Lecture 18: Image Synthesis

CMSC 472 / 672

Computer Vision





Lecture 18: Image Synthesis

CMSC 472 / 672

Computer Vision



SEARCH FOR THE BIGHT
TEMPLATE, MANUALLY
WRITE CAPTION, ALIGN
THE TEXT, PLACE THEM IN
THE BIGHT PLACE

TUBN TEXT
INTO MEMES

Supermemedi



Machine Learning Problems



SUPERVISED LEARNING

Training time

data:

$$\{\mathbf{x}^{(t)}, y^{(t)}\}$$

setting:

$$\mathbf{x}^{(t)}, y^{(t)} \leftarrow p(\mathbf{x}, y)$$

Test time

data:

$$\{\mathbf{x}^{(t)}, y^{(t)}\}$$

setting:

$$\mathbf{x}^{(t)}, y^{(t)} \leftarrow p(\mathbf{x}, y)$$

Example

Input: $x^{(t)}$ is an image

 $Output: y^{(t)}$ is an image

category

ONE-SHOT LEARNING

Training time

data:

$$\{\mathbf{x}^{(t)}, y^{(t)}\}$$

setting:

$$\mathbf{x}^{(t)}, y^{(t)} \leftarrow p(\mathbf{x}, y)$$

subject to $y^{(t)} \ 2 \ \{1, ..., C\}$

Test time

data:

$$\{\mathbf{x}^{(t)}, y^{(t)}\}$$

setting:

$$\mathbf{x}^{(t)}, y^{(t)} \leftarrow p(\mathbf{x}, y)$$

subject to $y^{(t)} 2 \{C + 1, ..., C + M\}$

- side information :
 - a single labeled example from each of the **M** new classes

Example

recognizing a person based on a single picture of them

ZERO-SHOT LEARNING

Training time

data:

$$\{\mathbf{x}^{(t)}, y^{(t)}\}$$

setting:

$$\mathbf{x}^{(t)}, y^{(t)} \leftarrow p(\mathbf{x}, y)$$

subject to $y^{(t)} 2 \{1, ..., C\}$

- side information :
 - description vector \mathbf{Z}_c of each of the \mathbf{C} classes

• Test time

data:

$$\{\mathbf{x}^{(t)}, y^{(t)}\}$$

setting:

$$\mathbf{x}^{(t)}, y^{(t)} \leftarrow p(\mathbf{x}, y)$$

subject to $y^{(t)} 2 \{C + 1, ..., C + M\}$

- side information :
 - description vector Z_c of each of the new M classes

Example

 recognizing an object based on a sentence description of it

UNSUPERVISED LEARNING

• Training time

data:

$$\{\mathbf{x}^{(t)}\}$$

setting:

$$\mathbf{x}^{(t)} \leftarrow p(\mathbf{x})$$

Test time

• data:

$$\{\mathbf{X}^{(t)}\}$$

setting:

$$\mathbf{x}^{(t)} \leftarrow p(\mathbf{x})$$

How can we train models "unsupervised"?

How can we train models "unsupervised"?

This is the focus of representation learning and generative models

This is the focus of representation learning

There's an entire co



ICLR has become one of the top CS and Engineering (not just AI) publication although it just started in 2013.

In fact top-10 in ALL OF SCIENCE

lop pu	blications		Q
Categories *			English ▼
	Publication	h5-index	h5-median
1.	Nature	<u>488</u>	745
2.	IEEE/CVF Conference on Computer Vision and Pattern Recognition	440	689
3.	The New England Journal of Medicine	<u>434</u>	897
4.	Science	<u>409</u>	633
5.	Nature Communications	<u>375</u>	492
6.	The Lancet	<u>368</u>	678
7.	Neural Information Processing Systems	<u>337</u>	614
8.	Advanced Materials	<u>327</u>	420
9.	Cell	320	482
10.	International Conference on Learning Representations	<u>304</u>	584

How

This is the

There's an entire



ICLR has become one of the top CS and Engineering (not just AI) publication although it just started in 2013.

In fact top-10 in ALL OF SCIENCE

International Conference on Learning Representations		Q
h5-index:304 h5-median:584 #2 Artificial Intelligence #4 Engineering & Computer Science		
Title / Author	Cited by	
An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. A Dosovitskiy, L Beyer, A Kolesnikov, D Weissenborn, X Zhai, ICLR	<u>38519</u>	
Decoupled Weight Decay Regularization. I Loshchilov, F Hutter ICLR (Poster)	<u>18046</u>	
Measuring and Improving the Use of Graph Information in Graph Neural Networks. Y Hou, J Zhang, J Cheng, K Ma, RTB Ma, H Chen, MC Yang ICLR	<u>7849</u>	
ALBERT: A Lite BERT for Self-supervised Learning of Language Representations. Z Lan, M Chen, S Goodman, K Gimpel, P Sharma, R Soricut ICLR	<u>7127</u>	
Large Scale GAN Training for High Fidelity Natural Image Synthesis. A Brock, J Donahue, K Simonyan ICLR	<u>5675</u>	
DARTS: Differentiable Architecture Search. H Liu, K Simonyan, Y Yang ICLR (Poster)	<u>4888</u>	
LoRA: Low-Rank Adaptation of Large Language Models. EJ Hu, Y Shen, P Wallis, Z Allen-Zhu, Y Li, S Wang, L Wang, W Chen ICLR	<u>4881</u>	
Deformable DETR: Deformable Transformers for End-to-End Object Detection. X Zhu, W Su, L Lu, B Li, X Wang, J Dai ICLR	4444	
BERTScore: Evaluating Text Generation with BERT. T Zhang, V Kishore, F Wu, KQ Weinberger, Y Artzi ICLR	4143	
ELECTRA: Pre-training Text Encoders as Discriminators Rather Than Generators. K Clark, MT Luong, QV Le, CD Manning ICLR	<u>3903</u>	

English *

h5-median

745

689

897

633

492

678

614

420

482

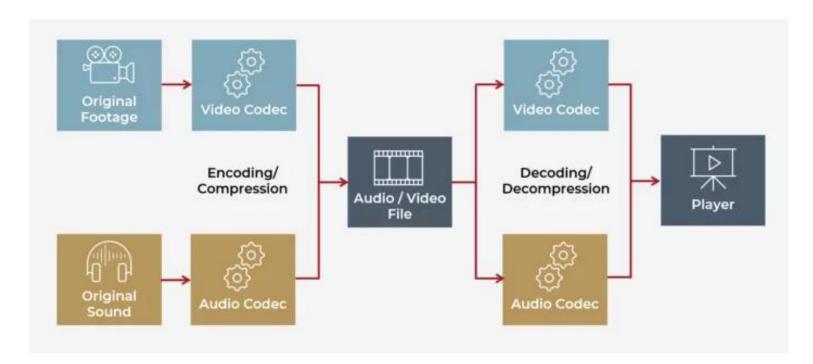
584

Warning

I might use the terms "latent", "embedding", "representation", "feature" interchangeably.

Motivation (kind of): Compression

- The idea is similar to compression (signal processing) or hashing (data structures):
 - o encode an image into a smaller vector s.t. you can decode it back to its original form
 - Example: images, audio, video are stored in a compressed form on your computer using compression algorithms like JPEG, MP3, MPEG etc. The computer has software to decode it back so that you can view it (everytime you "open" a JPEG file to view an image, the decoder runs and converts code to RGB)



Motivation (kind of): Compression

- The idea is similar to compression (signal processing) or hashing (information theory):
 - o encode an image into a smaller vector s.t. you can decode it back to its original form
 - Example: images, audio, video are stored in a compressed form on your computer using compression algorithms like JPEG, MP3, MPEG etc. The computer has software to decode it back so that you can view it (everytime you "open" a JPEG file to view an image, the decoder runs and converts code to RGB)

Representation Learning

- ~ convert inputs automatically into "codes" (called representations/embeddings / features) s.t. the representations are:
 - Useful for downstream tasks (e.g. classification, regression, ...)
 - o "explain the data" and are "meaningful"
- Main difference: "meaningful" representation spaces to do "tasks"
 - (The goal for compression is only efficient storage not data classification/clustering etc.)

Representation Learning Paradigm

Raw Data

(e.g. text, image, audio, video, ...)

Representation
Learning
(Encoding)

Do tasks / actions with these representations

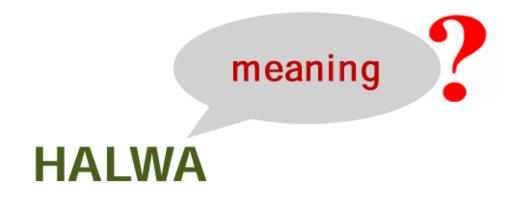
Similar representations for similar concepts







A Semblance of "Context" should be encoded ...



If you know the answer, don't share it with the class yet.

People from lands between Greece and India might know the answer ...

A Semblance of "Context" should be encoded ...



I am very <u>hungry</u>, I will <u>eat</u> HALWA

A Semblance of "Context" should be encoded ...



I am very *hungry*, I will *eat* HALWA

If you speak Marathi, this word has two meanings depending on context

Halwa (1): a food item

Halwa (2): (an instruction to) move (something)

derived from: Farsi

derived from: Sanskrit

A Semblance of "Context" should be encoded ...



Is it a good idea to eat HALWA after a meal?

A Semblance of "Context" should be encoded ...



Is it a good idea to eat HALWA after a meal?

A Semblance of "Context" should be encoded ...



Oh no! I forgot to put sugar in the HALWA

Parts, properties, attributes, ontology?

has subcategories

• "bird"

```
"bird" has "wing", "beak", "feathers"
"bird" can "fly"
"bird" is under category "animal"
```

"eagle", "peacock", "sparrow", "seagull", "pigeon"

Representation Learning is a Philosophy for Learning

Key assumptions in this philosophy:

- You can convert a high-dimensional input space into a low-dimensional representation space
 - \circ Example: RGB images \rightarrow 100 dim vectors
- A good representation space will have a "structure"
 - Example: Similarity, Symmetry, Relations will be easy to understand
 - Why? So that we can do arithmetic in representation space to do tasks
- Representations can be learned from data
- Representations can be leveraged for doing tasks

Parallel Work in Cog.Sci.

Trends in Cognitive Sciences



Volume 28, Issue 9, September 2024, Pages 844-856

Review

Why concepts are (probably) vectors

Steven T. Piantadosi ¹² O M, Dyana C.Y. Muller ², Joshua S. Rule ¹, Karthikeya Kaushik ¹, Mark Gorenstein ², Elena R. Leib ¹, Emily Sanford ¹

For decades, cognitive scientists have debated what kind of representation might characterize human concepts. Whatever the format of the representation, it must allow for the computation of varied properties, including similarities, features, categories, definitions, and relations. It must also support the development of theories, ad hoc categories, and knowledge of procedures. Here, we discuss why vector-based representations provide a compelling account that can meet all these needs while being plausibly encoded into neural architectures. This view has become especially promising with recent advances in both large language models and vector symbolic architectures. These innovations show how vectors can handle many properties traditionally thought to be out of reach for neural models, including compositionality, definitions, structures, and symbolic computational processes.

Highlights

Modern language models and vectorsymbolic architectures show that vector-based models are capable of handling the compositional, structured, and symbolic properties required for human concepts.

Vectors are also able to handle key phenomena from the psychology, including computation of features and similarities, reasoning about relations and analogies, and representation of theories.

Language models show how vector representation of word semantics and sentences can interface between concepts and language, as seen in definitional theories of concepts or ad hoc concepts.

