

Deep Learning and Applications

DSA 5204 • Lecture 10
Dr Low Yi Rui (Aaron)
Department of Mathematics



Last Time

We introduced an important class of generative models known as *variational autoencoders*

Important ideas:

- Latent variables
- Reparameterization trick

$$x \sim p_{\theta}(x) \rightarrow x = g(u, \theta), u \sim p_0(u)$$

Today, we will discuss some other generative models that also shares some of these ideas



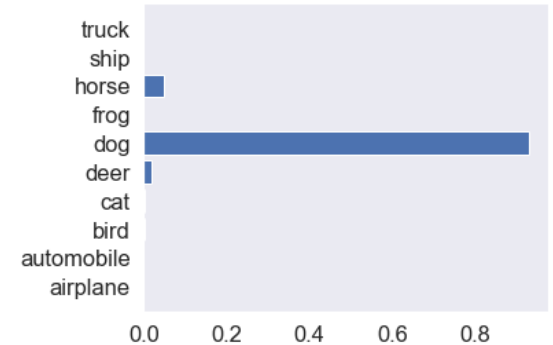
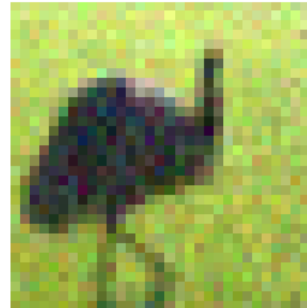
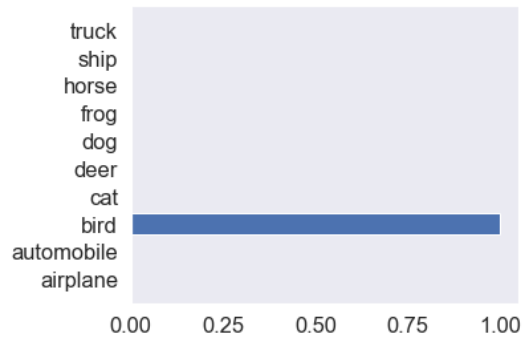
Generative Adversarial Networks

Review: Adversarial Training (Lecture 7)



Given a learned model $\hat{f}(x)$, we can fool the model by finding an *adversarial example*

$$x' = \operatorname{argmax}_{z, \|z-x\| \leq \delta} L(\hat{f}(z), y)$$



Review: Adversarial Training (Lecture 7)



To ameliorate this problem, adversarial training is proposed, in the form of a min-max optimization problem

Adversarial Training via FGSM (Fast Gradient Sign Method)

Hyper-parameters: δ (adv size), J (#adv train) ϵ_1, ϵ_2 (learning rates)

For $k = 0, 1, \dots$ do

$$\mathbf{z}_0 = \mathbf{x}$$

For $j = 0, 1, \dots, J - 1$ do

$$\mathbf{z}' = \mathbf{z}_j + \epsilon_2 \text{Sign} \left(\nabla_{\mathbf{z}} L(\hat{y}(\mathbf{z}_j; \boldsymbol{\theta}_k), y) \right)$$

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \epsilon_1 \nabla_{\boldsymbol{\theta}} L(\hat{y}(\mathbf{z}_J; \boldsymbol{\theta}_k), y)$$

The Role of Adversaries



The key idea of adversarial training is game-theoretic:

Use the adversary to challenge the training, so that if training wins, then we can beat adversaries.

Adversarial Training via FGSM (Fast Gradient Sign Method)

Hyper-parameters: δ (adv size), J (#adv train) ϵ_1, ϵ_2 (learning rates)

For $k = 0, 1, \dots$ do

$\mathbf{z}_0 = \mathbf{x}$

For $j = 0, 1, \dots, J - 1$ do

$$\mathbf{z}' = \mathbf{z}_j + \epsilon_2 \text{Sign} \left(\nabla_{\mathbf{z}} L(\hat{y}(\mathbf{z}_j; \boldsymbol{\theta}_k), y) \right)$$

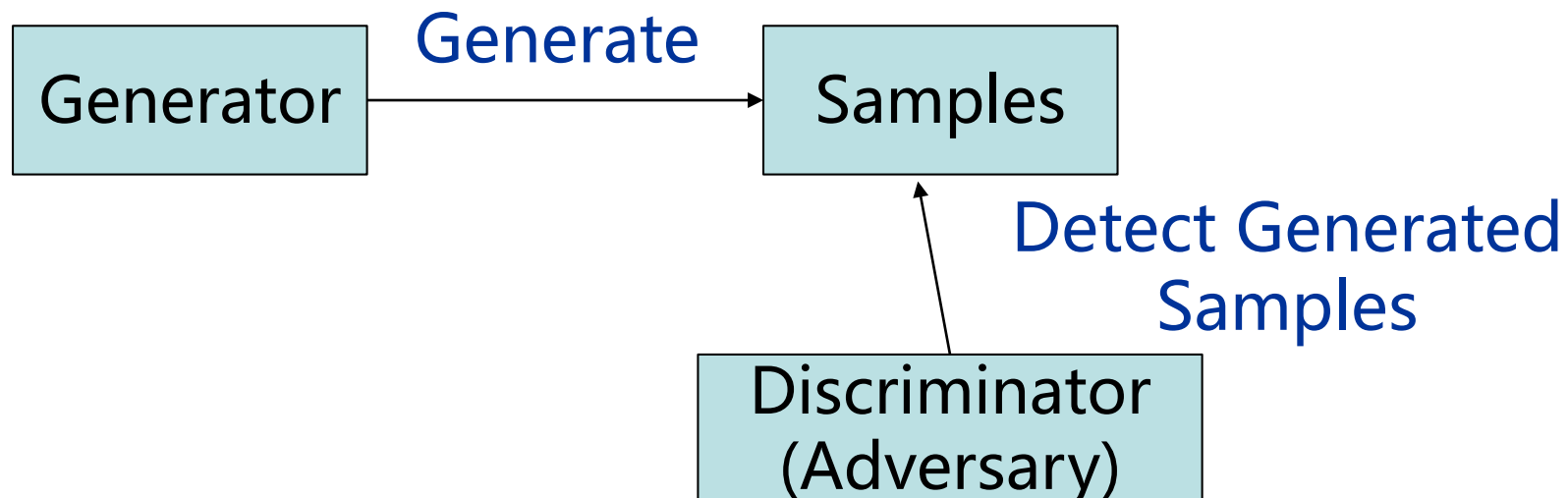
$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \epsilon_1 \nabla_{\boldsymbol{\theta}} L(\hat{y}(\mathbf{z}_J; \boldsymbol{\theta}_k), y)$$

