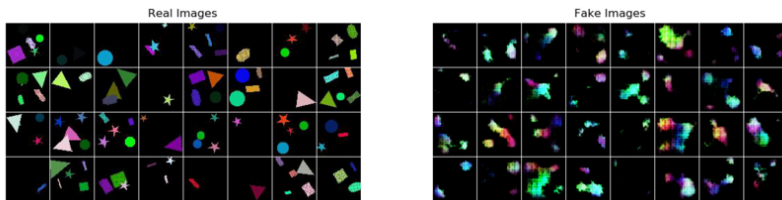
```

Losses vs. Iterations for WGAN



Critic and Generator losses over 500 epochs of training

Comparing Real and Fake Images for WGAN



At the end of 500 epochs of training, shown at left is a batch of real images and, at right, the images produced by the Generator from noise vectors

An Animated GIF of the Generator Output for WGAN

The following animated GIF shows how the Generator's output evolves over 30 epochs using the same set of noise vectors for the case of a DCGAN with relatively minor alterations.

https://engineering.purdue.edu/DeepLearn/pdf-kak/WGAN_generation_animation.gif

Outline

1	Distance Between Two Probability Distributions	10
2	Total Variation (TV) Distance	12
3	Kullback-Leibler Divergence	16
4	Jensen-Shannon Divergence and Distance	23
5	Earth Mover's Distance	27
6	Wasserstein Distance	37
7	A Random Experiment for Studying Differentiability	44
8	Differentiability of Distance Functions	47
9	PurdueShapes5GAN Dataset for Adversarial Learning	55
10	DCGAN Implementation in DLStudio	62
11	Making Small Changes to the DCGAN Architecture	84
12	Wasserstein GAN Implementation in DLStudio	90
13	Improving Wasserstein GAN with Gradient Penalty	106

WGAN-GP: Improving WGAN with Gradient Penalty

- As you would guess, the name extension "-GP" stands for "Gradient Penalty".
- It was shown by the authors Gulrajani, Ahmed, Arjovsky, Dumouli, and Courville of the paper "Improved Training of Wasserstein GANs" that implementing a 1-Lipschitz constraint with weight clipping as discussed in the previous section biases the Critic towards learning rather simple probability distribution functions.
- In WGAN-GP, the performance of a WGAN is improved by putting to use the theoretical property that the optimal WGAN critic C has unit gradient norm almost everywhere under P and Q . [See Proposition 1, Corollary 1 of the paper cited above.]

WGAN-GP (contd.)

- On the basis of the property mentioned at the bottom of the previous slide, in a WGAN-GP, we add a Gradient Penalty term to the Critic Loss that was shown earlier in Eq. (37):

$$\text{Critic Loss} = \underbrace{E_{z \sim p_z}[C(G(z))] - E_{x \sim P}[C(x)]}_{\text{The original critic loss}} + \underbrace{\lambda[\|\nabla_{\hat{x}} C(\hat{x})\|^2 - 1]^2}_{\text{The Gradient Penalty (GP)}} \quad (38)$$

- To explain what the symbol \hat{x} is doing in the GP term, note that the gradient is of the output of the 1-Lipschitz function (meaning the output of the Critic network) with respect to its input. Since the Critic network sees both the training samples and those produced by the Generator at its input, for the purpose of calculating this gradient, we first construct a fictitious sample by taking a weighted sum of a sample drawn from the training data and one produced by the Generator using a randomly chosen fractional number ϵ :

$$\hat{x} = \epsilon x + (1 - \epsilon)\tilde{x} \quad (39)$$

WGAN-GP (contd.)

- Shown below is the Tensorflow code for calculating the Gradient Penalty as posted by the authors of the "Improved Training of Wasserstein GANs" paper:

```
if MODE == 'wgan-gp':
    epsilon = tf.random_uniform( shape=[BATCH_SIZE,1], minval=0., maxval=1. )
    interpolates = epsilon*real_data + ((1-epsilon)*fake_data)
    disc_interpolates = Discriminator(interpolates)
    gradients = tf.gradients(disc_interpolates, [interpolates])[0]
    slopes = tf.sqrt(tf.reduce_sum(tf.square(gradients), reduction_indices=[1]))
    gradient_penalty = tf.reduce_mean((slopes-1)**2)
```

- Shown below is the PyTorch version of the same code as posted by Marvin Cao (caogang) at GitHub:

```
def calc_gradient_penalty(netC, real_data, fake_data):
    epsilon = torch.rand(batch_size, 1).cuda()
    epsilon = epsilon.expand(real_data.size())
    interpolates = epsilon * real_data + ((1 - epsilon) * fake_data)
    interpolates = interpolates.requires_grad_(True).cuda()
    critic_interpolates = netC(interpolates)
    gradients = autograd.grad(outputs=critic_interpolates, inputs=interpolates,
                              grad_outputs=torch.ones(critic_interpolates.size()).cuda(),
                              create_graph=True, retain_graph=True, only_inputs=True)[0]
    gradient_penalty = ((gradients.norm(2, dim=1) - 1) ** 2).mean() * LAMBDA
    return gradient_penalty
```

WGAN-GP for Learning a Point Distribution in 2D

- In order to demonstrate how effective the gradient penalty is in improving the performance of a WGAN, I'll use the 8-Gaussian example from the "Improved Training of Wasserstein GANs" paper. 8-Gaussian refers to a multi-Gaussian distribution of points in an xy-plane. The centers of the eight Gaussians are equispaced on a unit circle around the origin of the plane. The width of each Gaussian is specified by the user. The code snippet shown below returns a batch of 256 points in the xy-plane each time the function `multi_gaussian_source()` is called:

```
def multi_gaussian_source():
    """
    A Python 'generator' function: Each call to this function with the built-in "next()" will yield
    a fresh BATCH_SIZE (typically 256) number of points in the xy-plane.
    """
    scale = 2.
    centers = [ (1, 0), (-1, 0), (0, 1), (0, -1), (1. / np.sqrt(2), 1. / np.sqrt(2)),
               (1. / np.sqrt(2), -1. / np.sqrt(2)), (-1. / np.sqrt(2), 1. / np.sqrt(2)),
               (-1. / np.sqrt(2), -1. / np.sqrt(2))
             ]
    centers = [(scale * x, scale * y) for x, y in centers]
    while True:
        dataset = []
        #spread = 0.02
        spread = 0.1                    ## controls the spread of each Gaussian
        for i in range(BATCH_SIZE):
            point = np.random.randn(2) * spread
            center = random.choice(centers)
            point[0] += center[0]
            point[1] += center[1]
            dataset.append(point)
        dataset = np.array(dataset, dtype='float32')
        dataset /= 1.414 # stdev
        yield dataset
```

WGAN-GP for Learning a Point Distro (contd.)

- Given the data source shown on the previous slide for the ground-truth, we want to train a WGAN so that its Generator would transform noise into data samples (points in the xy-plane) that look like they came from the 8-Gaussians distribution. For the WGAN, we will use the Generator and the Critic classes as shown below:

```

class Generator(nn.Module):
    def __init__(self):
        super(Generator, self).__init__()
        main = nn.Sequential(
            nn.Linear(2, DIM),
            nn.ReLU(True),
            nn.Linear(DIM, DIM),
            nn.ReLU(True),
            nn.Linear(DIM, DIM),
            nn.ReLU(True),
            nn.Linear(DIM, 2),
        )
        self.main = main

    def forward(self, noise):
        output = self.main(noise)
        return output

class Critic(nn.Module):
    def __init__(self):
        super(Critic, self).__init__()
        main = nn.Sequential(
            nn.Linear(2, DIM),
            nn.ReLU(True),
            nn.Linear(DIM, DIM),
            nn.ReLU(True),
            nn.Linear(DIM, DIM),
            nn.ReLU(True),
            nn.Linear(DIM, 1),
        )
        self.main = main