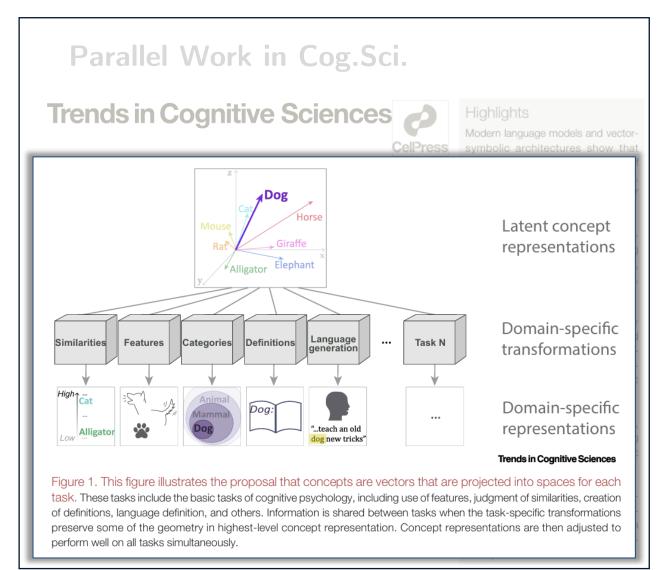
The idea of Church encoding, from logic, allows us to understand how meaning can arise in vector-based or symbolic systems.

By combining these recent computational results with classic findings in psychology, vector-based models provide a compelling account of human conceptual representation.

Representation Learning is a Philosophy for Learning

Key assumptions in this philosophy:

- You can convert a high-dimensional input space into a low-dimensional representation space
 - \circ Example: RGB images \rightarrow 100 dim vectors
- A good representation space will have a "structure"
 - Example: Similarity, Symmetry, Relations will be easy to understand
 - Why? So that we can do arithmetic in representation space to do tasks
- Representations can be learned from data
- Representations can be leveraged for doing tasks



Ok whatever. Tell us how it works ...

Types of Modeling (Probabilistic Interpretation)

Data: x; Label: y



"cat"

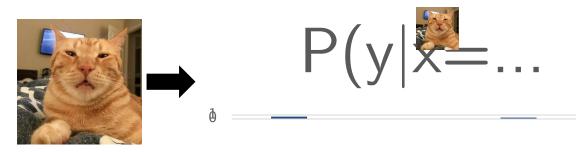
Density Function: p(x)

$$\int_X p(x) dx = 1$$

(probabilities of all inputs sum to 1)

Discriminative Model

Learn Prob. Dist. P(y|x)



Cat Harsa Timer



$$\forall x, \sum_{c} P(y = c | x) = 1$$

Types of Modeling (Probabilistic Interpretation)

Data: x; Label: y



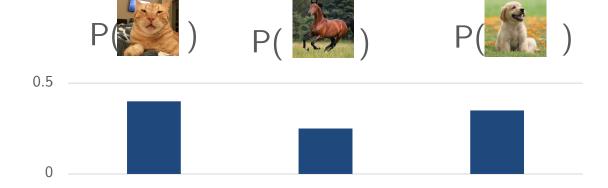
Density Function: p(x)

$$\int_X p(x) \ dx = 1$$

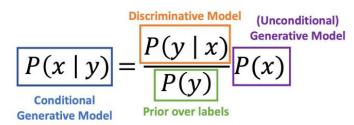
(probabilities of all inputs sum to 1)

Generative Model

Learn Marginal Prob. Dist. P(x)



Conditional Generative Model



Types of Modeling (Probabilistic Interpretation)

Data: x; Label: y



<u>"cat"</u>

Density Function: p(x)

$$\int_X p(x) dx = 1$$

(probabilities of all inputs sum to 1)

Discriminative Model

Learn Prob. Dist. P(y|x)

Generative Model

Learn Marginal Prob. Dist. P(x)

Conditional Generative Model

Types of Modeling (Probabilistic Interpretation) APPLICATIONS

Classification, Regression, Representation Learning (with labels)



Discriminative Model

Learn Prob. Dist. P(y|x)

Generative Model

Learn Marginal Prob. Dist. P(x)

Conditional Generative Model

Types of Modeling (Probabilistic Interpretation) APPLICATIONS

Discriminative Model

Learn Prob. Dist. P(y|x)

Data Generation
Outlier Detection
Representation Learning
(without labels)

Generative Model

Learn Marginal Prob. Dist. P(x)

Conditional Generative Model

Types of Modeling (Probabilistic Interpretation) APPLICATIONS

Discriminative Model

Learn Prob. Dist. P(y|x)

Generative Model

Learn Marginal Prob. Dist. P(x)

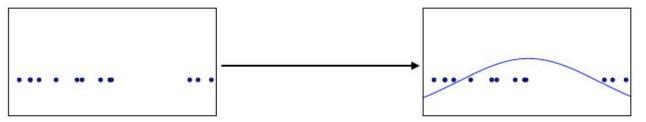
Machine Translation
Text-to-image generation

(pretty much every "GenAI" product you see is a conditional generative model)

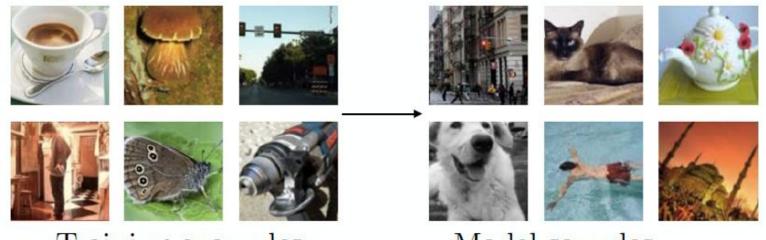
Conditional Generative Model

Generative Modeling

• Density estimation



• Sample generation



Training examples

Model samples

(Goodfellow 2016)

Why Generative Models?

- We've only seen discriminative models so far
 - Given an image X, predict a label Y
 - Estimates P(Y|X)
- Discriminative models have several key limitations
 - Can't model P(X), i.e. the probability of seeing a certain image
 - Thus, can't sample from P(X), i.e. can't generate new images
- · Generative models (in general) cope with all of above
 - Can model P(X)
 - Can generate new images

Generative Models

- What's a Generative Model?
 - \circ A model for the probability distribution of data x
 - A model that can be used to "generate" data

P(x)
marketing term "genAl"





- Generative Models can be learned
 - You are given some observed data X

(e.g. face images)

- \circ You choose a function (e.g. neural network) to model $P(x;\theta)$ using parameters θ
- \circ You estimate θ s.t. $P(x;\theta)$ best fits the observations X

Ok whatever. Tell us how it works ...

Let's start simple ...

The idea of an "Auto-encoder"

NN trained to reproduce the input

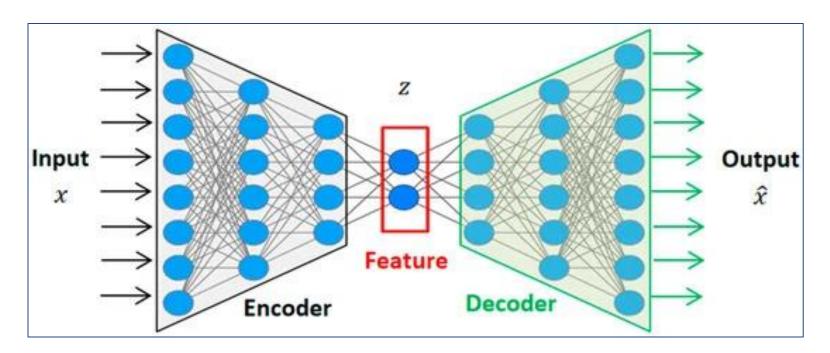
 $\widehat{x} = F(x)$

- F() is a composition of two functions:
 - Embedding / Feature / Latent
 - Output

encoder E() and decoder D()

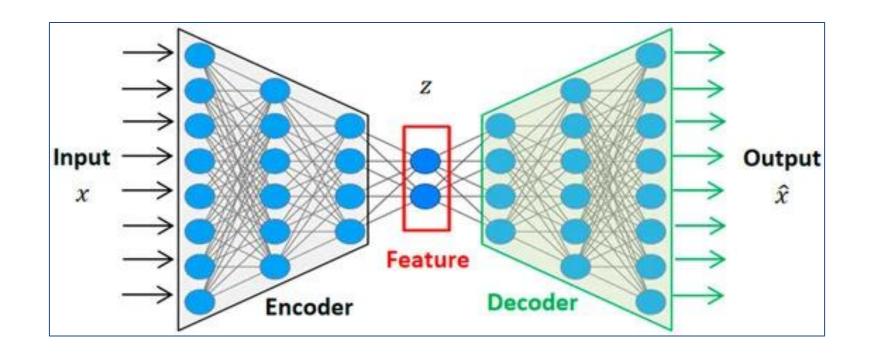
$$z = E(x)$$

$$\hat{x} = D(z) = D(E(x))$$



How would you train an autoencoder?

Loss Function?



Autoencoder: Loss Function

- The objective is to minimize the "distance" between x and \hat{x} • If $d(x,\hat{x}) = 0$ then we get perfect reconstruction
- Mean squared error!

$$l(f(\mathbf{x})) = \frac{1}{2} \sum_{k} (\widehat{x}_k - x_k)^2$$

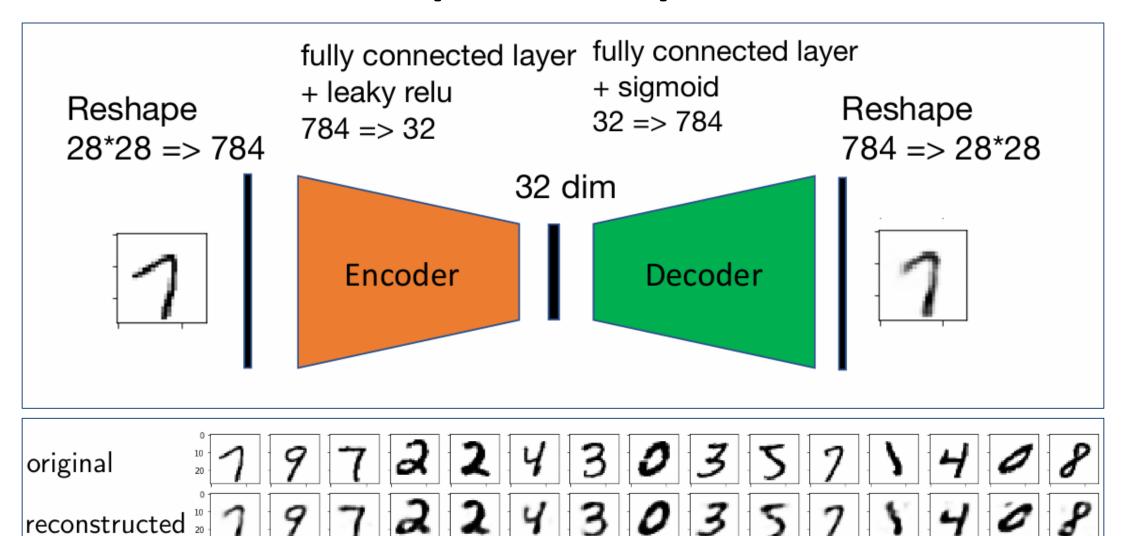
Cross Entropy (for binary inputs)

$$l(f(\mathbf{x})) = -\sum_{k} (x_k \log(\widehat{x}_k) + (1 - x_k) \log(1 - \widehat{x}_k))$$