Adversary

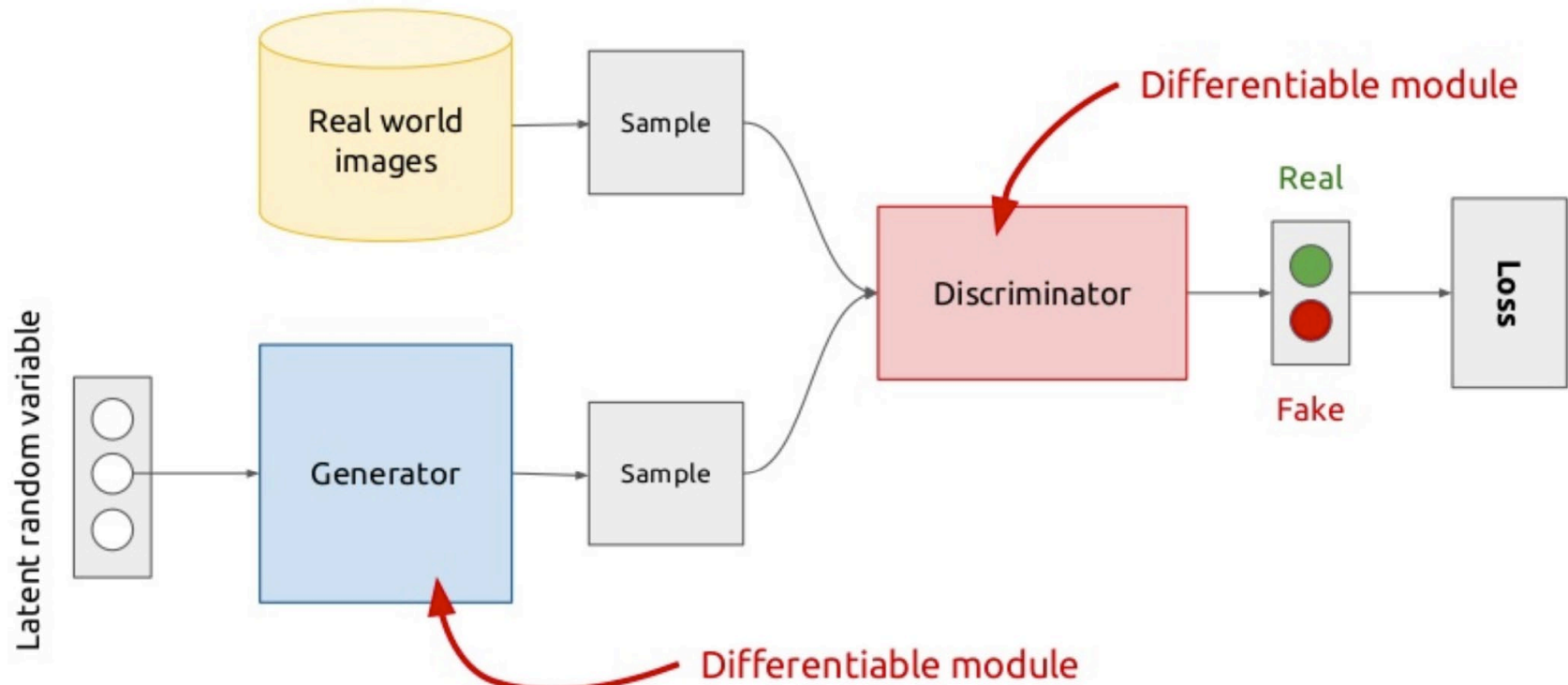
Adversaries for Generative Models



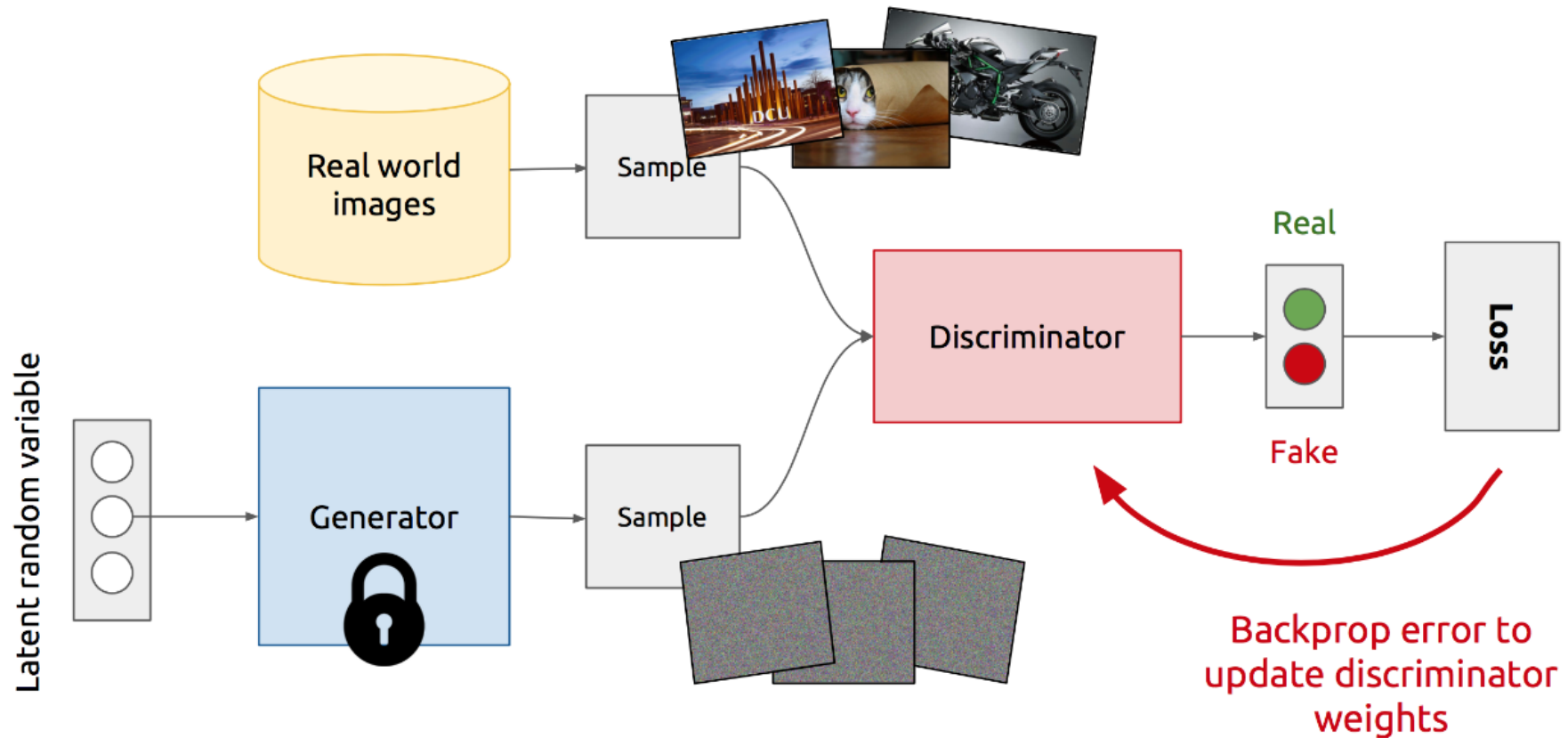
We make use of the previous idea of adversaries to generate outputs. This is known as *Generative Adversarial Networks (GAN)*



