

TTIC 31230, Fundamentals of Deep Learning

David McAllester, Autumn 2020

Generative Adversarial Networks (GANs)

Modeling Probability Distributions on Images

Suppose we want to train a model of the probability distribution of natural images using cross-entropy loss.

$$\Phi^* = \operatorname{argmin}_{\Phi} E_{y \sim \text{pop}} - \ln p_{\Phi}(y)$$

Images are continuous structured objects — a continuous value at every pixel.

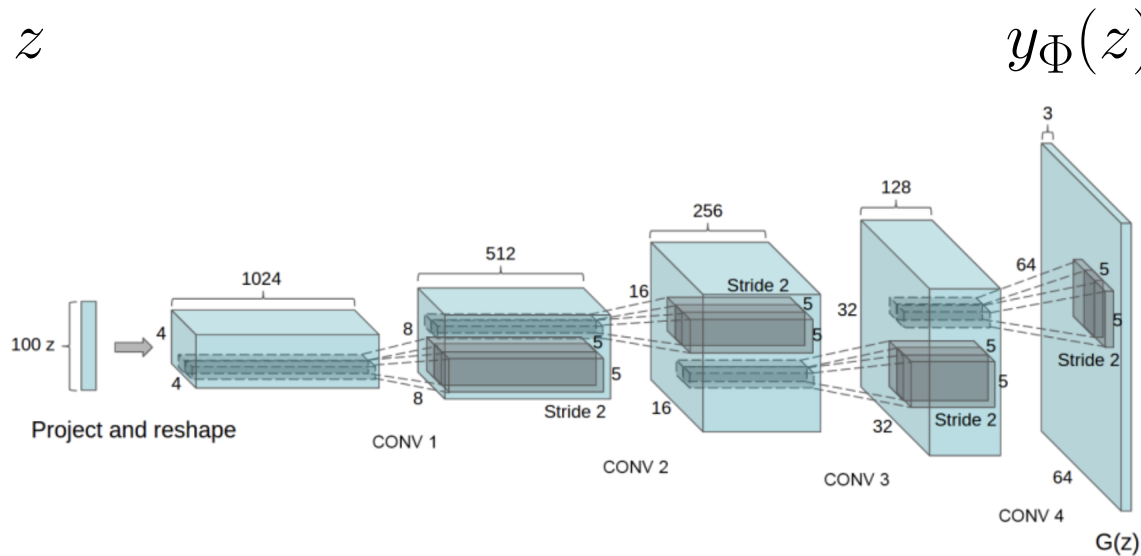
It is difficult to build probability models for images (or other continuous structured values) that both accurately model the distribution and also allow us to calculate $p_{\Phi}(y)$.

Generative Adversarial Networks (GANs)

GANs represent $p_{\Phi}(y)$ implicitly by constructing an image generator and abandon the ability to compute $p_{\Phi}(y)$.

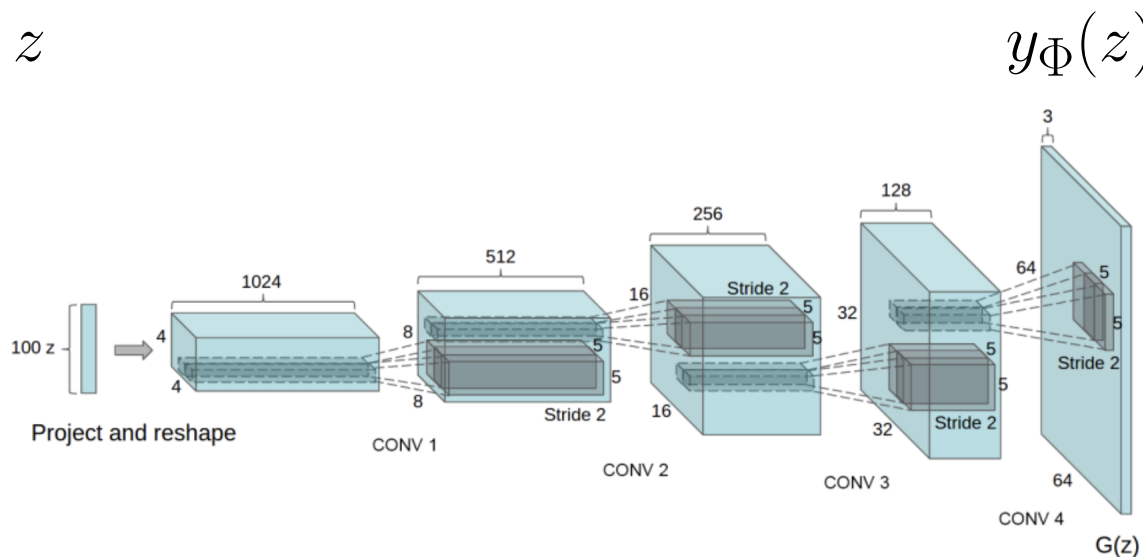
The cross-entropy loss is replaced by an adversarial discriminator which tries to distinguish between generated images and real images.