    def forward(self, inputs):
        output = self.main(inputs)
        return output

```

WGAN-GP for Learning a Point Distro (contd.)

- As you can see on the previous slide, except for the last layer, the network layout for both the Generator and the Critic are identical. The output of the Generator is 2D because it is supposed to generate points in the xy -plane. On the other hand, the output of the Critic is a 1D value that expresses the Critic's confidence that the input is genuine or fake.
- The next code segment that follows is about the training of the WGAN. A highlight of the code shown is three calls to `backward()` for estimating the gradients of the Critic weights in lines (F), (K), and (M) and one call to the same for the Generator in line (U).
- To elaborate the code shown for WGAN training, the main loop starts in line (A). Each iteration of the main training loop involves training the Generator network once. At the same time, it requires that the Critic be taken through multiple updates in keeping with the requirements of the expectation operator in Eq. (37) and Eq. (38).

WGAN-GP for Learning a Point Distro (contd.)

- The multiple updates of the Critic in the inner loop start in line (B). The needs of the expectation operator in Eq. (38) are met by averaging both over multiple iterations of the inner loop that starts in line (B) and, in each of those iterations, by averaging over all the samples in a batch, as you will soon see.
- Each inner-loop update of the Critic entails first feeding it a batch of real data (typically 256 points in the xy-plane) in line (D). In keeping with the requirements of the expectation in Eq. (38), we find the mean of the output of the Critic for all the samples in the batch in line (E). The call to `backward()` in line (F) updates the gradients of the Critic weights for this phase of learning for the Critic. Note the target gradient of “-1” in the call to `backward()` in line (F).
- For the next phase of Critic learning, we feed it a batch of the fakes produced by the Generator in line (H).

WGAN-GP for Learning a Point Distro (contd.)

- With regard to phase of Critic learning at the bottom of the previous slide, for the same reason as mentioned earlier, this output of the Critic is averaged over the batch in line (J) and subject to a call to `backward()` in line (K).
- For the third and final phase of Critic learning, using the implementation shown earlier, we first estimate the gradient penalty in line (L) and then make the call on `backward()` in line (M) for the final updating the gradients of the Critic weights in this iteration of the inner loop for Critic training.

```

for iteration in range(ITERS):                                     ## (A)
    ## Update the Critic network:
    for p in netC.parameters():                                  # reset the requires_grad attribute
        p.requires_grad = True                                  # this attribute is set to False in netG update

    for iter_d in range(CRITIC_ITERS):                           ## (B)
        ## The data_source supplies one BATCH_SIZE of 2D points from true distribution
        real_data = next(data_source)                            ## (C)
        real_data = torch.Tensor(real_data).requires_grad_(True).cuda()
        netC.zero_grad()
        # Train Critic network with real data:
        critic_for_reals = netC(real_data)                       ## (D)
        critic_for_reals_mean = critic_for_reals.mean()         ## (E)
        critic_for_reals_mean = torch.unsqueeze(critic_for_reals_mean, 0)
        critic_for_reals_mean.backward(minus_one)                ## (F)
        # Train Critic with Generator output:
        noise = torch.randn(BATCH_SIZE, 2).requires_grad_(True).cuda() ## (G)
        fakes = netG(noise)                                     ## (H)
        critic_for_fakes = netC(fakes)                           ## (I)
        critic_for_fakes_mean = critic_for_fakes.mean()         ## (J)
        critic_for_fakes_mean = torch.unsqueeze(critic_for_fakes_mean, 0)
        critic_for_fakes_mean.backward(one)                       ## (K)

```

WGAN-GP for Learning a Point Distro (contd.)

(..... continued from the previous slide)

```

if MODE == "wgan-gp":
    # Train Critic with gradient penalty:
    gradient_penalty = calc_gradient_penalty(netC, real_data, fakes.data)      ## (L)
    gradient_penalty.backward()                                               ## (M)
    lossCritic = critic_for_fakes_mean - critic_for_reals_mean + gradient_penalty ## (N)
elif MODE == "wgan":
    lossCritic = critic_for_fakes_mean - critic_for_reals_mean                ## (O)
wasser_dist = critic_for_reals_mean - critic_for_fakes_mean                  ## (P)
optim_critic.step()