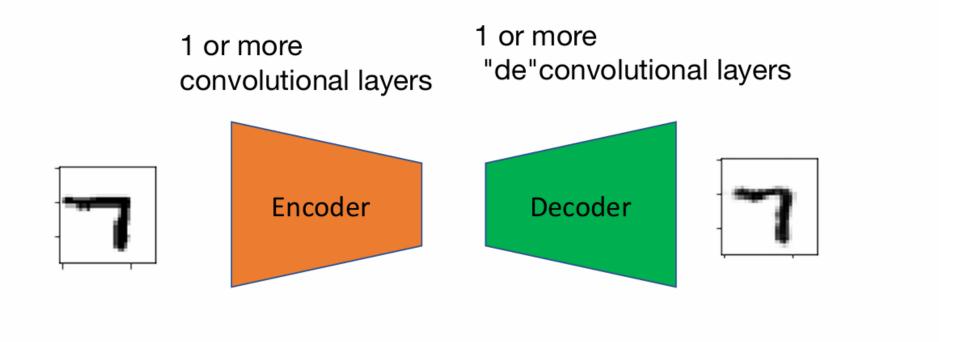
• For both cases, gradient is very simple:

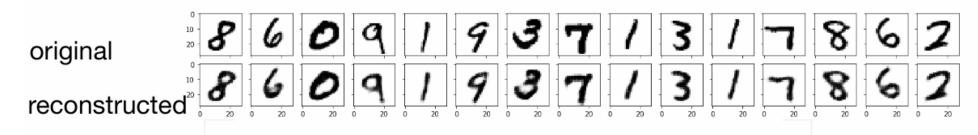
$$\nabla_{x} L(f(x), x) = \hat{x} - x$$

Autoencoder: Simple Example

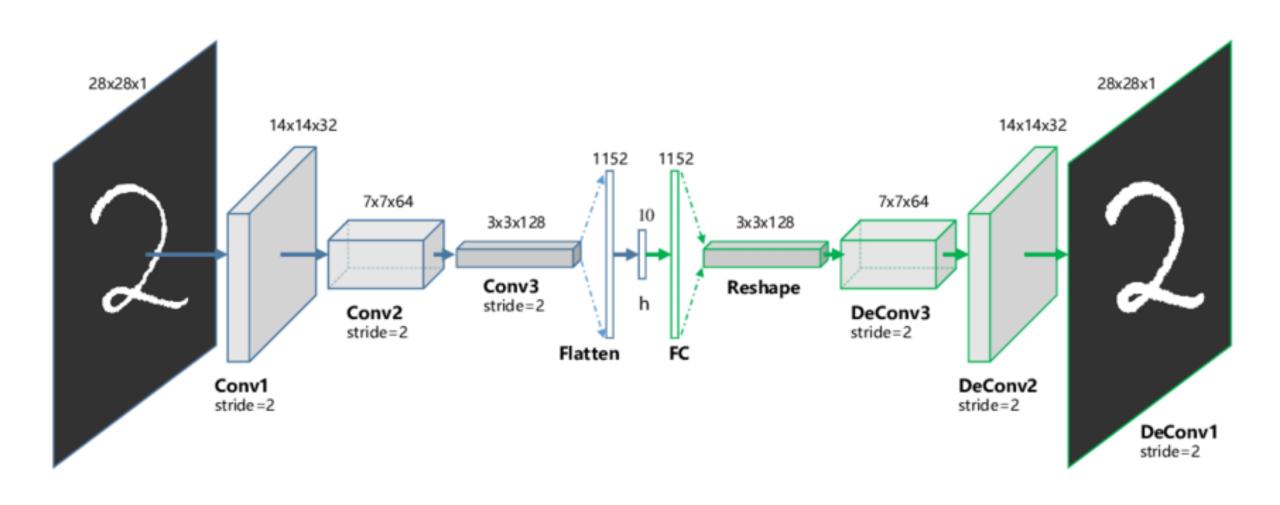


Convolutional Autoencoder





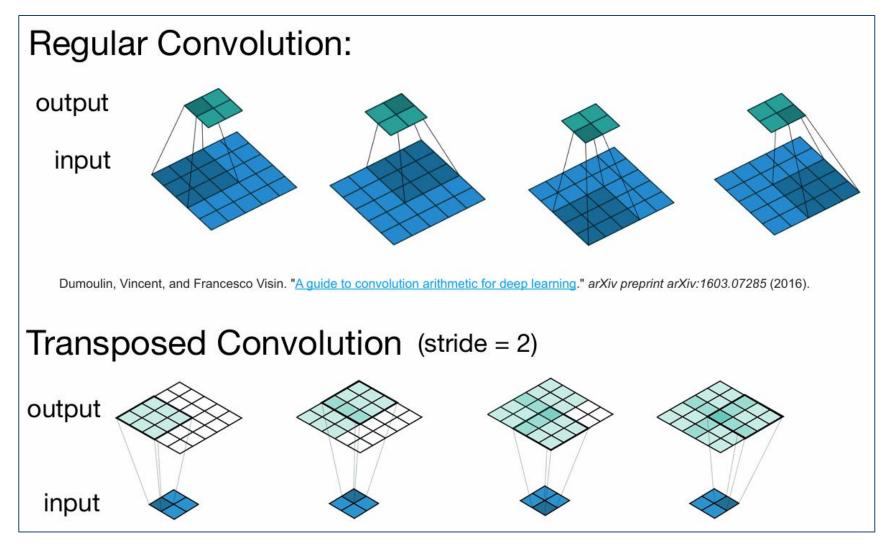
Convolutional Autoencoder



Convolutional Autoencoder: Expand Dimensions? Transposed Convolution!

- The decoder needs to "expand dimensions"
 - Convert a small feature z into a large input x
- Use transposed convolution! A.K.A. fractionally stride convolution
 - Often (incorrectly) called "de"convolution
 - This is an incorrect term because mathematically "deconvolution" is "inverse of convolution"

Convolutional Autoencoder: Expand Dimensions? Transposed Convolution!



Transposed Conv in PyTorch

: import torch

torch.manual_seed(123)

a = torch.rand(4).view(1, 1, 2, 2)

```
conv_t = torch.nn.ConvTranspose2d(in_channels=1,
                                    out_channels=1,
                                    kernel size=(3, 3),
                                    padding=0,
                                    stride=1)
  # output = s(n-1)+k-2p = 1*(2-1)+3-2*0 = 4
  conv t(a)
tensor([[[-0.2863, -0.2766, -0.1478, -0.3274],
            [-0.3522, -0.5356, -0.1591, -0.2911],
            [-0.3054, -0.4644, -0.3286, -0.2444],
            [-0.2332, -0.2557, -0.1876, -0.3970]]]]
         grad_fn=<ThnnConvTranspose2DBackward>)
  torch.manual_seed(123)
  a = torch.rand(16).view(1, 1, 4, 4)
  conv_t = torch.nn.ConvTranspose2d(in_channels=1,
                                    out channels=1,
                                    kernel_size=(3, 3),
                                    padding=0,
                                    stride=1)
  # output = s(n-1)+k-2p = 1*(4-1)+3-2*0 = 6
  conv_t(a).size()
  torch.Size([1, 1, 6, 6])
```

```
output = s(n-1) + k - 2p
```

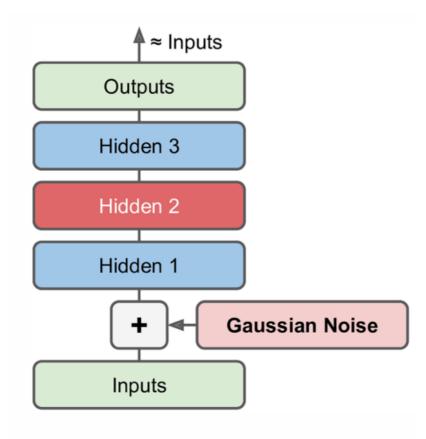
Denoising Autoencoder

- The input is "noisy" \tilde{x} . The expected output is a clean image (denoised image)
- Noise Examples:

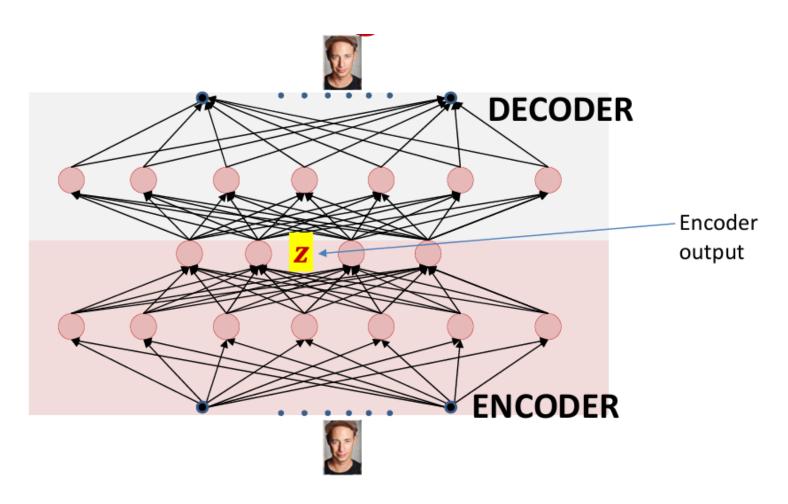
o Gaussian:
$$\tilde{x} = x + z$$
; $z \sim N(0, \sigma^2 I)$

- Masking: Zero-out some of the components of x
 (for images, make some pixels 0)
 - Can be random masks
 - Can be square masks
- Adding noise makes representations more robust

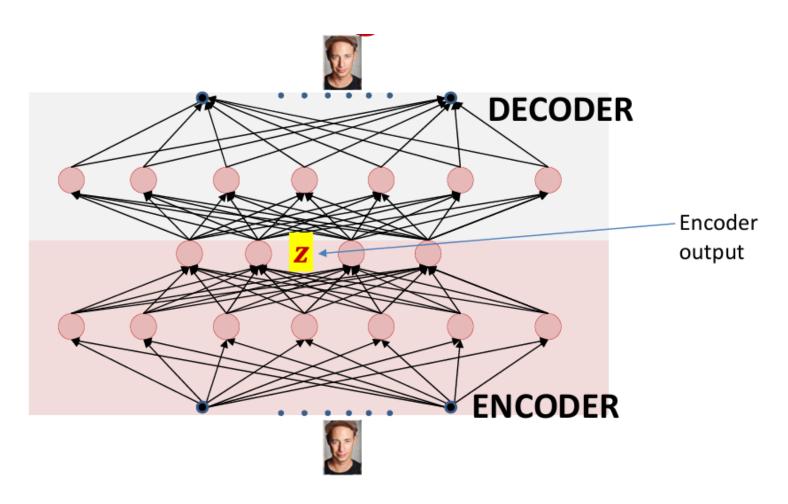
$$\circ$$
 Expect $D(E((\tilde{x})) = x$ for all z



Once trained, what can you do with this model?



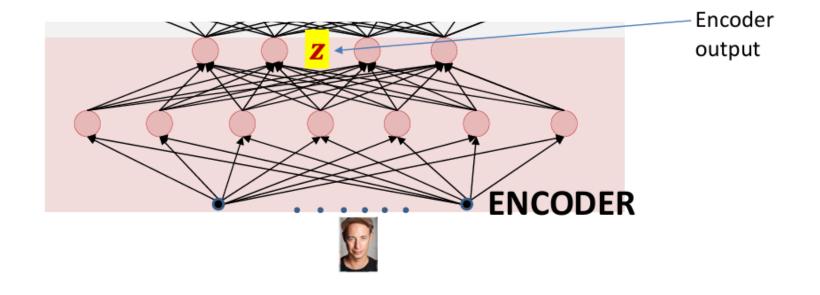
Once trained, what can you do with this model?



Once trained, what can you do with this model?