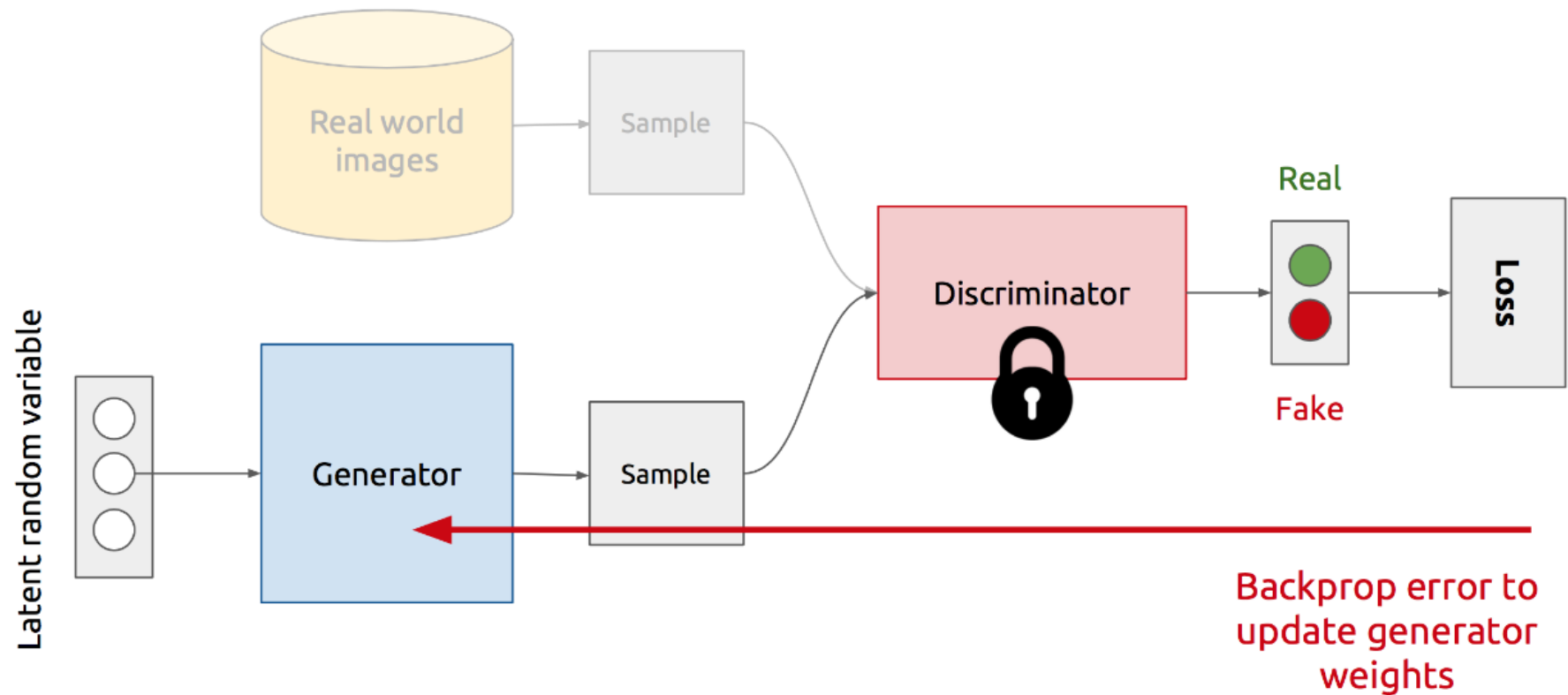
Detailed Architecture of GANs



Training Discriminator



Training Generator



Mathematical Formulation of The Training Problem

Latent variable (n -dimensional)

$$\mathbf{z} \sim p_0 = \mathcal{N}(0, I_n)$$

Training data

$$\{\mathbf{x}^{(i)} \in \mathbb{R}^d : i = 1, \dots, N, \mathbf{x}^{(i)} \sim p^*\}$$

Generator:

$$G_{\theta}: \mathbb{R}^n \rightarrow \mathbb{R}^d$$

Discriminator:

$$D_{\phi}: \mathbb{R}^d \rightarrow [0,1]$$

where $D_{\phi}(x)$ is the predicted probability of x being a real sample.

Mathematical Formulation of The Training Problem

Discriminator's cost function is to minimize the cross-entropy loss for real vs generated samples

Recall, given true binary labels $y \in \{0,1\}$, the (binary) cross-entropy loss is

$$-y \log D_{\phi}(x) - (1 - y) \log (1 - D_{\phi}(x))$$

- **If $y = 1$ (real sample), then the average loss is**

$$\mathbb{E}_{x \sim p^*} [-\log D_{\phi}(x)]$$

- **If $y = 0$ (generated sample), then the average loss is**

$$\mathbb{E}_{z \sim p_0} [-\log (1 - D_{\phi}(G_{\theta}(z)))]$$

Hence, the total (negative) loss the discriminator should maximize is

$$L(\theta, \phi) = \mathbb{E}_{x \sim p^*} [\log D_{\phi}(x)] + \mathbb{E}_{z \sim p_0} [\log (1 - D_{\phi}(G_{\theta}(z)))]$$

Min-Max Problem

Just like adversarial training, the adversary tries to maximize the loss, and the training then minimizes the maximized loss:

$$\min_{\theta} \max_{\phi} L(\theta, \phi)$$

Writing this out (and negating), we have

$$\min_{\theta} \max_{\phi} \mathbb{E}_{x \sim p^*} [\log D_{\phi}(x)] + \mathbb{E}_{z \sim p_0} [\log (1 - D_{\phi}(G_{\theta}(z)))]$$

This is called a

- *Min-max* (or minimax) formulation
- *Saddle-point* formulation
- *Zero-sum game* between D_{ϕ} and G_{θ}

Relationship with Game Theory



When interpreted as a two-player zero-sum game, we can ask where the *Nash equilibrium* is (or are).

A Nash equilibrium is where neither players have incentive to change its strategy.