Representing a Distribution with a Generator



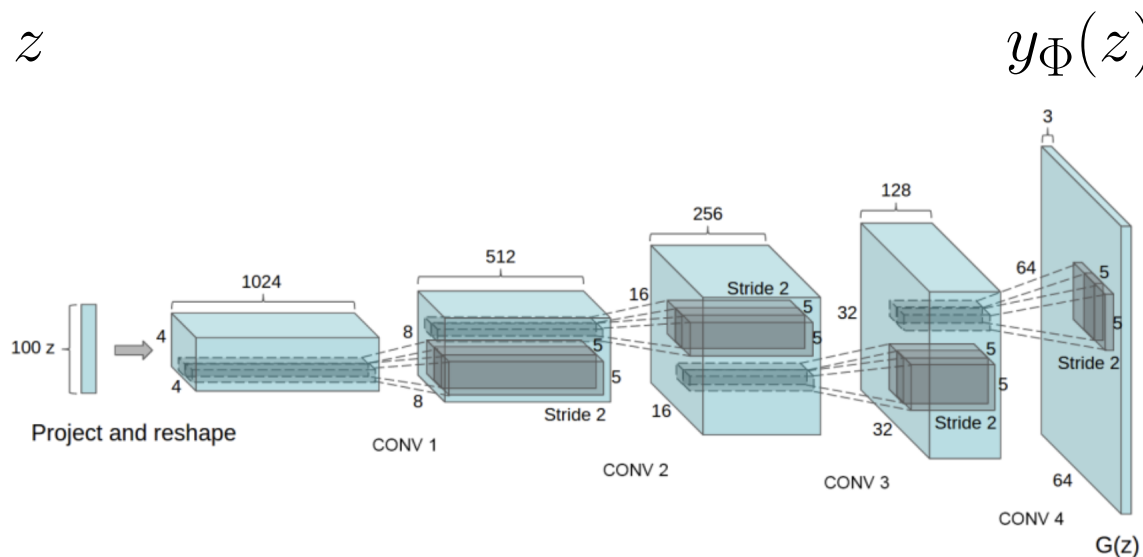
The random input z defines a probability density on images $y_{\Phi}(z)$. We will write this as $p_{\Phi}(y)$ for the image y .

Representing a Distribution with a Generator



We want $p_{\Phi}(y)$ to model a natural image distribution such as the distribution over human faces.

Representing a Distribution with a Generator



We can sample from $p_{\Phi}(y)$ by sampling z . But we cannot compute $p_{\Phi}(y)$ for y sampled from the population.

Increasing Spatial Dimension

Reducing spatial dimension with strided convolution:

For $x, y, j, \Delta x, \Delta y, i$

$$L_{\ell+1}[\textcolor{red}{x}, \textcolor{red}{y}, j] += W[\Delta x, \Delta y, i, j] L_{\ell}[\textcolor{red}{s} * \textcolor{red}{x} + \Delta x, \textcolor{red}{s} * \textcolor{red}{y} + \Delta y, i]$$

Increasing spatial dimension with PyTorch ConvTranspose2d:

For $x, y, j, \Delta x, \Delta y, i$

$$L_{\ell+1}[\textcolor{red}{s} * \textcolor{red}{x} + \Delta x, \textcolor{red}{s} * \textcolor{red}{y} + \Delta y, i] += W[\Delta x, \Delta y, i, j] L_{\ell}[\textcolor{red}{x}, \textcolor{red}{y}, j]$$

Generative Adversarial Networks (GANs)

Let y range over images. We have a generator p_Φ . For $i \in \{-1, 1\}$ we define a probability distribution over pairs $\langle i, y \rangle$ by

$$\begin{aligned}\tilde{p}_\Phi(i = 1) &= 1/2 \\ \tilde{p}_\Phi(y|i = 1) &= \text{pop}(y) \\ \tilde{p}_\Phi(y|i = -1) &= p_\Phi(y)\end{aligned}$$

We also have a discriminator $P_{\text{Disc}}(i|y)$ that tries to determine the source i given the image y .

The generator tries to fool the discriminator.

$$\text{Gen}^* = \underset{\text{Gen}}{\text{argmax}} \min_{\text{Disc}} E_{\langle i, y \rangle \sim \tilde{p}_{\text{Gen}}} - \ln P_{\text{Disc}}(i|y)$$

GANs

The generator tries to fool the discriminator.

$$\text{Gen}^* = \underset{\text{Gen}}{\operatorname{argmax}} \min_{\text{Disc}} E_{\langle i, y \rangle \sim \tilde{p}_{\text{Gen}}} - \ln P_{\text{Disc}}(i|y)$$

Assuming universality of both the generator p_{Gen} and the discriminator P_{Disc} we have $p_{\text{Gen}}^* = \text{pop}$.

Note that this involves only discrete cross-entropy.

GANs

To make the gradient descent clearer we write

$$E_{\langle i, y \rangle \sim \tilde{p}_{\text{Gen}}} - \ln P_{\text{Disc}}(i|y)$$

as

$$\frac{1}{2} E_{y \sim \text{pop}} - \ln P_{\text{Disc}}(1|y) \quad + \quad \frac{1}{2} E_{z \sim \mathcal{N}(0, I)} - \ln P_{\text{Disc}}(-1|y_{\text{Gen}}(z))$$

Generative Adversarial Nets

Goodfellow et al., June 2014



The rightmost column (yellow borders) gives the nearest neighbor in the training data to the adjacent column.

GAN Mode Collapse

A major concern is “mode collapse” where the learned distribution omits a significant fraction of the population distribution.

There is no quantitative performance measure that provides a meaningful guarantee against mode collapse.

The Fréchet Inception Score (FID)

The main problem with GANs is the lack of a meaningful quantitative evaluation metric.

A standard quantitative performance measure is Fréchet Inception Distance (FID).

This measures statistics of the features of the inception image classification model (trained on imagenet) for images generated by the generator.

It then compares those statistics to the same statistics for images drawn from the population.

But the FID score provides no guarantees against mode collapse.

END