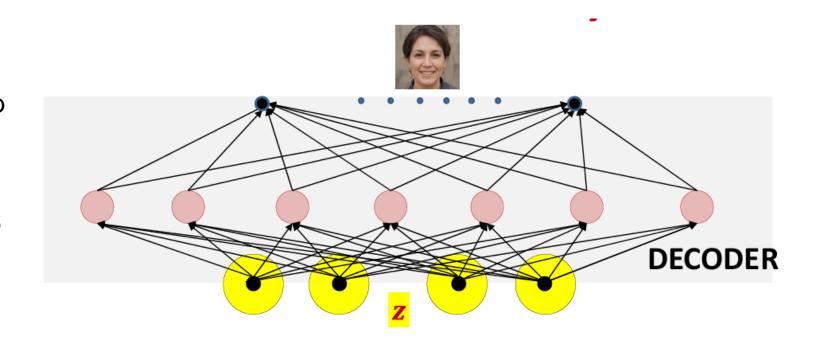
(1) Encode images into vectors

(throw away the decoder ...)



Once trained, what can you do with this model?

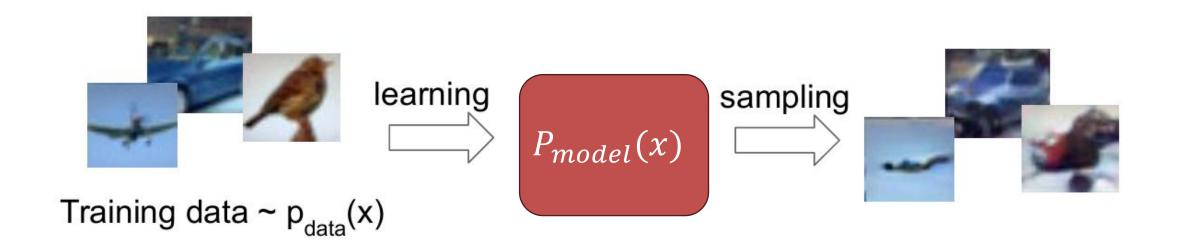
(1) Encode images into vectors



(2) Generate new faces ...

(throw away the encoder)

With Generative Models, there are 2 objectives:



Objectives:

- 1. Learn $p_{model}(x)$ that approximates $p_{data}(x)$
- 2. Sampling new x from $p_{model}(x)$

An Auto-Encoder is a Generative Model

Probabilistic Interpretation:

$$P_{\theta_E}(z|x)$$

• Decoder
$$D()$$
 estimates

$$P_{\theta_D}(x|z)$$

 The marginal Bayes/Chain Rule ...

$$P(x) = \int_{z} P(x,z) dz = \int_{z} P(z)P(x|z)$$

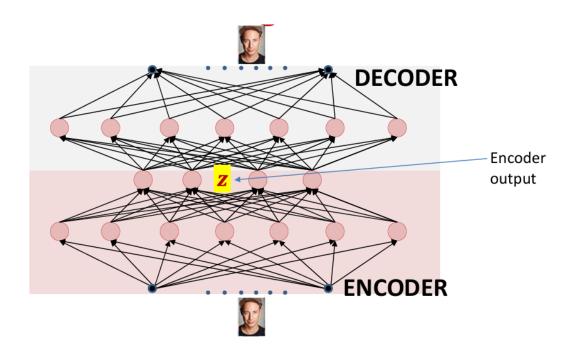
- Once the AE is "trained"
 - o you get a generative model that generates "x" given a latent code "z"
 - A conditional generative model takes an additional input "y"

- E.g. generating images from text

$$y=text, x = image$$

More on this later ...

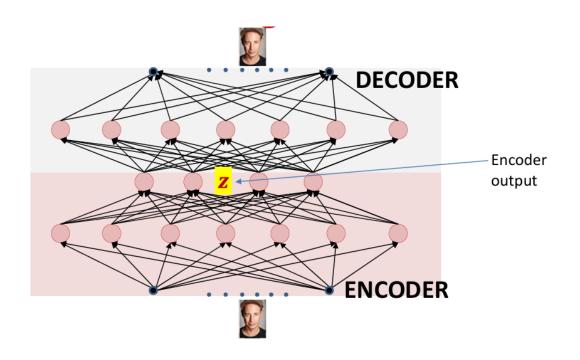
Types of Autoencoders



So far, we have not enforced any "structure" on the latents z

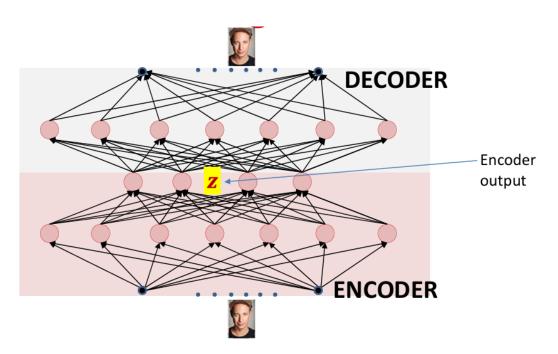
But a structure is desirable

(Remember our motivations / goals for representation learning)



So far, we have not enforced any "structure" on the latents z

We can't generate $\underline{\textit{new images}}$ from D() if we don't understand the z-space



So far, we have not enforced any "structure" on the latents z

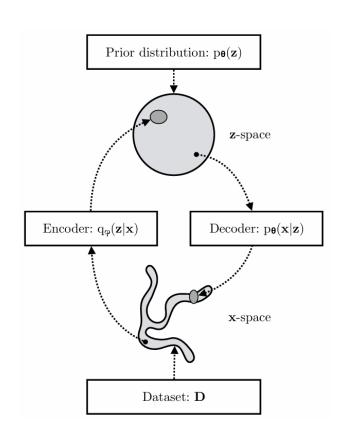
We can't generate $\underline{\textit{new images}}$ from D() if we don't understand the z-space

For example, if I ask you to generate a "face with beard, glasses, brown hair" which z would you choose?

VAE: Variational Autoencoder

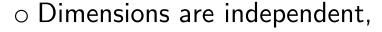




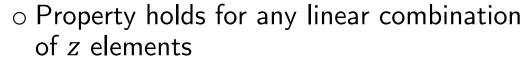


- Force a "prior" distribution on the latent space
 - \circ Example: Gaussian N(0,I)
- Gaussians are nice because they are perfectly symmetrical in every dimension

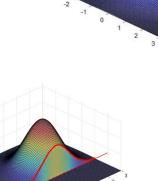


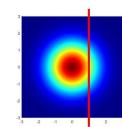


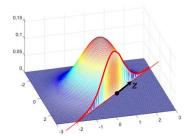
i.e.
$$P(z_1|z_2) = P(z_1) = N(0,I) \forall z_1, z_2$$



- i.e.
$$P(z_1|az_2 + bz_3) = N(0, I)$$





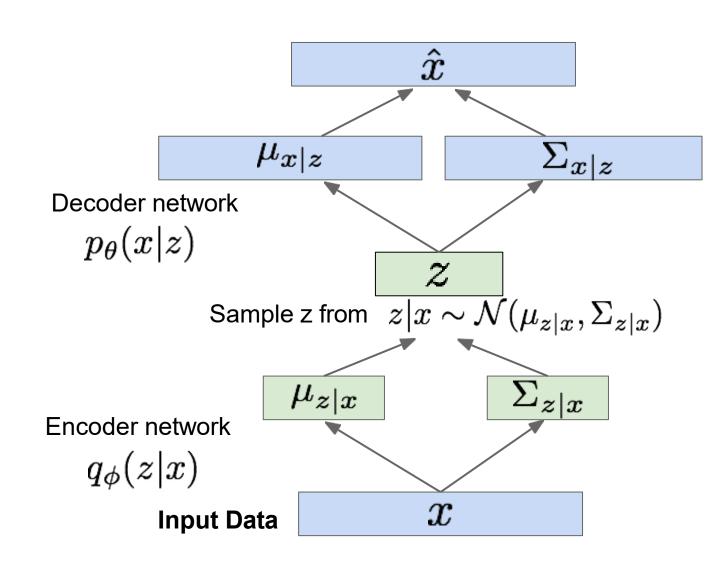


Variational Autoencoders

Key idea:

Force a Gaussian distribution on the latent space *z*

Details: take my NN class!



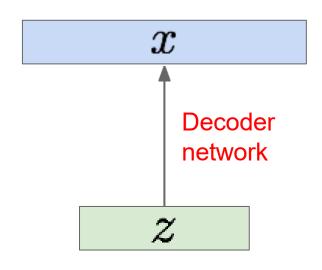
Our assumption about data generation process

Sample from true conditional

$$p_{\theta^*}(x \mid z^{(i)})$$

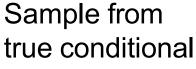
Sample from true prior

$$z^{(i)} \sim p_{ heta^*}(z)$$



Our assumption about data generation process

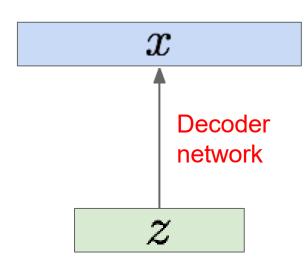
Now given a trained VAE: use decoder network & sample z from prior!

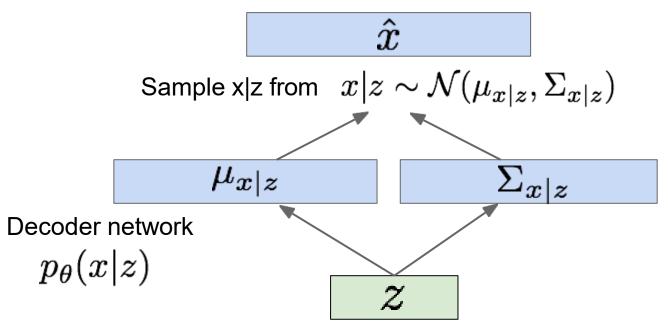


$$p_{\theta^*}(x \mid z^{(i)})$$

Sample from true prior

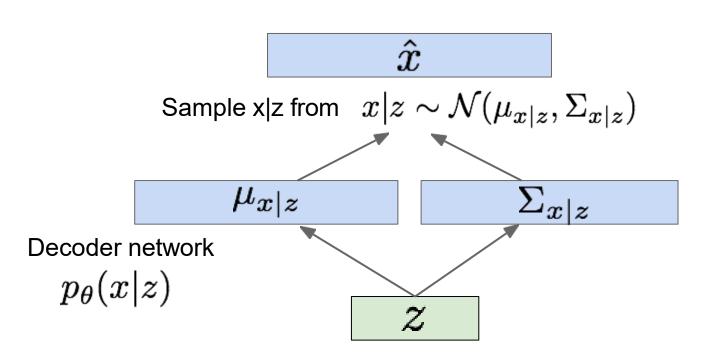
$$z^{(i)} \sim p_{ heta^*}(z)$$





Sample z from $\,z \sim \mathcal{N}(0,I)\,$

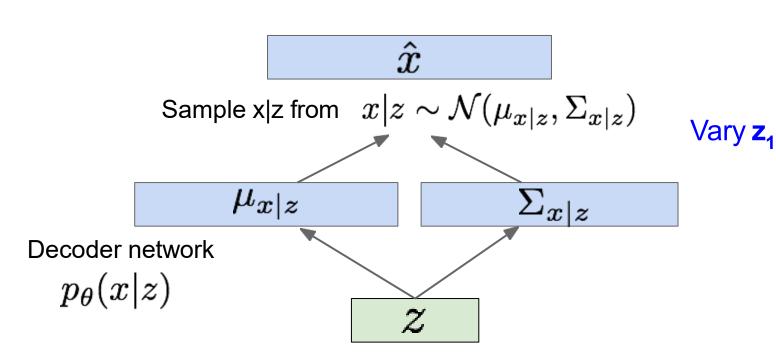
Use decoder network. Now sample z from prior!