A Nash equilibrium

- The distribution of $G_{\theta}(\mathbf{z})$ equals p^*
- $D_{\phi}(\mathbf{x}) = \frac{1}{2}$ for all \mathbf{x}

Relationship with KL Divergence

<Lecture Notebook>

Let us forget about the explicit parameterizations as write

$$\min_G \max_D \mathbb{E}_{\mathbf{x} \sim p^*} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_0} [\log(1 - D(G(\mathbf{z})))]$$

We can then show:

- Denote by p_G the distribution of $G(\mathbf{z})$. Then, for *fixed* G , the optimal D is

$$D_G^*(\mathbf{x}) = \frac{p^*(\mathbf{x})}{p^*(\mathbf{x}) + p_G(\mathbf{x})}$$

- Thus, the minimax problem can be framed as

$$\min_G D_{KL} \left(p^* \parallel \frac{p^* + p_G}{2} \right) + D_{KL} \left(p_G \parallel \frac{p^* + p_G}{2} \right) - 2 \log 2$$

KL Divergence and JS Divergence

Recall the KL divergence

$$D_{KL}(p\|q) = \int p(x) \log \left(\frac{p(x)}{q(x)} \right) dx$$

The Jensen-Shannon (JS) divergence is defined as

$$D_{JS}(p\|q) = \frac{1}{2} D_{KL} \left(p \parallel \frac{p+q}{2} \right) + \frac{1}{2} D_{KL} \left(q \parallel \frac{p+q}{2} \right)$$

Properties

- Like KL, JS divergence is non-negative and equals 0 if and only if $p = q$
- Unlike KL, JS divergence is *symmetric*
- In fact, $\sqrt{D_{JS}}$ is a metric on the space of distributions

Training GAN from the Divergence Point of View

The GAN training problem, interpreted from the divergence point of view, can be stated as

$$\min_G D_{JS}(p_G \| p^*) - \cancel{2 \log 2}$$

That is, we are minimizing the JS divergence (think “distance”) between the generation distribution p_G and the true data distribution p^ .*

The Basic Training Algorithm for GAN

Repeat:

Train Discriminator

- For k steps do
 - Sample m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\} \sim p_0$
 - Sample m real data samples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\} \in \mathcal{D}$
 - Update

$$\boldsymbol{\phi} \leftarrow \boldsymbol{\phi} + \epsilon_1 \nabla_{\boldsymbol{\phi}} \frac{1}{m} \sum_{i=1}^m \left[\log D_{\boldsymbol{\phi}}(\mathbf{x}^{(i)}) + \log \left(1 - D_{\boldsymbol{\phi}} \left(G_{\boldsymbol{\theta}}(\mathbf{z}^{(i)}) \right) \right) \right]$$

- Sample m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\} \sim p_0$ and update

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \epsilon_2 \nabla_{\boldsymbol{\theta}} \frac{1}{m} \sum_{i=1}^m \log \left(1 - D_{\boldsymbol{\phi}} \left(G_{\boldsymbol{\theta}}(\mathbf{z}^{(i)}) \right) \right)$$

Train Generator

Problem I: Vanishing Gradients

Loss function

$$L(\theta, \phi) = \mathbb{E}_{x \sim p^*} [\log D_\phi(x)] + \mathbb{E}_{z \sim p_0} [\log (1 - D_\phi(G_\theta(z)))]$$

The gradient with respect to θ controls the generator updates

$$\nabla_\theta L(\theta, \phi) = \nabla_\theta \mathbb{E}_{z \sim p_0} [\log (1 - D_\phi(G_\theta(z)))]$$

Observe if $\sigma = \text{Sigmoid}$, then

- Since $\nabla_a \log(1 - \sigma(a)) = -\sigma(a)$
- Then, $\nabla_\theta L(\theta, \phi) \propto \mathbb{E}_{z \sim p_0} D_\phi(G_\theta(z))$

Problem I: Vanishing Gradients

We have shown that

$$\nabla_{\theta} L(\theta, \phi) \propto \mathbb{E}_{\mathbf{z} \sim p_0} D_{\phi}(G_{\theta}(\mathbf{z})) \approx 0 \text{ if } D_{\phi} \approx 0$$

This means that

- Early in training, G_{θ} generates very bad samples
- Then, D_{ϕ} rejects them easily with high confidence, i.e. $D_{\phi}(G_{\theta}(\mathbf{z})) \approx 0$
- This causes $\nabla_{\theta} L$ to vanish, preventing further learning of θ

A Heuristic to Solve the Gradient Vanishing Problem

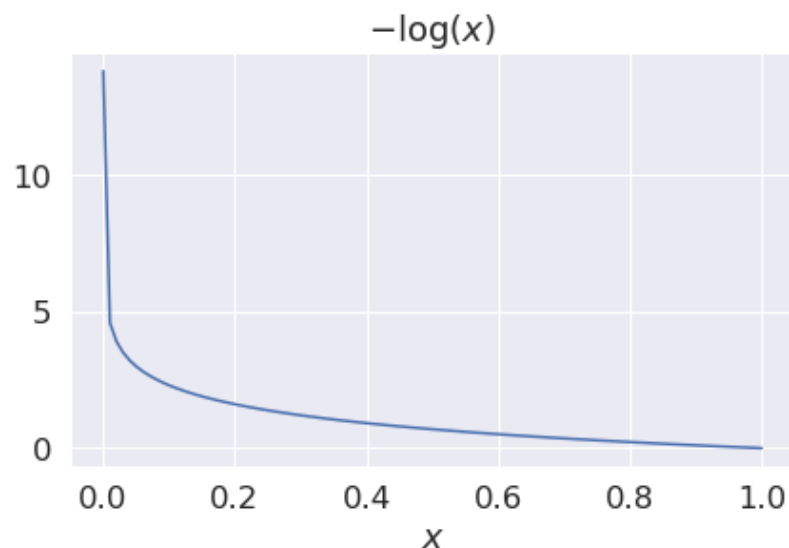
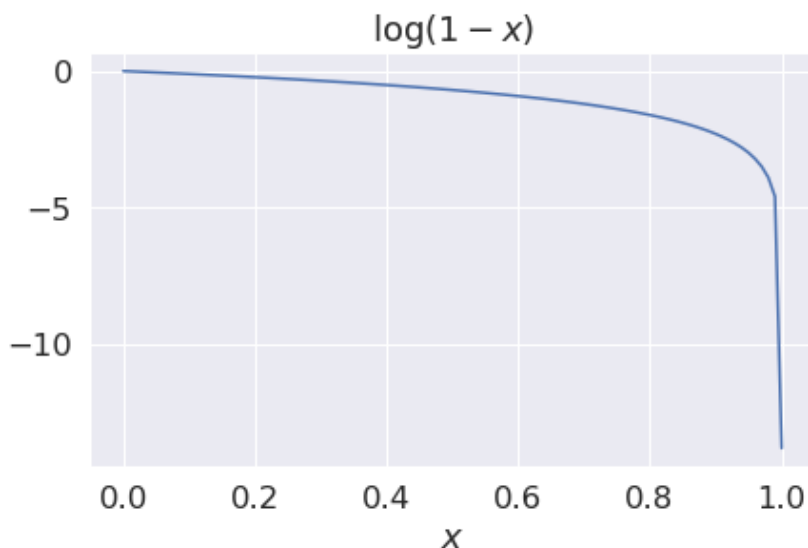


To solve the gradient vanishing problem for θ , we can replace

$$\nabla_{\theta} \mathbb{E}_{\mathbf{z} \sim p_0} \left[\log \left(1 - D_{\phi}(G_{\theta}(\mathbf{z})) \right) \right]$$

by

$$\nabla_{\theta} \mathbb{E}_{\mathbf{z} \sim p_0} \left[-\log \left(D_{\phi}(G_{\theta}(\mathbf{z})) \right) \right]$$



Problem II: Non-convergence

<Lecture Notebook>

Besides the gradient vanishing problem, in general GD is not designed to solve minimax problems.