## Update the Generator network:
for p in netC.parameters():                                                 ## (Q)
    p.requires_grad = False
netG.zero_grad()
noise = torch.randn(BATCH_SIZE, 2).requires_grad_(True).cuda()             ## (R)
## A BATCH_SIZE of fakes coming from the Generator
fakes = netG(noise)
critic_for_fakes_g = netC(fakes)                                           ## (S)
critic_for_fakes_g_mean = critic_for_fakes_g.mean()                        ## (T)
critic_for_fakes_g_mean = torch.unsqueeze(critic_for_fakes_g_mean,0)
critic_for_fakes_g_mean.backward(minus_one)                                 ## (U)
lossGen = -critic_for_fakes_g_mean
optim_gen.step()

```

- The code file that follows is the full implementation of the WGAN code. On the code shown, what remains unexplained is the implementation of the function `display_distributions()` that plays an important role in depicting the effectiveness of using gradient penalty for training a WGAN. That implementation will be explained later in this section.

Implementation of `wgan_for_point_distros.py`

```

## wgan_for_point_distros.py

import random
import matplotlib
import matplotlib.pyplot as plt
import numpy as np
import torch
import torch.autograd as autograd
import torch.nn as nn
import torch.optim as optim
import sys, os, glob, time

MODE = 'wgan-gp'          # Choose one of the two
#MODE = 'wgan'

DIM = 512                 # Dimensionality of nn.Linear layers
LAMBDA = .1              # For estimating the contribution of GP to overall loss
CRITIC_ITERS = 5         # How many critic iterations per generator iteration
BATCH_SIZE = 256        # Batch size
ITERS = 100000           # how many generator iterations to train for
#ITERS = 10000           # how many generator iterations to train for
dir_name_for_results = 'results' + "_" + MODE

# ===== Refresh directory for results =====
if os.path.exists(dir_name_for_results):
    files = glob.glob(dir_name_for_results + "/*")
    for file in files:
        if os.path.isfile(file):
            os.remove(file)
        else:
            files = glob.glob(file + "/*")
            list(map(lambda x: os.remove(x), files))
else:
    os.mkdir(dir_name_for_results)

since_beginning_dict = {'critic_loss': [], 'wasser_dist': [], 'gen_loss': []}
since_last_flush_dict = {'critic_loss': [], 'wasser_dist': [], 'gen_loss': []}

# ===== Class Definitions =====
class Generator(nn.Module):
    def __init__(self):
        super(Generator, self).__init__()
        main = nn.Sequential(
            nn.Linear(2, DIM),
            nn.ReLU(True),
            nn.Linear(DIM, DIM),
            nn.ReLU(True),
            nn.Linear(DIM, DIM),
            nn.ReLU(True),
            nn.Linear(DIM, 2),
        )
        self.main = main
    def forward(self, noise):
        output = self.main(noise)
        return output

class Critic(nn.Module):
    def __init__(self):
        super(Critic, self).__init__()
        main = nn.Sequential(
            nn.Linear(2, DIM),
            nn.ReLU(True),
            nn.Linear(DIM, DIM),
            nn.ReLU(True),
            nn.Linear(DIM, DIM),
            nn.ReLU(True),
            nn.Linear(DIM, 1),
        )
        self.main = main
    def forward(self, inputs):
        output = self.main(inputs)
        return output

```

wgan_for_point_distros.py (contd.)

(..... continued from the previous slide)

```

def weights_init(m):
    """
    This function is used to initialize the learnable weights in the Critic and
    the Generator networks
    """
    classname = m.__class__.__name__
    if classname.find('Linear') != -1:
        m.weight.data.normal_(0.0, 0.02)
        m.bias.data.fill_(0)

frame_index = [0]

# ===== Utility Functions =====
def display_distributions(real_data, netC, netG):
    """
    This very useful visualization function, written originally by the authors
    of the celebrated paper "Improved Training of Wasserstein GANs", does the
    following three things simultaneously:

    1) Creates a 128x128 array of points in a [-3,3]x[-3x3] box in the xy-plane.
    These points can subsequently be fed into the Critic for the values it
    would yield at each point in the array. The value returned by the Critic
    at point (x,y) in the array would be Critic's confidence whether that
    (x,y) point belongs to the probability distribution for the training data.
    The surface formed by such Critic values is best visualized through
    equi-valued contours.