Sample z from $\,z \sim \mathcal{N}(0,I)\,$

```
6666666666666
```

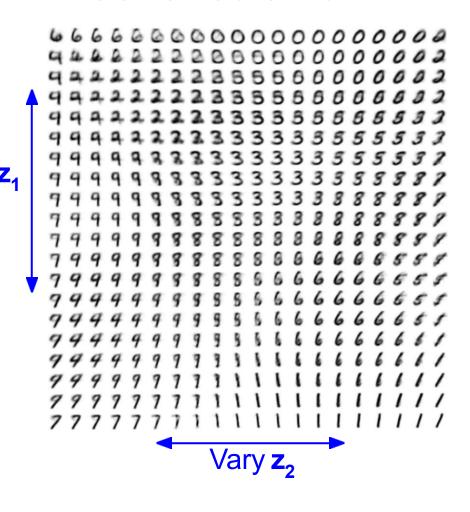
Use decoder network. Now sample z from prior!



Sample z from $z \sim \mathcal{N}(0, I)$

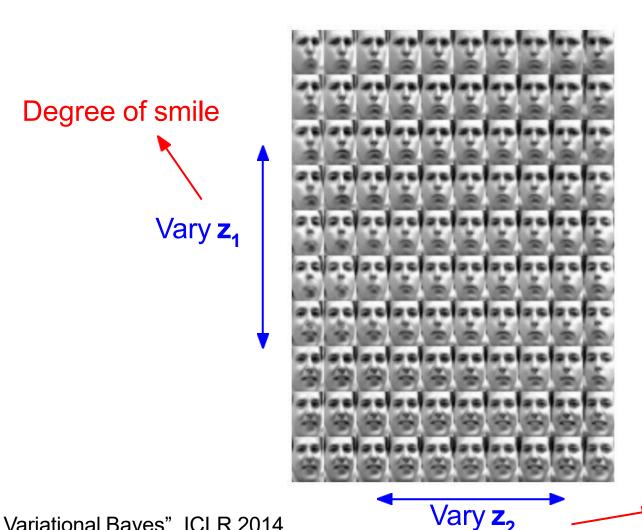
Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Data manifold for 2-d z



Diagonal prior on **z** => independent latent variables

Different dimensions of **z** encode interpretable factors of variation



Head pose

Diagonal prior on **z** => independent latent variables

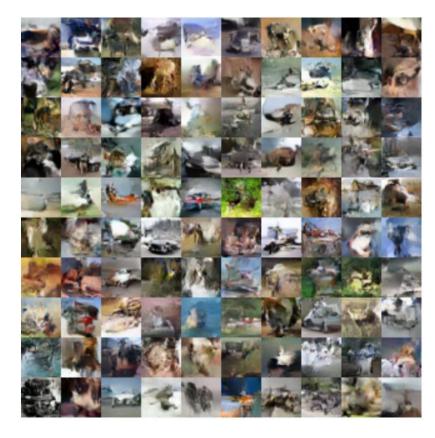
Different dimensions of **z** encode interpretable factors of variation

Also good feature representation that can be computed using $q_{\phi}(z|x)!$

Degree of smile Vary **z**₁

Vary z,

Head pose

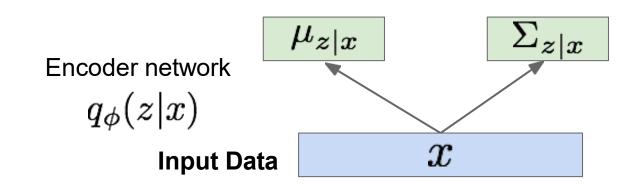


32x32 CIFAR-10

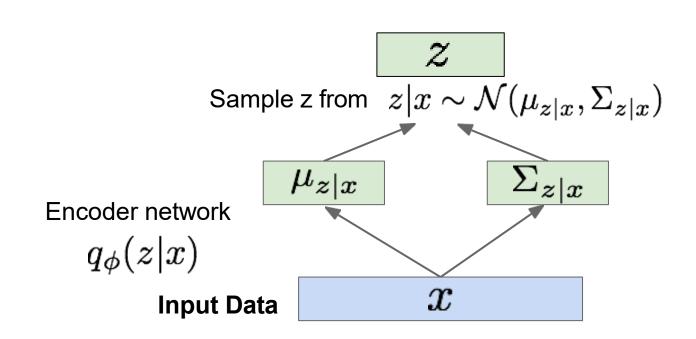


Labeled Faces in the Wild

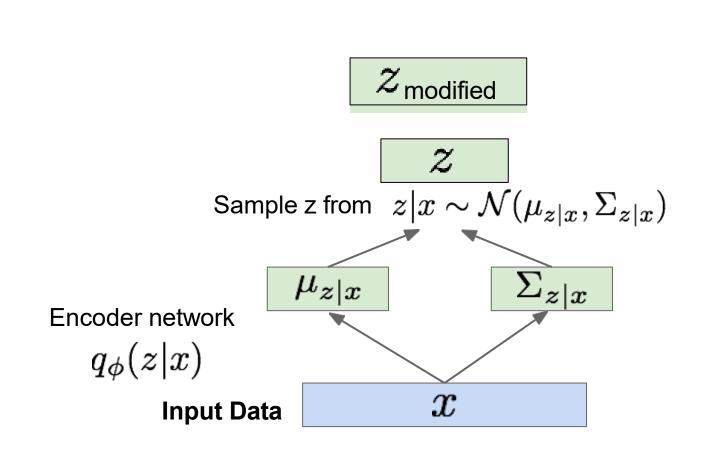
1. Run input data through encoder to get a distribution over latent codes



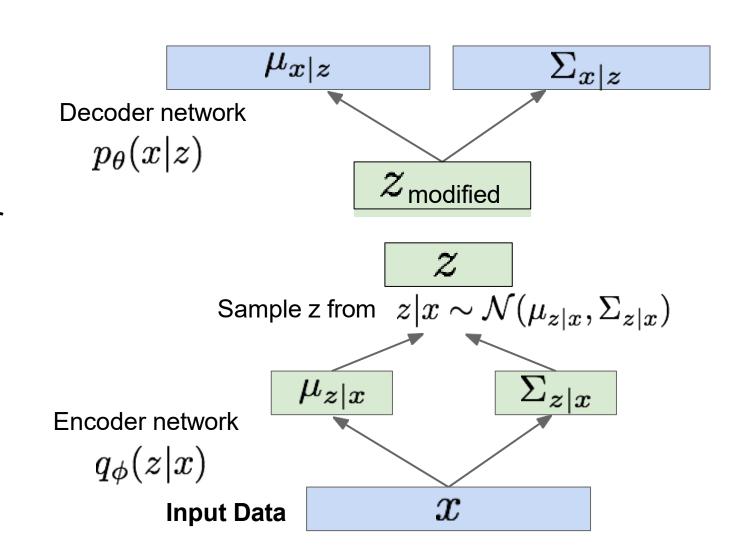
- Run input data through encoder to get a distribution over latent codes
- 2. Sample code z from encoder output



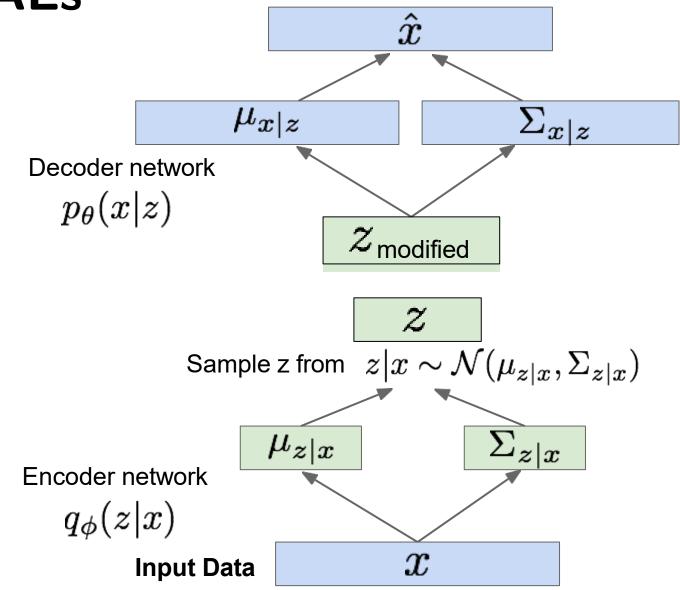
- Run input data through encoder to get a distribution over latent codes
- 2. Sample code z from encoder output
- 3. Modify some dimensions of sampled code

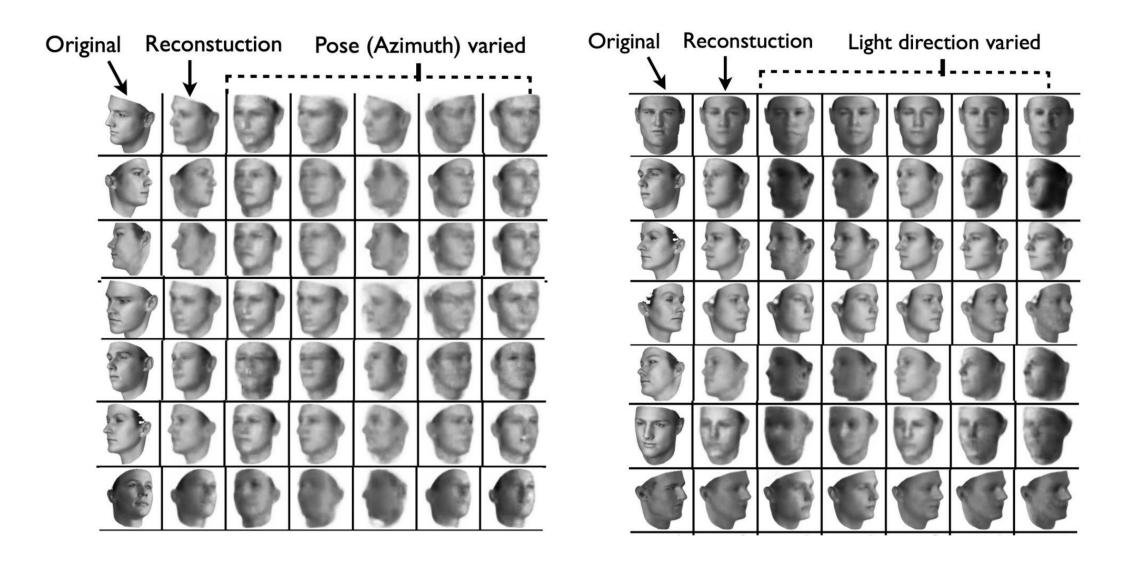


- Run input data through encoder to get a distribution over latent codes
- 2. Sample code z from encoder output
- 3. Modify some dimensions of sampled code
- 4. Run modified z through decoder to get a distribution over data sample



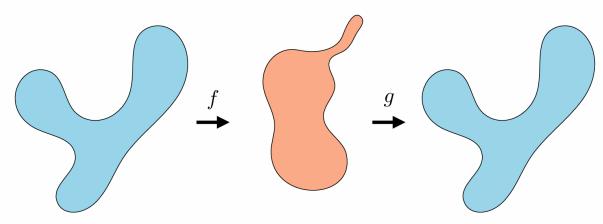
- Run input data through encoder to get a distribution over latent codes
- 2. Sample code z from encoder output
- Modify some dimensions of sampled code
- 4. Run modified z through decoder to get a distribution over data sample
- 5. Sample new data from (4)





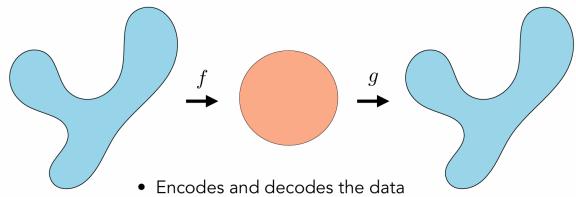
To Summarize: AE, VAE

Autoencoder



- Encodes and decodes the data
- Low-dimensional bottleneck

Variational Autoencoder

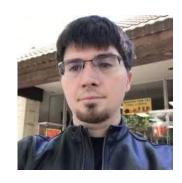


- Low-dimensional bottleneck
- Gaussian bottleneck (can sample; disentangled)

Generative Adversarial Networks

Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, Yoshua Bengio[‡]

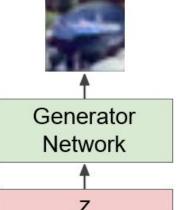
Département d'informatique et de recherche opérationnelle Université de Montréal Montréal, QC H3C 3J7



NIPS 2014

- Problem: We want to sample from a high-dimensional training distribution p(x)
 - But there is no direct way to do this ...
 - We don't know which z maps to which image (so we can't use autoencoders)
- We know how to sample from a random distribution (e.g. Gaussian)
 - \circ Can we map a random distribution <u>directly</u> to p(x)?