Consider solving

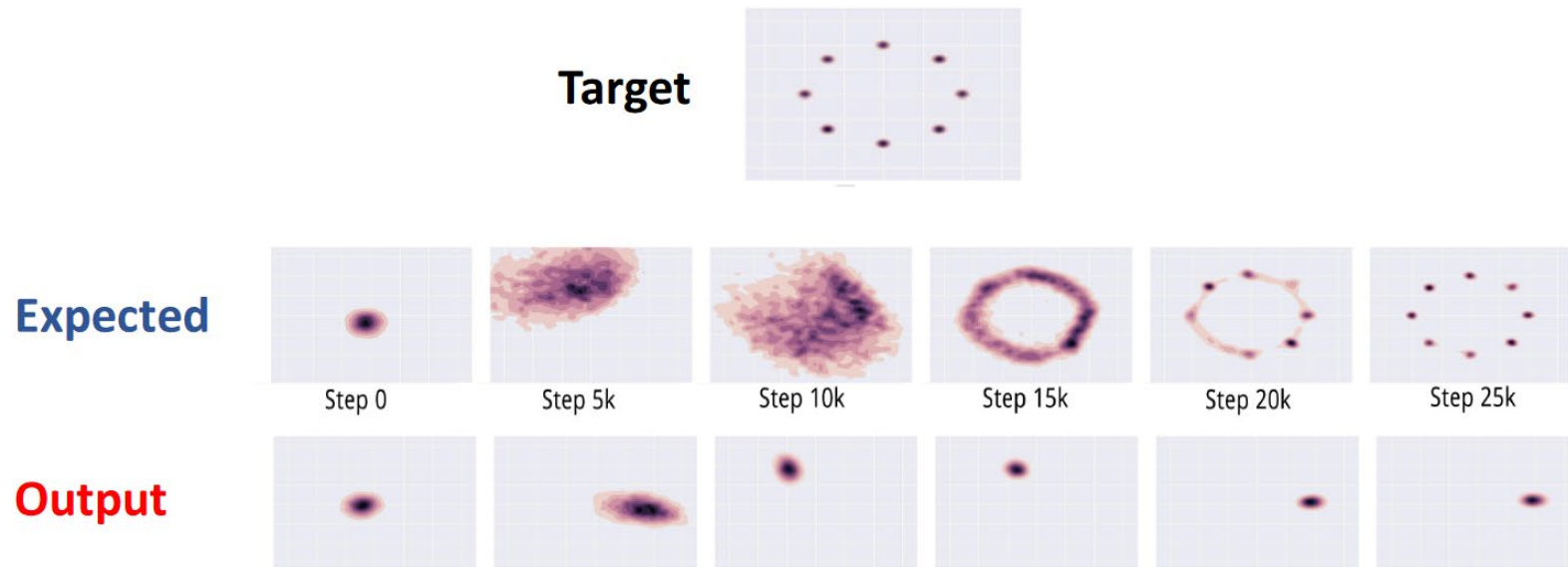
$$\min_x \max_y L(x, y) = xy$$

- The solutions are $x = 0, y \in \mathbb{R}$
- However, GD by alternating inner and outer loop will not find this solution!

Problem III: Mode Collapse



Mode collapse refers to the phenomenon that the generated samples are not varied and becomes very similar to a training sample.



Using Mini-Batch Statistics



Mode collapse can be detected, since generated samples are of small variance

To combat this, we can use *Mini-batch GAN*

- Generation happens in mini-batches, with batch statistics (e.g. variance, pair-wise L^2 norms) fed into discriminator
- Discriminator can tag an entire batch of examples as false if the variation is low
- Hence, generator is forced to produce varied samples



Wasserstein GANs

The Distance Viewpoint

Let us return to the viewpoint of training GANs as *minimizing “distances”* between distributions.

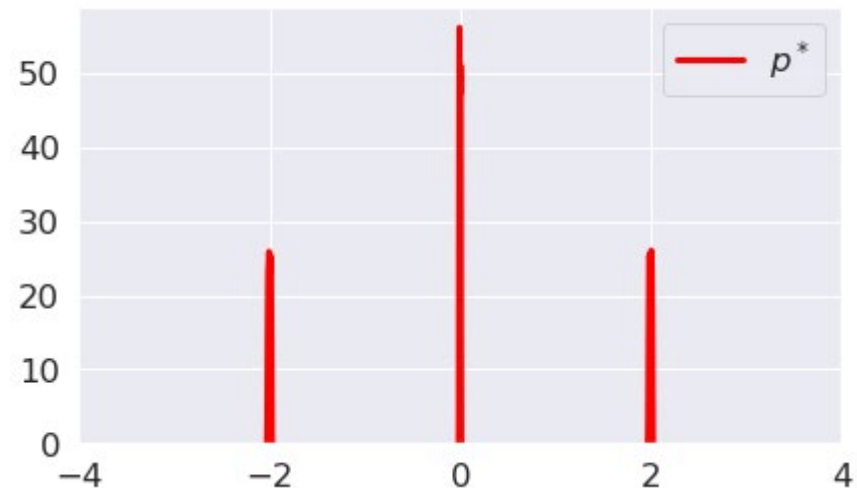
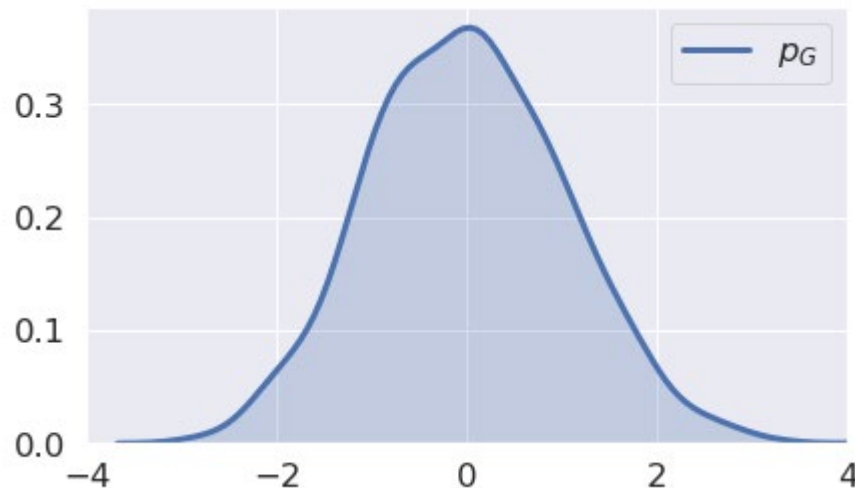
Recall that training GAN can be formulated as

$$\begin{aligned} & \min_G D_{JS}(p_G \| p^*) \\ &= \min_G \frac{1}{2} D_{KL} \left(p_G \| \frac{p^* + p_G}{2} \right) + \frac{1}{2} D_{KL} \left(p^* \| \frac{p^* + p_G}{2} \right) \end{aligned}$$

What's wrong with KL/JS divergences?