    2) The 'real_data' that is the first argument to this function is a batch-full
    (typically 256) points in the xy-plane that were produced by the function
    multi_gaussian_source(). This source represents the true training data
    for training the GAN. These points are shown in the xy-plane by orange
    '+' points.

    3) It takes a batch-full (typically 256) 2D noise vectors and sends them
    through the Generator network netG. The Generator network produces a
    2D point in the xy-plane for each 2D noise input. The 256 points returned
    by the Generator are displayed as green 'x' points in the same xy-plane
    that is used for the above two items.
    """
    NPPOINTS = 128
    RANGE = 3
    POINTS = np.zeros((NPPOINTS, NPPOINTS, 2), dtype='float32')
    points[:, :, 0] = np.linspace(-RANGE, RANGE, NPPOINTS)[1:,None]
    points[:, :, 1] = np.linspace(-RANGE, RANGE, NPPOINTS)[None, 1:]
    points = points.reshape((-1, 2))
    points = torch.Tensor(points).requires_grad_(False).cuda()
    ## Generate the Critic's value at each point at each (x,y) point
    ## created above:
    critic_map = netC(points).cpu().data.numpy()
    ## Now we need a batch-full of 2D noise vectors for feeding into
    ## the Generator:
    noise = torch.randn(BATCH_SIZE, 2).requires_grad_(False).cuda()
    fakes = netG(noise).cpu().detach().numpy()
    plt.clf()
    x = y = np.linspace(-RANGE, RANGE, NPPOINTS)
    ## Display the Critic output value surface through contours:
    plt.contour(x, y, critic_map.reshape((1en(x), 1en(y))).transpose())
    ## Display the 256 first-arg real_data points that were previously
    ## generated by the ground-truth source multi_gaussian_source():
    plt.scatter(real_data[:,0], real_data[:,1], c='orange', marker='+')
    ## Now display the 256 'fake' points returned by the Generator:
    plt.scatter(fakes[:,0], fakes[:,1], c='green', marker='x')
    plt.savefig(dir_name_for_results + '/' + 'frame' + str(frame_index[0]) + '.jpg')
    frame_index[0] += 1

```

wgan_for_point_distros.py (contd.)

(..... continued from the previous slide)

```

def multi_gaussian_source():
    """
    A Python 'generator' function: Each call to this function with the built-in 'next()' will yield
    a fresh BATCH_SIZE (typically 256) number of points in the xy-plane.
    """
    scale = 2.
    centers = [
        (1, 0),
        (-1, 0),
        (0, 1),
        (0, -1),
        (1. / np.sqrt(2), 1. / np.sqrt(2)),
        (1. / np.sqrt(2), -1. / np.sqrt(2)),
        (-1. / np.sqrt(2), 1. / np.sqrt(2)),
        (-1. / np.sqrt(2), -1. / np.sqrt(2))
    ]
    centers = [(scale * x, scale * y) for x, y in centers]
    while True:
        dataset = []
        spread = 0.02
        spread = 0.1 # controls the spread of each Gaussian
        for i in range(BATCH_SIZE):
            point = np.random.randn(2) * spread
            center = random.choice(centers)
            point[0] += center[0]
            point[1] += center[1]
            dataset.append(point)
        dataset = np.array(dataset, dtype='float32')
        dataset /= 1.414 # stdev
        yield dataset

def calc_gradient_penalty(netC, real_data, fake_data):
    """
    Implementation by Marvin Cao at GitHub
    """
    epsilon = torch.rand(BATCH_SIZE, 1).cuda()
    epsilon = epsilon.expand(real_data.size())
    interpolates = epsilon * real_data + ((1 - epsilon) * fake_data)
    interpolates = interpolates.requires_grad_(True).cuda()
    critic_interpolates = netC(interpolates)
    gradients = autograd.grad(outputs=critic_interpolates, inputs=interpolates,
                              grad_outputs=torch.ones(critic_interpolates.size()).cuda(),
                              create_graph=True, retain_graph=True, only_inputs=True)[0]
    gradient_penalty = ((gradients.norm(2, dim=1) - 1) ** 2).mean() * LAMBDA
    return gradient_penalty