Output: Sample from training distribution

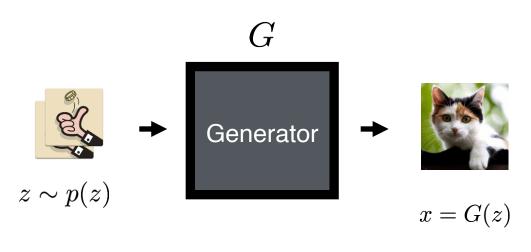


Input: Random noise

Generative Adversarial Networks

Goal: Map all z to some realistic-looking x

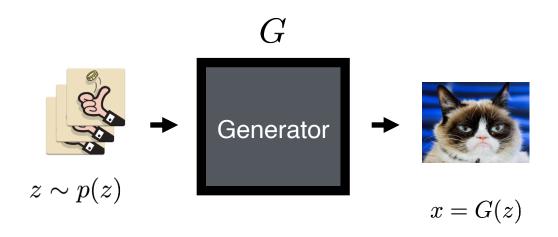
Image synthesis from "noise"



Sampler

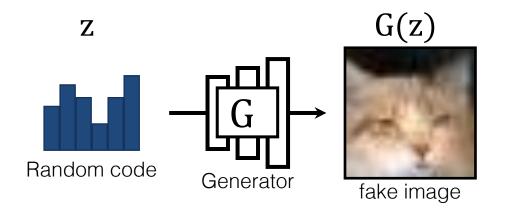
$$G: \mathcal{Z} \to \mathcal{X}$$
$$z \sim p(z)$$
$$x = G(z)$$

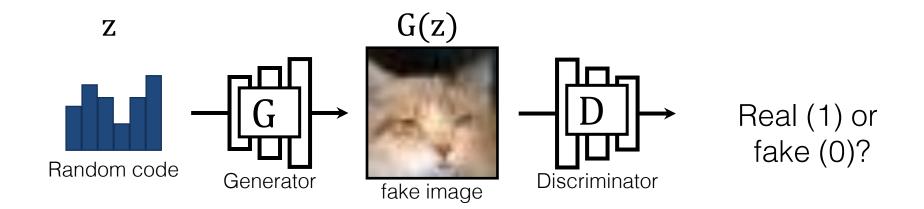
Image synthesis from "noise"



Sampler

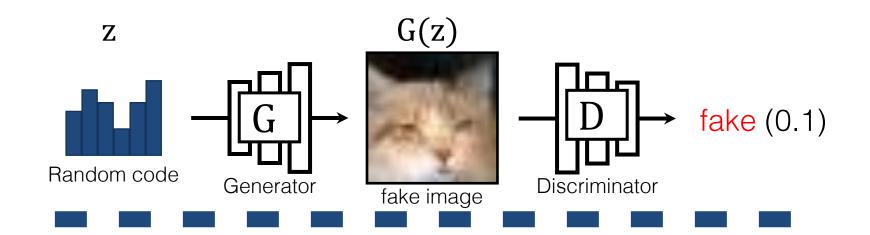
$$G: \mathcal{Z} \to \mathcal{X}$$
$$z \sim p(z)$$
$$x = G(z)$$



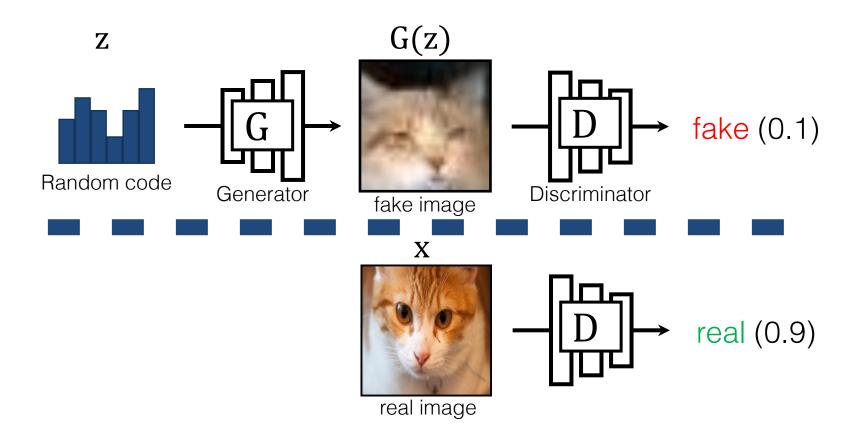


A two-player game:

- G tries to generate fake images that can fool D.
- D tries to detect fake images.

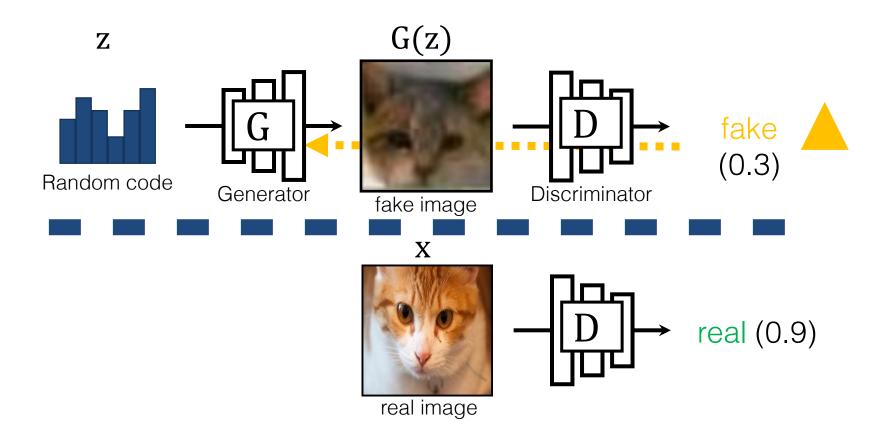


Learning objective (GANs)
$$\min_{G} \max_{D} \mathbb{E}_z[\log(1 - D(G(z))]$$



Learning objective (GANs)

$$\min_{G} \max_{D} \mathbb{E}_z[\log(1 - D(G(z))] + \mathbb{E}_x[\log D(x)]$$



Learning objective (GANs)

$$\min_{G} \max_{D} \mathbb{E}_{z}[\log(1 - D(G(z)))] + \mathbb{E}_{x}[\log D(x)]$$

GAN Training Breakdown

- From the discriminator D's perspective:
 - o binary classification: real vs. fake.
 - Nothing special: similar to 1 vs. 7 or cat vs. dog

$$\max_{D} \mathbb{E}[\log(1 - D(\square))] + \mathbb{E}[\log D(\square)]$$

GAN Training Breakdown

- From the discriminator D's perspective:
 - o binary classification: real vs. fake.
 - Nothing special: similar to 1 vs. 7 or cat vs. dog

$$\max_{D} \mathbb{E}[\log(1-D(\mathbf{Q}))] + \mathbb{E}[\log D(\mathbf{Q})]$$

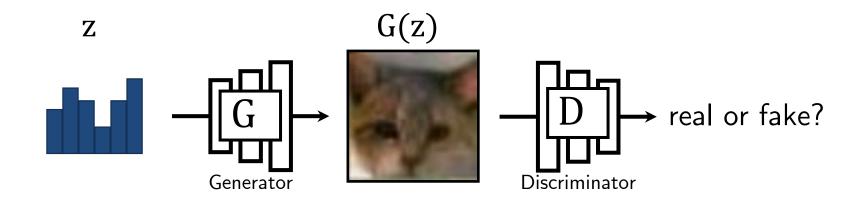
- o From the generator G's perspective:
 - Optimizing a loss that depends on a classifier D

$$\min_{G} \mathbb{E}_{z}[\mathcal{L}_{D}(G(z))] \qquad \min_{G} \mathbb{E}_{(x,y)}||F(G(x)) - F(y)||$$

GAN loss for G

Perceptual Loss for G

GAN Training Breakdown

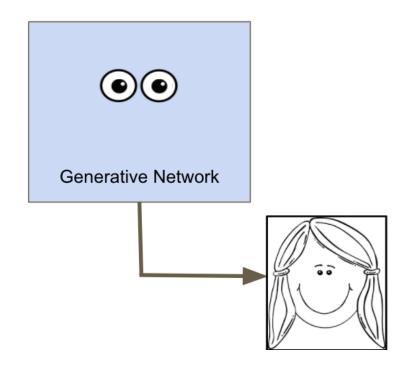


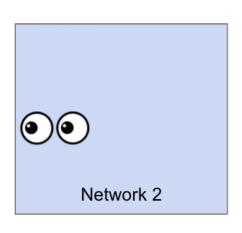
G tries to synthesize fake images that fool D

D tries to identify the fakes

- Training: iterate between training D and G with backprop.
- Global optimum when G reproduces data distribution.

Discriminator network: try to distinguish between real and fake images **Generator network**: try to fool the discriminator by generating real-looking images



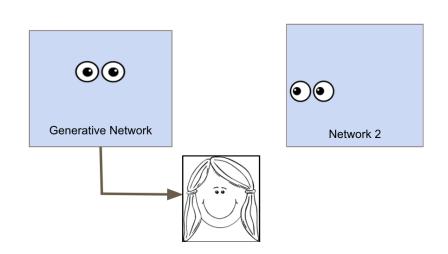




Discriminator network:

Generator network:

try to distinguish between real and fake images try to fool the discriminator by generating real-looking images





Connection to Game Theory: Zero-Sum "Minimax" Game

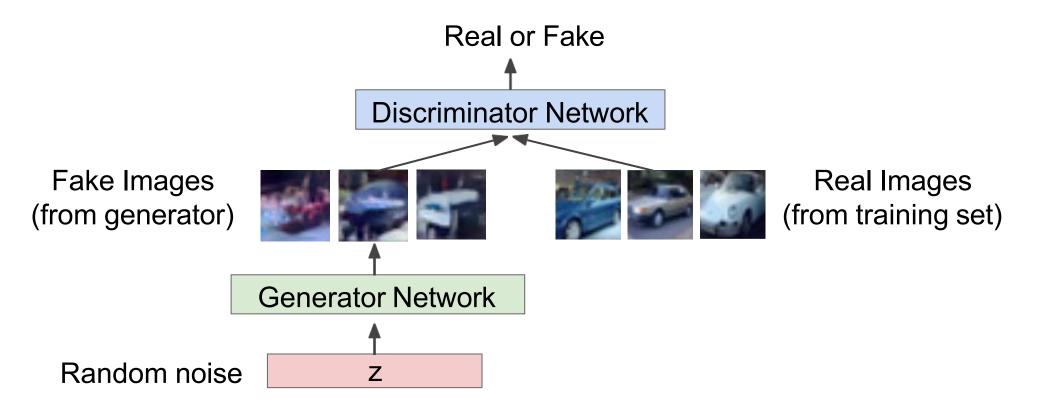
- Each player trying to minimize the opponent's profits
- Each player trying to maximize their own profits

Discriminator network:

Generator network:

try to distinguish between real and fake images

try to fool the discriminator by generating real-looking images

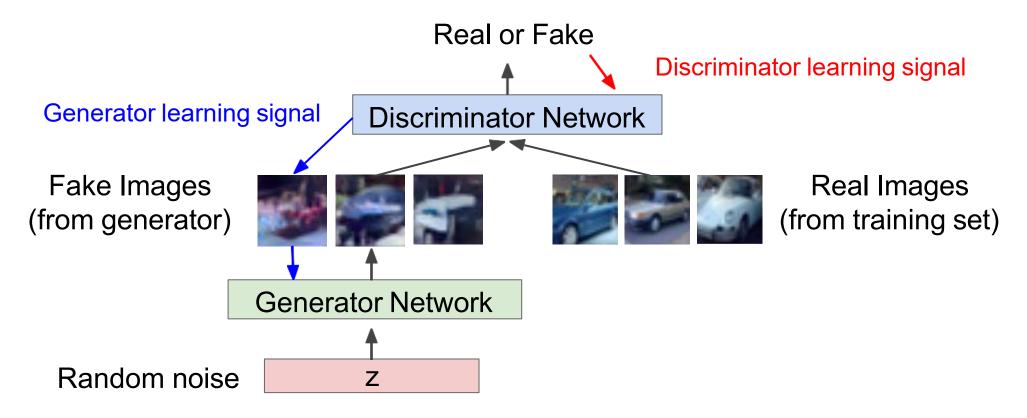


Fake and real images copyright Emily Denton et al. 2015. Reproduced with permission.