KL/JS divergences don't play nice when the distributions are *singular*



Example: Singular KL Divergence

Let $p(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2}$ **and** $q(x) = \frac{1}{2} \mathbb{I}_{x=1} + \frac{1}{2} \mathbb{I}_{x=-1}$

Then,

$$D_{KL}(p\|q) = \int p(x) \log \left(\frac{p(x)}{q(x)} \right) dx$$

is undefined since for all $x \neq \pm 1$, $q(x) = 0$ **but** $p(x) \neq 0$.

In general, $D_{KL}(p\|q)$ **is finite if and only if** p **is** *absolutely continuous* **with respect to** q , **meaning that for any*** subset A ,

$$\int_A q(x) dx = 0 \Rightarrow \int_A p(x) dx = 0$$

Otherwise, the D_{KL} **is infinite, and we say that** p **is** *singular* **with respect to** q .

Continuity and Distance



As in VAE and GAN, our primary idea of representing a distribution is to

- Sample latent $\mathbf{z} \sim p_0$ (simple distribution)
- Generate samples $G(\mathbf{z})$ with desired distribution p_G by adjusting G
- With parameterization: $G(\mathbf{z}) = G(\mathbf{z}; \boldsymbol{\theta})$ so $p_G = p_{\boldsymbol{\theta}}$

Then, the GAN training problem is

$$\min_{\boldsymbol{\theta}} \text{Distance}(p_{\boldsymbol{\theta}}, p^*)$$

In the usual GAN, the Distance function is the JS-divergence

Continuity and Distance



It is our hope to minimize $\text{Distance}(p_\theta, p^*)$ by gradient descent

$$\theta \leftarrow \theta - \epsilon \nabla_\theta \text{Distance}(p_\theta, p^*)$$

For this to work, it is essential that the following *continuity property* be satisfied

If we change θ by a little bit, then $\text{Distance}(p_\theta, p^)$ changes also by only a little bit*

It turns out, that the JS-divergence (or KL-divergence) does not satisfy this property!

Example: Learning Parallel Lines

<Lecture Notebook>



Let $z \sim \mathcal{N}(0,1)$ and consider $x \in \mathbb{R}^2$

- Let p^* be the distribution of $(0, z)$
- Let p_θ be the distribution of (θ, z)

Then, we can show

- $D_{KL}(p_\theta \| p^*) = \mathbb{I}_{\theta \neq 0}(+\infty)$
- $D_{JS}(p_\theta \| p^*) = \mathbb{I}_{\theta \neq 0}(\log 2)$

Both are not continuous at $\theta = 0$!

Modifying the GAN Problem

The natural question, is then, can we replace JS-divergence by another measure of distance so that the discontinuity issue is satisfied?

The answer is yes!

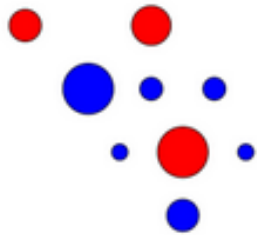
Consider the *Wasserstein distance*, or the *earth-mover's distance* (EMD)

$$W(p, q) = \inf_{\gamma \in \Pi(p, q)} \mathbb{E}_{(x, y) \sim \gamma} [\|x - y\|]$$

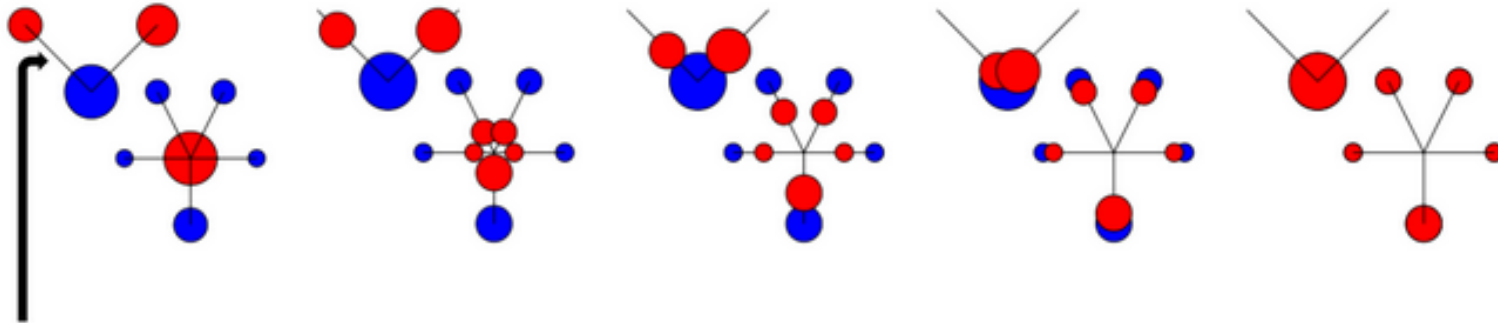
where $\Pi(p, q)$ is the set of joint distributions whose marginals are p and q respectively, i.e. $\gamma \in \Pi(p, q)$ if and only if

$$\int \gamma(x, y) dy = p(x) \quad \text{and} \quad \int \gamma(x, y) dx = q(y)$$

EMD Visualized



- red distribution: "dirt"
- blue distribution: "holes"



The distance between points (ground distance) can be Euclidean distance, Manhattan...

Wasserstein GAN

The Wasserstein GAN transforms the problem

$$\min_{\theta} D_{JS}(p_{\theta} \| p^*) \rightarrow \min_{\theta} W(p_{\theta}, p^*)$$

Advantages

- Unlike D_{JS} (or D_{KL}), W is a true distance, i.e. a metric
- The mapping $\theta \mapsto W(p_{\theta}, p^*)$ is continuous (for the line problem, we have $W(p_{\theta}, p^*) = |\theta|$)
- The mapping $\theta \mapsto W(p_{\theta}, p^*)$ is differentiable almost everywhere if $\theta \mapsto p_{\theta}$ is Lipschitz

How to Compute W ?