# ===== Create Network Instances and Train WGAN =====
netG = Generator().cuda()
netC = Critic().cuda()
netC.apply(weights_init)
netG.apply(weights_init)

# print(netG)
# print(netC)

optim_critic = optim.Adam(netC.parameters(), lr=1e-4, betas=(0.5, 0.9))
optim_gen = optim.Adam(netG.parameters(), lr=1e-4, betas=(0.5, 0.9))

one = torch.FloatTensor([1])
minus_one = one * -1

one = one.cuda()
minus_one = minus_one.cuda()

data_source = multi_gaussian_source() # returns one BATCH_SIZE collection 2D points
start_time = time.perf_counter()

```

wgan_for_point_distros.py (contd.)

(..... continued from the previous slide)

```

for iteration in range(ITERs):
    ## Update the Critic network
    for p in netC.parameters():
        p.requires_grad = True # reset the requires_grad attribute
                                # this attribute is set to False in netG update

    for iter_d in range(CRITIC_ITERs):
        ## The data_source supplies one BATCH_SIZE of 2D points from true distribution
        real_data = next(data_source)
        real_data = torch.Tensor(real_data).requires_grad_(True).cuda()
        netC.zero_grad()
        # train Critic network with real data
        critic_for_reals = netC(real_data)
        critic_for_reals_mean = critic_for_reals.mean()
        critic_for_reals_mean = torch.unsqueeze(critic_for_reals_mean, 0)
        critic_for_reals_mean.backward(minus_one)
        # train Critic with fakes
        noise = torch.randn(BATCH_SIZE, 2).requires_grad_(True).cuda()
        fakes = netG(noise)
        critic_for_fakes = netC(fakes)
        critic_for_fakes_mean = critic_for_fakes.mean()
        critic_for_fakes_mean = torch.unsqueeze(critic_for_fakes_mean, 0)
        critic_for_fakes_mean.backward(one)

    if MODE == "wgan-gp":
        # train with gradient penalty
        gradient_penalty = calc_gradient_penalty(netC, real_data, fakes.data)
        gradient_penalty.backward()
        lossCritic = critic_for_fakes_mean - critic_for_reals_mean + gradient_penalty
    elif MODE == "wgan":
        lossCritic = critic_for_fakes_mean - critic_for_reals_mean
        wasser_dist = critic_for_reals_mean - critic_for_fakes_mean
        optim_critic.step()

    ## Update the Generator network
    for p in netG.parameters():
        p.requires_grad = False
    netG.zero_grad()
    noise = torch.randn(BATCH_SIZE, 2).requires_grad_(True).cuda()
    ## A BATCH_SIZE of fakes coming from the Generator
    fakes = netG(noise)
    critic_for_fakes_g = netC(fakes)
    critic_for_fakes_g_mean = critic_for_fakes_g.mean()
    critic_for_fakes_g_mean = torch.unsqueeze(critic_for_fakes_g_mean, 0)
    critic_for_fakes_g_mean.backward(minus_one)
    lossGen = -critic_for_fakes_g_mean
    optim_gen.step()

    ## Update the dicts for the losses and distance:
    since_last_flush_dict['critic_loss'].append(lossCritic.cpu().data.numpy()[0])
    since_last_flush_dict['wasser_dist'].append(wasser_dist.cpu().data.numpy()[0])
    since_last_flush_dict['gen_loss'].append(lossGen.cpu().data.numpy()[0])