Discriminator network:

Generator network:

try to distinguish between real and fake images try to fool the discriminator by generating real-looking images



Fake and real images copyright Emily Denton et al. 2015. Reproduced with permission.

Discriminator network: try to distinguish between real and fake images

Generator network: try to fool the discriminator by generating real-looking images

Train jointly in minimax game

Minimax objective function:

$$\min_{\substack{\theta_g \text{ Mode objective}}} \max_{\substack{\theta_d \text{ Discriminator objective}}} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log (1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

try to distinguish between real and fake images **Discriminator network:**

try to fool the discriminator by generating real-looking images **Generator network:**

Train jointly in **minimax game**

Minimax objective function:

Discriminator outputs likelihood in (0,1) of real image

$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log (1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$
 Discriminator output for for real data x generated fake data G(z)

- Discriminator (θ_d) wants to **maximize objective** s.t. D(x) is close to 1 (real) and D(G(z)) is close to 0 (fake)
- Generator (θ_{q}) wants to **minimize objective** s.t. D(G(z)) is close to 1 (discriminator is fooled into thinking generated G(z) is real)

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log (1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Alternate between:

1. Gradient ascent on discriminator

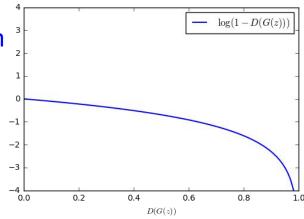
$$\max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log (1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

2. Gradient descent on generator

$$\min_{ heta_g} \mathbb{E}_{z \sim p(z)} \log (1 - D_{ heta_d}(G_{ heta_g}(z)))$$
 fake, want to learn from

In practice, optimizing this generator objective does not work well!

When sample is likely fake, want to learn from it to improve generator (move to the right on X axis).



Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log (1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Alternate between:

1. Gradient ascent on discriminator

$$\max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log (1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Gradient descent on generator

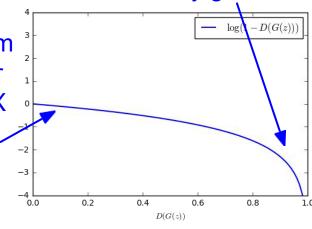
$$\min_{ heta_g} \mathbb{E}_{z \sim p(z)} \log (1 - D_{ heta_d}(G_{ heta_g}(z)))$$
 fake, want to learn from

In practice, optimizing this generator objective does not work well!

When sample is likely fake, want to learn from it to improve generator (move to the right on X axis).

But gradient in this region is relatively flat!

Gradient signal dominated by region where sample is already good



Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log (1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Alternate between:

1. Gradient ascent on discriminator

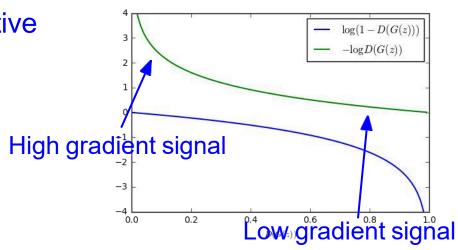
$$\max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log (1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

2. Instead: Gradient ascent on generator, different objective

$$\max_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(D_{\theta_d}(G_{\theta_g}(z)))$$

Instead of minimizing likelihood of discriminator being correct, now maximize likelihood of discriminator being wrong.

Same objective of fooling discriminator, but now higher gradient signal for bad samples => works much better! Standard in practice.



Putting it together: GAN training algorithm

for number of training iterations do

for k steps do

- Sample minibatch of m noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Sample minibatch of m examples $\{x^{(1)}, \ldots, x^{(m)}\}$ from data generating distribution $p_{\text{data}}(x)$.
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D_{\theta_d}(x^{(i)}) + \log(1 - D_{\theta_d}(G_{\theta_g}(z^{(i)}))) \right]$$

end for

- Sample minibatch of m noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_q(z)$.
- Update the generator by ascending its stochastic gradient (improved objective):

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log(D_{\theta_d}(G_{\theta_g}(z^{(i)})))$$

end for

Putting it together: GAN training algorithm

for number of training iterations do for k steps do

- Sample minibatch of m noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Sample minibatch of m examples $\{x^{(1)}, \ldots, x^{(m)}\}$ from data generating distribution $p_{\text{data}}(x)$.
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D_{\theta_d}(x^{(i)}) + \log(1 - D_{\theta_d}(G_{\theta_g}(z^{(i)}))) \right]$$

Some find k=1 more stable, others use k > 1, no best rule.

Followup work (e.g. Wasserstein GAN, BEGAN) alleviates this problem, better stability!

end for

- Sample minibatch of m noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_q(z)$.
- Update the generator by ascending its stochastic gradient (improved objective):

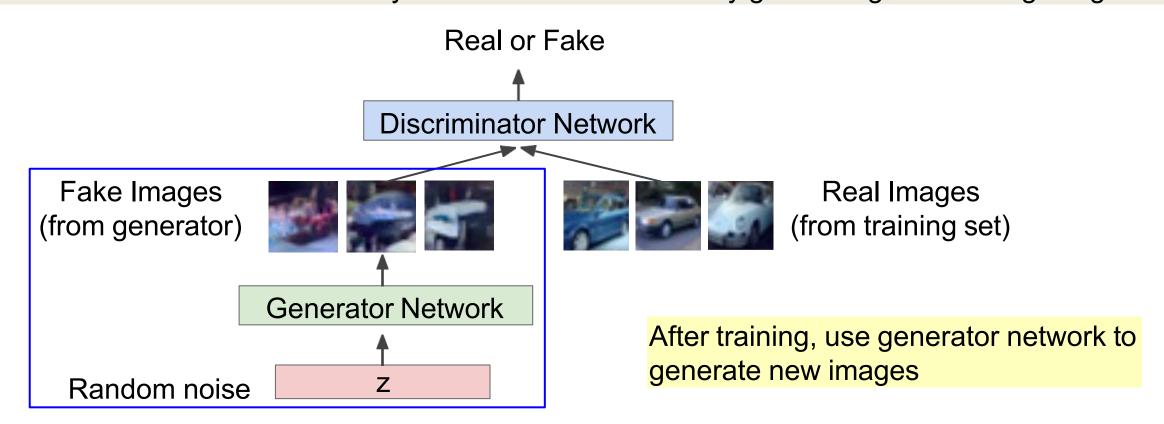
$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log(D_{\theta_d}(G_{\theta_g}(z^{(i)})))$$

end for

Discriminator network:

Generator network:

try to distinguish between real and fake images try to fool the discriminator by generating real-looking images



Fake and real images copyright Emily Denton et al. 2015. Reproduced with permission.

Generative Adversarial Nets: Convolutional Architectures

Generator is an upsampling network with fractionally-strided convolutions Discriminator is a convolutional network

Architecture guidelines for stable Deep Convolutional GANs

- Replace any pooling layers with strided convolutions (discriminator) and fractional-strided convolutions (generator).
- Use batchnorm in both the generator and the discriminator.
- Remove fully connected hidden layers for deeper architectures.
- Use ReLU activation in generator for all layers except for the output, which uses Tanh.
- Use LeakyReLU activation in the discriminator for all layers.

Radford et al, "Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks", ICLR 2016

Generative Adversarial Nets: Convolutional Architectures

Samples from the model look much better!

Generative Adversarial Nets: Convolutional Architectures

Interpolating between random points in laten space

Generative Adversarial Nets: Interpretable Vector Math

Samples from the model

Smiling woman Neutral woman



Neutral man



Generative Adversarial Nets: Interpretable Vector Math

Neutral woman Neutral man Smiling woman Samples from the model Average Z vectors, do arithmetic

Generative Adversarial Nets: Interpretable Vector Math

Neutral woman Neutral man Smiling woman Samples from the model Average Z vectors, do arithmetic

Radford et al, ICLR 2016

Smiling Man



Generative Adversarial Nets: Interpretable Vector Math

Glasses man

No glasses man

No glasses woman

Radford et al, **ICLR 2016**











Woman with glasses



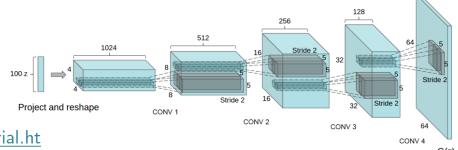








GAN in PyTorch



https://pytorch.org/tutorials/beginner/dcgan_faces_tutorial.ht

```
ml
 class Discriminator(nn.Module):
     def __init__(self, ngpu):
         super(Discriminator, self).__init__()
         self.ngpu = ngpu
         self.main = nn.Sequential(
             # input is ``(nc) x 64 x 64``
             nn.Conv2d(nc, ndf, 4, 2, 1, bias=False),
             nn.LeakyReLU(0.2, inplace=True),
             # state size. ''(ndf) x 32 x 32''
             nn.Conv2d(ndf, ndf * 2, 4, 2, 1, bias=False),
             nn.BatchNorm2d(ndf * 2),
             nn.LeakyReLU(0.2, inplace=True),
             # state size. ``(ndf*2) x 16 x 16``
             nn.Conv2d(ndf * 2, ndf * 4, 4, 2, 1, bias=False),
             nn.BatchNorm2d(ndf * 4),
             nn.LeakyReLU(0.2, inplace=True),
             # state size. ``(ndf*4) x 8 x 8``
             nn.Conv2d(ndf * 4, ndf * 8, 4, 2, 1, bias=False),
             nn.BatchNorm2d(ndf * 8),
             nn.LeakyReLU(0.2, inplace=True),
             # state size. ``(ndf*8) x 4 x 4``
             nn.Conv2d(ndf * 8, 1, 4, 1, 0, bias=False),
             nn.Sigmoid()
     def forward(self, input):
         return self.main(input)
```

```
class Generator(nn.Module):
   def __init__(self, ngpu):
        super(Generator, self).__init__()
        self.ngpu = ngpu
       self.main = nn.Sequential(
            # input is Z, going into a convolution
            nn.ConvTranspose2d( nz, ngf * 8, 4, 1, 0, bias=False),
            nn.BatchNorm2d(ngf * 8),
            nn.ReLU(True),
            # state size. ``(ngf*8) x 4 x 4``
            nn.ConvTranspose2d(ngf \star 8, ngf \star 4, 4, 2, 1, bias=False),
            nn.BatchNorm2d(ngf * 4),
            nn.ReLU(True),
            # state size. ``(ngf*4) x 8 x 8``
            nn.ConvTranspose2d( ngf * 4, ngf * 2, 4, 2, 1, bias=False),
            nn.BatchNorm2d(ngf * 2),
            nn.ReLU(True),
            # state size. ``(ngf*2) x 16 x 16``
            nn.ConvTranspose2d( ngf * 2, ngf, 4, 2, 1, bias=False),
            nn.BatchNorm2d(ngf),
            nn.ReLU(True),
            # state size. ''(ngf) x 32 x 32''
            nn.ConvTranspose2d( ngf, nc, 4, 2, 1, bias=False),
            nn.Tanh()
            # state size. ''(nc) x 64 x 64''
   def forward(self, input):
       return self.main(input)
```