We need to compute

$$W(p_\theta, p^*) = \inf_{\gamma \in \Pi(p_\theta, p^*)} \mathbb{E}_{(x,y) \sim \gamma} [\|x - y\|]$$

and possibly also $\nabla_\theta W(p_\theta, p^*)$ for GD.

Direct computation via optimization is *not* tractable.
However, we have a beautiful duality theory for W !

Kantorovich-Rubinstein Duality:

$$W(p_\theta, p^*) = \sup_{f: \|f\|_L \leq 1} \mathbb{E}_{x \sim p_\theta} [f(x)] - \mathbb{E}_{x \sim p^*} [f(x)]$$

Here, $\|f\|_L$ denotes the Lipschitz constant for f (Think maximum value of $\|\nabla f\|$)

Wasserstein GAN

Owing to Kantorovich duality, we can approximate

$$W(p_\theta, p^*) \approx \max_{\phi} \mathbb{E}_{x \sim p_\theta} [f(x; \phi)] - \mathbb{E}_{x \sim p^*} [f(x; \phi)]$$

with $f(x, \phi)$ a neural network (ϕ are the weights).

In particular, if $\hat{\phi}$ attains this maximum then the gradient is

$$\begin{aligned} \nabla_{\theta} W &= \nabla_{\theta} \mathbb{E}_{x \sim p_\theta} [f(x; \hat{\phi})] \\ &= \mathbb{E}_{z \sim p_0} \nabla_{\theta} f(G(z; \theta); \hat{\phi}) \end{aligned}$$

The last expectation can be computed using Monte-Carlo approximation.

Wasserstein GAN Algorithm

Train Dual

Repeat:

- For k steps do
 - Sample m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\} \sim p_0$
 - Sample m real data samples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\} \in \mathcal{D}$
 - Update

$$\boldsymbol{\phi} \leftarrow \boldsymbol{\phi} + \epsilon_1 \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\phi}} [f(\mathbf{x}^{(i)}; \boldsymbol{\phi}) - f(G_{\boldsymbol{\theta}}(\mathbf{z}^{(i)}); \boldsymbol{\phi})]$$

- Sample m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\} \sim p_0$ and update

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \epsilon_2 \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} f(G_{\boldsymbol{\theta}}(\mathbf{z}^{(i)}); \boldsymbol{\phi})$$

Train Generator

Some Remarks on GANs



There are many variants of GAN architectures and training methods, and we have only touch on a few
Training GANs is not a solved problem!

Some resources:

- <https://arxiv.org/pdf/2001.06937.pdf>
- <https://arxiv.org/abs/1906.01529.pdf>



Demo: GAN for Image Generation



Neural Style Transfer

Transferring Styles



The style transfer task

- Given a content image x_c
- And a style image x_s
- Generate image x_{cs} with content x_c and style x_s