    if iteration % 100 == 99:
        current_time = time.perf_counter()
        elapsed_time = int(current_time - start_time)
        prints = []
        for name, vals in since_last_flush_dict.items():
            prints.append("{} {:.3f}".format(name, np.mean(vals)))
        since_beginning_dict[name] += vals
        print("\t[iter: {:5d}] time: {:5d} secs)\t\t{}".format(iteration+1, elapsed_time, "\t".join(prints) ))
        since_last_flush_dict = {'critic_loss': [], 'wasser_dist': [], 'gen_loss': []}
        real_data = next(data_source)
        display_distributions(real_data, netC, netG)
    iteration += 1

```

(Continued on the next slide

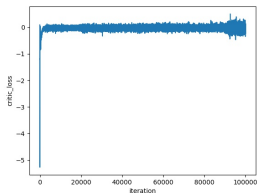
wgan_for_point_distros.py (contd.)

(..... continued from the previous slide)

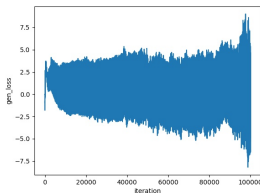
```
for name,vals in since_beginning_dict.items():
    x_vals = np.array(range(iteration))
    y_vals = [since_beginning_dict[name][x] for x in x_vals]
    plt.clf()
    plt.plot(x_vals, y_vals)
    plt.xlabel('iteration')
    plt.ylabel(name)
    plt.savefig(dir_name_for_results + "/" + name.replace(' ', '_')+'.jpg')
```

- Shown on the next slide are the Critic Loss, the Generator Loss, and the Wasserstein Distance calculated in the main training loop of the code.

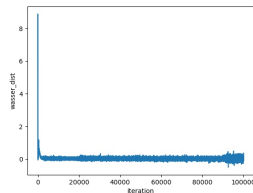
Losses vs. Iterations for WGAN-GP



Critic Loss



Generator Loss



Wasserstein Distance

Losses and distance based on 100,000 of training

Comparing GP with No-GP in Training a WGAN

- As to how effective using the Gradient Penalty is in improving the performance of a WGAN is best visualized by using the function `display_distributions()` whose implementation is presented next. The code for this function is as provided by the original authors of the paper "Improved Training of Wasserstein GANs". This function does the following three things simultaneously:
 - Creates a 128×128 array of points in a $[-3, 3] \times [-3, 3]$ box in the xy-plane. These points can subsequently be fed into the Critic for the values it would yield at each point. The value returned by the Critic at point (x,y) is Critic's confidence whether that point belongs to the probability distribution for the training data. The surface formed by such Critic values is best visualized through equi-valued contours.
 - The `real_data` that is the first argument to this function is a batch-full (typically 256) points in the xy-plane that were produced by the function `multi_gaussian_source()`. This source represents the ground-truth data for training the GAN. These points are shown in the xy-plane by orange '+' points.

Comparing GP with No-GP (contd.)

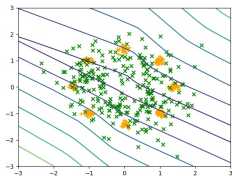
- It takes a batch-full (typically 256) 2D noise vectors and sends them through the Generator network netG. The Generator network produces a 2D point in the xy-plane for each 2D noise input. The 256 points returned by the Generator are displayed as green 'x' points in the same xy-plane that is used for the above two items.
- Shown below is the implementation of `display_distributions()`:

```
def display_distributions(real_data, netC, netG):
    NPOINTS = 128
    RANGE = 3
    points = np.zeros((NPOINTS, NPOINTS, 2), dtype='float32')
    points[:, :, 0] = np.linspace(-RANGE, RANGE, NPOINTS)[1:-1]
    points[:, :, 1] = np.linspace(-RANGE, RANGE, NPOINTS)[None, :]
    points = points.reshape((-1, 2))
    points = torch.Tensor(points).requires_grad_(False).cuda()
    ## Generate the Critic's value at each point at each (x,y) point
    ## created above:
    critic_map = netC(points).cpu().data.numpy()
    ## Now we need a batch-full of 2D noise vectors for feeding into
    ## the Generator:
    noise = torch.randn(BATCH_SIZE, 2).requires_grad_(False).cuda()
    fakes = netG(noise).cpu().detach().numpy()
    plt.clf()
    x = y = np.linspace(-RANGE, RANGE, NPOINTS)
    ## Display the Critic output value surface through contours:
    plt.contour(x, y, critic_map.reshape((len(x), len(y))).transpose())
    ## Display the 256 first-arg real_data points that were previously
    ## generated by the ground-truth source multi_gaussian_source():
    plt.scatter(real_data[:, 0], real_data[:, 1], c='orange', marker='+')
    ## Now display the 256 'fake' points returned by the Generator:
    plt.scatter(fakes[:, 0], fakes[:, 1], c='green', marker='x')
    plt.savefig(dir_name_for_results + "/" + 'frame' + str(frame_index[0]) + '.jpg')
    frame_index[0] += 1
```

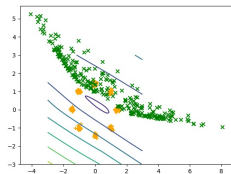
Comparing GP with No-GP (contd.)

- Shown in the next few slides is a side-by-side comparison of GP vs. no-GP on WGAN training at the same iteration index. The plots were produced by the function `display_distributions()` during a training session that consisted of 100,000 iterations.
- As mentioned earlier, the orange '+' marks denote the points in the xy-plane as produced by the true 8-Gaussian distribution and the green 'x' marks denote the points that the Generator produced from purely noise input. The greater the overlap between the green 'x' points and the orange '+' points, the superior the performance of the Generator. In addition, you would want the clusters formed by the green 'x' points to be as tight as those formed by the orange '+' points. Finally, you would want all the green 'x' points to fall inside the $[-3, 3] \times [-3, 3]$ box in the xy-plane.
- The contours depict the value surface for the Critic.

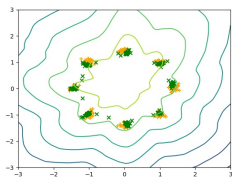
GP vs. No-GP Performance Comparison



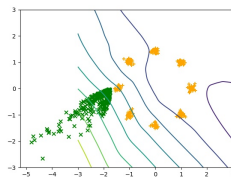
(a) With GP at iteration 10,000



(b) Without GP at iteration 10,000

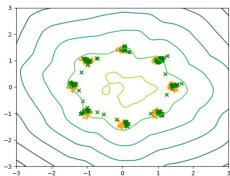


(a) With GP at iteration 20,000

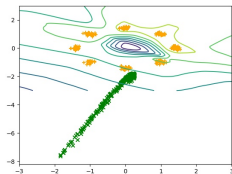


(b) Without GP at iteration 20,000

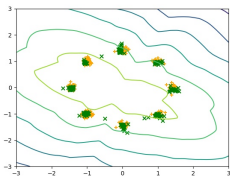
GP vs. No-GP Performance Comparison



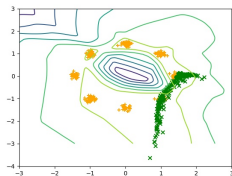
(a) With GP at iteration 30,000



(b) Without GP at iteration 30,000

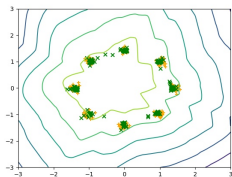


(a) With GP at iteration 40,000

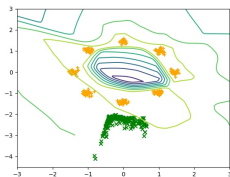


(b) Without GP at iteration 40,000

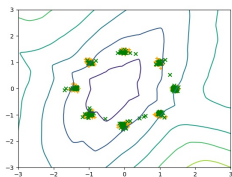
GP vs. No-GP Performance Comparison



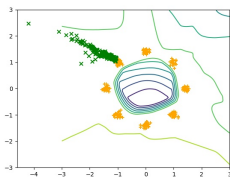
(a) With GP at iteration 50,000



(b) Without GP at iteration 50,000

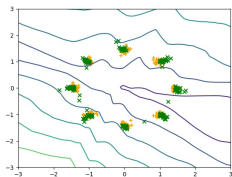


(a) With GP at iteration 60,000

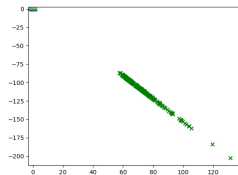


(b) Without GP at iteration 60,000

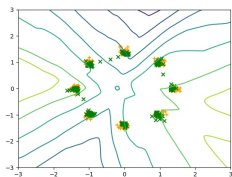
GP vs. No-GP Performance Comparison



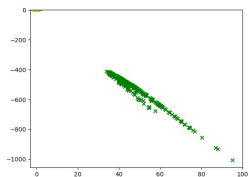
(a) With GP at iteration 70,000



(b) Without GP at iteration 70,000



(a) With GP at iteration 80,000



(b) Without GP at iteration 80,000