Content Image

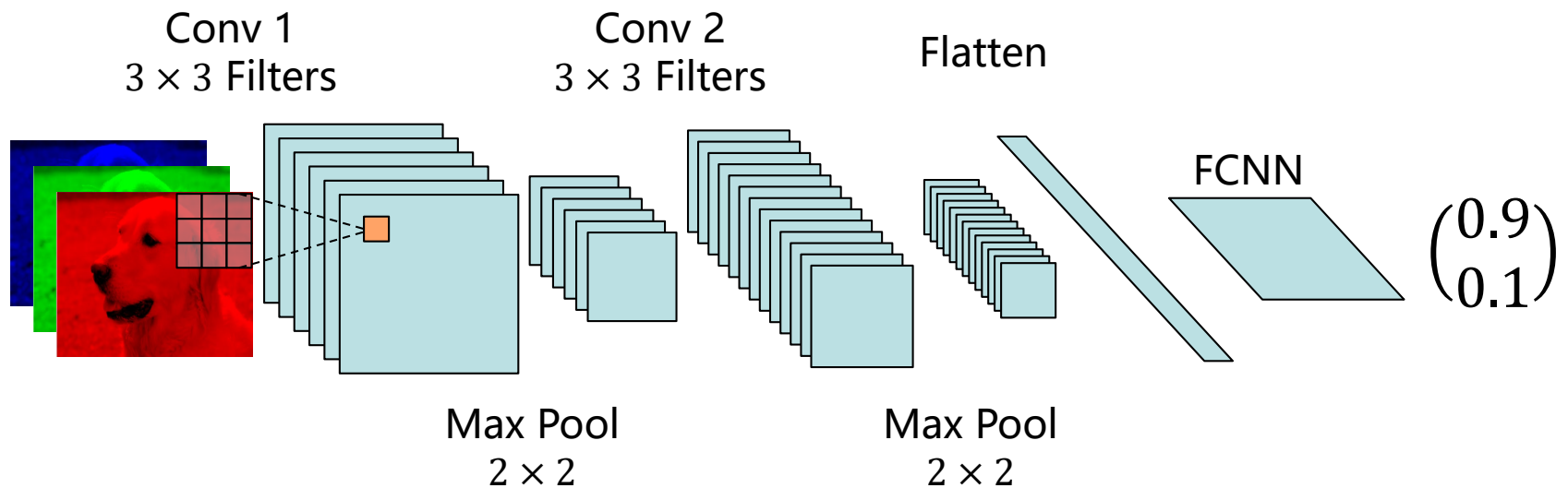
+



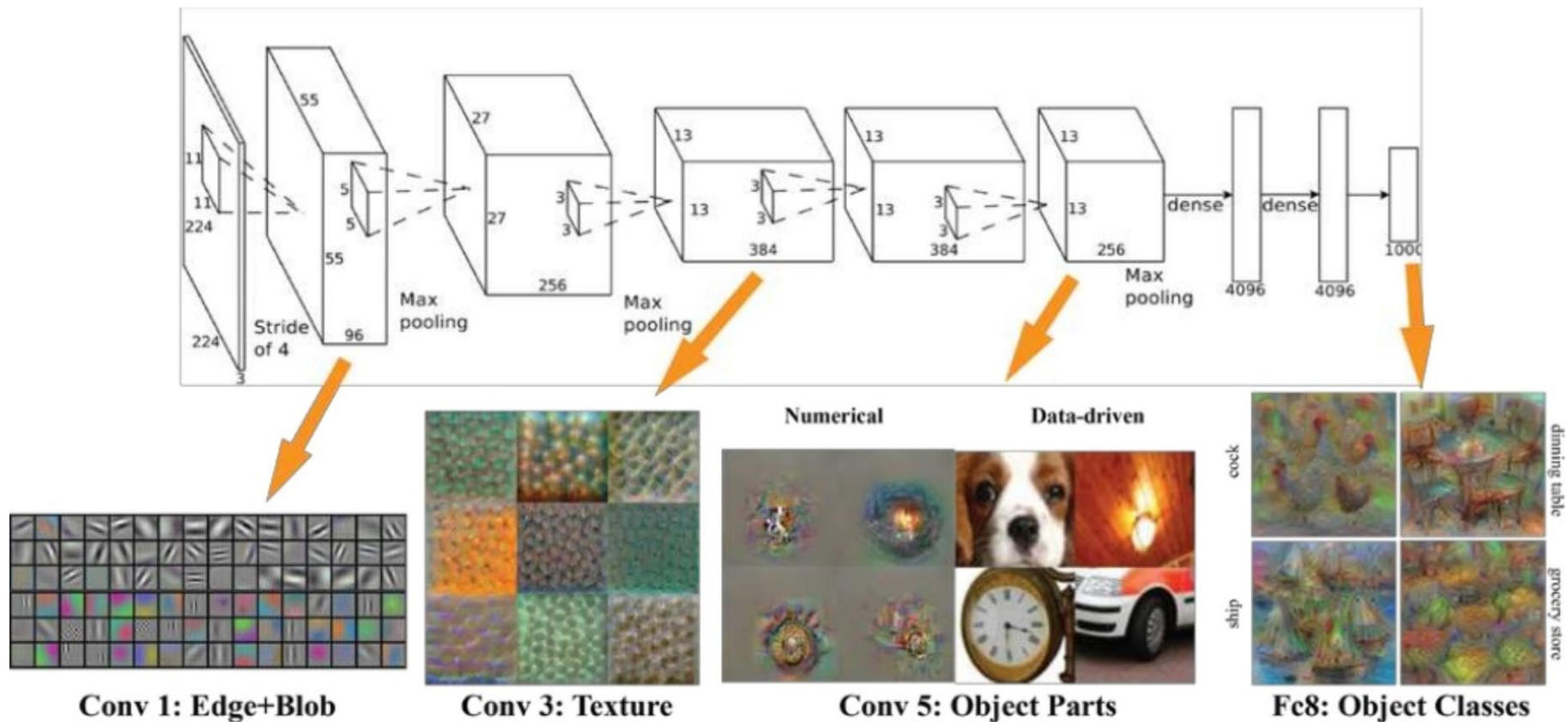
Style Image

Gatys et al, "A Neural Algorithm of Artistic Style", arXiv 2015

The Structure of CNNs



Features at Different Layers



Neural style transfer presentation done at Tensorflow summit Colombo 2018

A decorative network diagram in the top right corner, consisting of light blue nodes connected by thin lines, forming a complex web-like structure.

**The hierarchical
representation of
features in CNNs forms
the basis of style transfer**

Notation

We consider passing an input image x into a L -layer CNN.

The feature maps are

$$h^{(\ell+1)} = f^{(\ell)}(h^{(\ell)}; \theta^{(\ell)}), \quad h^{(0)} = x$$

The function $f^{(\ell)}$ is the ℓ^{th} convolutional block (e.g. Conv-BN-Activation-Pool).

The hidden feature maps are $\{h^{(\ell)}\}$. We denote

$$h_{ij}^{(\ell)} = j^{\text{th}} \text{ pixel of } i^{\text{th}} \text{ channel feature at } \ell^{\text{th}} \text{ layer}$$

To emphasize dependence on the initial input image x , we can write

$$h_{ij}^{(\ell)} = h_{ij}^{(\ell)}(x)$$

Combining Content and Style Losses

Our goal is to generate an image x with content x_c and style x_s

The idea is to perform minimization

$$\min_x \alpha L_c(x, x_c) + \beta L_s(x, x_s)$$

Here,

- L_c is the *content loss* that measures similarity in content
- L_s is the *style loss* that measures similarity in style

Content Loss

For the content loss we can take the square distance between the respective feature maps

$$E_c(\mathbf{x}, \mathbf{x}_c, \ell) = \frac{1}{2} \sum_{i,j} \left(h_{ij}^{(\ell)}(\mathbf{x}) - h_{ij}^{(\ell)}(\mathbf{x}_c) \right)^2$$

The overall content loss is a weighted sum

$$L_c(\mathbf{x}, \mathbf{x}_c) = \sum_{\ell \in \ell_c} E_c(\mathbf{x}, \mathbf{x}_c, \ell)$$

where ℓ_c is the set of layers relevant to content

Style Loss

To define the style loss, we first define the *Gram matrix*

Given a dataset matrix X (each row is a data vector), the Gram matrix is given by

$$G = XX^T$$

which has the dimensions $N \times N$

In fact, $G_{ij} = [x^{(i)}]^T x^{(j)}$, hence G measures the *correlations or similarities* amongst different data vectors, and is a summary of the statistics of our data

Style Loss

Therefore, we can define the layer-wise style loss

$$E_S(\mathbf{x}, \mathbf{x}_s, \ell) = \frac{1}{4N_l M_l} \sum_{i,j} \left(G_{ij}^{(\ell)}(\mathbf{x}) - G_{ij}^{(\ell)}(\mathbf{x}_s) \right)^2$$

where N_l is the number of channels, M_l is the image size and $G^{(\ell)}$ is the Gram matrix of the features

$$G_{ij}^{(\ell)}(\mathbf{x}) = \sum_k h_{ik}^{(\ell)}(\mathbf{x}) h_{jk}^{(\ell)}(\mathbf{x})$$

The total style loss is

$$L_S(\mathbf{x}, \mathbf{x}_s) = \sum_{\ell \in \ell_s} E_S(\mathbf{x}, \mathbf{x}_s; \ell)$$

Further Reading on NST

- <https://arxiv.org/pdf/1705.04058.pdf> (Review)
- <https://arxiv.org/pdf/1508.06576.pdf> (Original)



Demo: Neural Style Transfer

Summary



In this class, we introduced two other types of generative models

- Generative Adversarial Networks
- Neural Style Transfer

They represent a combination of ideas we have introduced previously

- Training with adversaries
- Generating samples using latent variables
- Using hierarchical representations in CNNs

Summary of DSA5204

A decorative network graph in the top right corner, consisting of numerous light blue circular nodes connected by thin, light blue lines, forming a complex web-like structure.

This class is an overall introduction of deep learning

- Architectures (FCNNs, CNNs, RNNs)
- Optimization methods (SGD and variants)
- Regularization methods (batch-norm, data augmentation, dropout, etc.)

We have also seen a variety of applications

- Supervised learning (classification, regression)
- Unsupervised learning (compression, denoising)
- Generative modelling

Next Steps

A decorative network graph in the top right corner, consisting of numerous light blue circular nodes connected by thin, light blue lines, forming a complex web-like structure.

This class is just an introduction, and there are many topics that are not covered, e.g. Bayesian NNs, attention models and many more!

How to learn more about deep learning?

- Read papers (NIPS, ICLR, ICML, CVPR, JMLR, etc)
- Follow ArXiv preprints
- Tensorflow/pytorch demonstrations
- Do an application project that interests you!