Since then: Explosion of GANs

"The GAN Zoo"

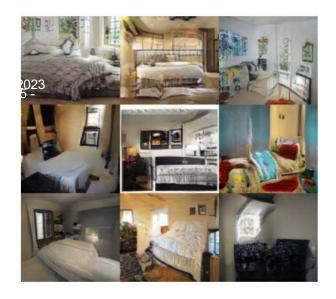
- GAN Generative Adversarial Networks
- 3D-GAN Learning a Probabilistic Latent Space of Object Shapes via 3D Generative-Adversarial Modeling
- acGAN Face Aging With Conditional Generative Adversarial Networks
- AC-GAN Conditional Image Synthesis With Auxiliary Classifier GANs
- AdaGAN AdaGAN: Boosting Generative Models
- · AEGAN Learning Inverse Mapping by Autoencoder based Generative Adversarial Nets
- AffGAN Amortised MAP Inference for Image Super-resolution
- AL-CGAN Learning to Generate Images of Outdoor Scenes from Attributes and Semantic Layouts
- ALI Adversarially Learned Inference
- · AM-GAN Generative Adversarial Nets with Labeled Data by Activation Maximization
- AnoGAN Unsupervised Anomaly Detection with Generative Adversarial Networks to Guide Marker Discovery
- · ArtGAN ArtGAN: Artwork Synthesis with Conditional Categorial GANs
- b-GAN b-GAN: Unified Framework of Generative Adversarial Networks
- Bayesian GAN Deep and Hierarchical Implicit Models
- BEGAN BEGAN: Boundary Equilibrium Generative Adversarial Networks
- BiGAN Adversarial Feature Learning
- BS-GAN Boundary-Seeking Generative Adversarial Networks
- CGAN Conditional Generative Adversarial Nets
- CaloGAN CaloGAN: Simulating 3D High Energy Particle Showers in Multi-Layer Electromagnetic Calorimeters with Generative Adversarial Networks
- · CCGAN Semi-Supervised Learning with Context-Conditional Generative Adversarial Networks
- · CatGAN Unsupervised and Semi-supervised Learning with Categorical Generative Adversarial Networks
- CoGAN Coupled Generative Adversarial Networks

- Context-RNN-GAN Contextual RNN-GANs for Abstract Reasoning Diagram Generation
- C-RNN-GAN C-RNN-GAN: Continuous recurrent neural networks with adversarial training
- CS-GAN Improving Neural Machine Translation with Conditional Sequence Generative Adversarial Nets
- · CVAE-GAN CVAE-GAN: Fine-Grained Image Generation through Asymmetric Training
- CycleGAN Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks
- DTN Unsupervised Cross-Domain Image Generation
- . DCGAN Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks
- DiscoGAN Learning to Discover Cross-Domain Relations with Generative Adversarial Networks
- . DR-GAN Disentangled Representation Learning GAN for Pose-Invariant Face Recognition
- DualGAN DualGAN: Unsupervised Dual Learning for Image-to-Image Translation
- EBGAN Energy-based Generative Adversarial Network
- f-GAN f-GAN: Training Generative Neural Samplers using Variational Divergence Minimization
- FF-GAN Towards Large-Pose Face Frontalization in the Wild
- GAWWN Learning What and Where to Draw
- GeneGAN GeneGAN: Learning Object Transfiguration and Attribute Subspace from Unpaired Data
- Geometric GAN Geometric GAN
- GoGAN Gang of GANs: Generative Adversarial Networks with Maximum Margin Ranking
- GP-GAN GP-GAN: Towards Realistic High-Resolution Image Blending
- IAN Neural Photo Editing with Introspective Adversarial Networks
- iGAN Generative Visual Manipulation on the Natural Image Manifold
- IcGAN Invertible Conditional GANs for image editing
- ID-CGAN Image De-raining Using a Conditional Generative Adversarial Network
- Improved GAN Improved Techniques for Training GANs
- · InfoGAN InfoGAN: Interpretable Representation Learning by Information Maximizing Generative Adversarial Nets
- LAGAN Learning Particle Physics by Example: Location-Aware Generative Adversarial Networks for Physics Synthesis
- LAPGAN Deep Generative Image Models using a Laplacian Pyramid of Adversarial Networks

https://github.com/hindupuravinash/the-gan-zoo

2017: Explosion of GANs

Better training and generation



LSGAN, Zhu 2017.



Wasserstein GAN, Arjovsky 2017. Improved Wasserstein GAN, Gulrajani 2017.





Progressive GAN, Karras 2018.

Some challenges with GANs ...

Challenges with GANs

Vanishing gradients:

o the discriminator becomes too good and the generator gradient vanishes.

Non-Convergence:

o the generator and discriminator oscillate without reaching an equilibrium.

Mode Collapse:

o the generator distribution collapses to a small set of examples.

Mode Dropping:

o the generator distribution doesnt fully cover the data distribution.

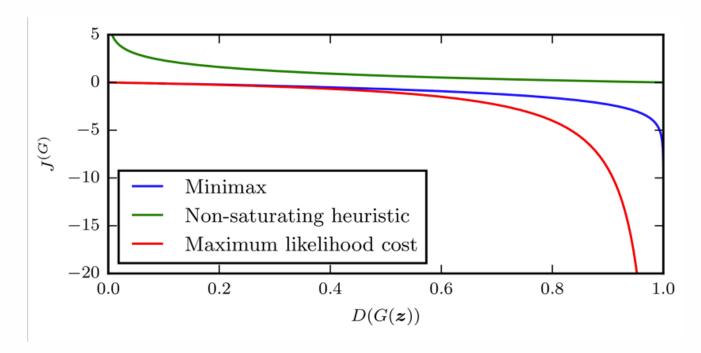
Challenges with GANs: Vanishing Gradients

• The **minimax** objective saturates when D_{θ_d} is close to perfect:

$$V(\theta_d, \theta_g) = \mathbb{E}_{p_{\mathsf{data}}} \left[\log D_{\theta_d}(\mathbf{x}) \right] + \mathbb{E}_{p_{\mathbf{z}}(\mathbf{z})} \left[\log \left(1 - D_{\theta_d}(G_{\theta_g}(\mathbf{z})) \right) \right].$$

• A **non-saturating heuristic** objective for the generator is

$$J(G_{\theta_g}) = -\mathbb{E}_{p_{\mathbf{z}}(\mathbf{z})} \left[\log \left(D_{\theta_d}(G_{\theta_g}(\mathbf{z})) \right) \right].$$



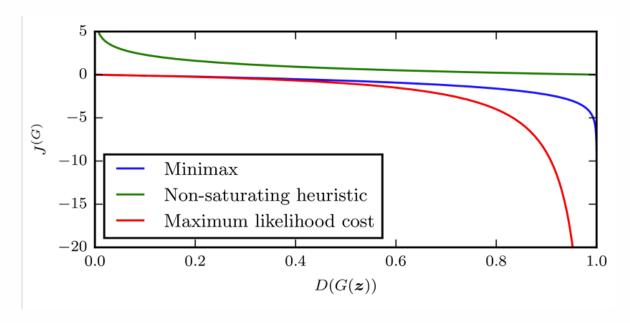
Challenges with GANs: Vanishing Gradients

• The **minimax** objective saturates when D_{θ_d} is close to perfect:

$$V(\theta_d, \theta_g) = \mathbb{E}_{p_{\text{data}}} \left[\log D_{\theta_d}(\mathbf{x}) \right] + \mathbb{E}_{p_{\mathbf{z}}(\mathbf{z})} \left[\log \left(1 - D_{\theta_d}(G_{\theta_g}(\mathbf{z})) \right) \right].$$

A non-saturating heuristic objective for the generator is

$$J(G_{\theta_g}) = -\mathbb{E}_{p_{\mathbf{z}}(\mathbf{z})} \left[\log \left(D_{\theta_d}(G_{\theta_g}(\mathbf{z})) \right) \right].$$



Potential Solutions:

- 1. Explore other training objectives?
- 2. Discriminator Capacity:
 - make it small?
 - train it less?
 - slow learning rate?
- 3. Learning Schedule:
 - try to balance trainingG and D

Problems: Nonconvergence

Deep Learning models (in general) involve a single player

- The player tries to maximize its reward (minimize its loss).
- Use SGD (with Backpropagation) to find the optimal parameters.
- SGD has convergence guarantees (under certain conditions).
- Problem: With non-convexity, we might converge to local optima.

$$\min_{G} L(G)$$

GANs instead involve two (or more) players

- Discriminator is trying to maximize its reward.
- Generator is trying to minimize Discriminator's reward.

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

- SGD was not designed to find the Nash equilibrium of a game.
- Problem: We might not converge to the Nash equilibrium at all.

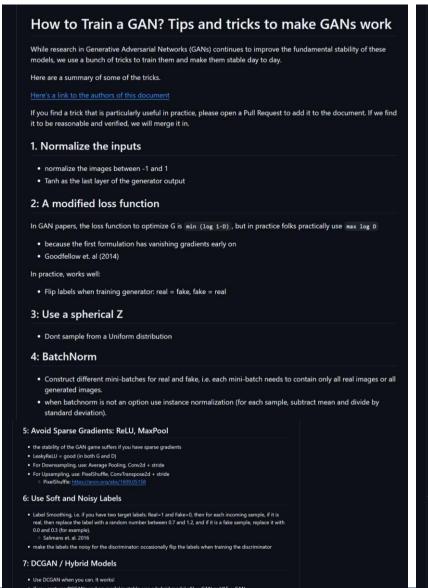
Challenges with GANs: Non-Convergence

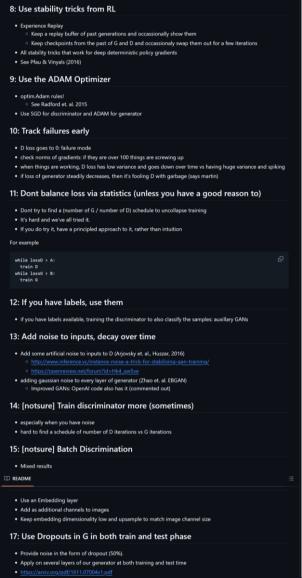
- Simultaneous gradient descent is not guaranteed to converge for minimax objectives.
- Goodfellow et al. only showed convergence when updates are made in the function space.
- The parameterization of D and G results in highly non-convex objective.
- In practice, training tends to oscillate – updates undo each other!

Challenges with GANs: Non-Convergence

- Simultaneous gradient descent is not guaranteed to converge for minimax objectives.
- Goodfellow et al. only showed convergence when updates are made in the function space.
- The parameterization of D and G results in highly non-convex objective.
- In practice, training tends to oscillate – updates undo each other!

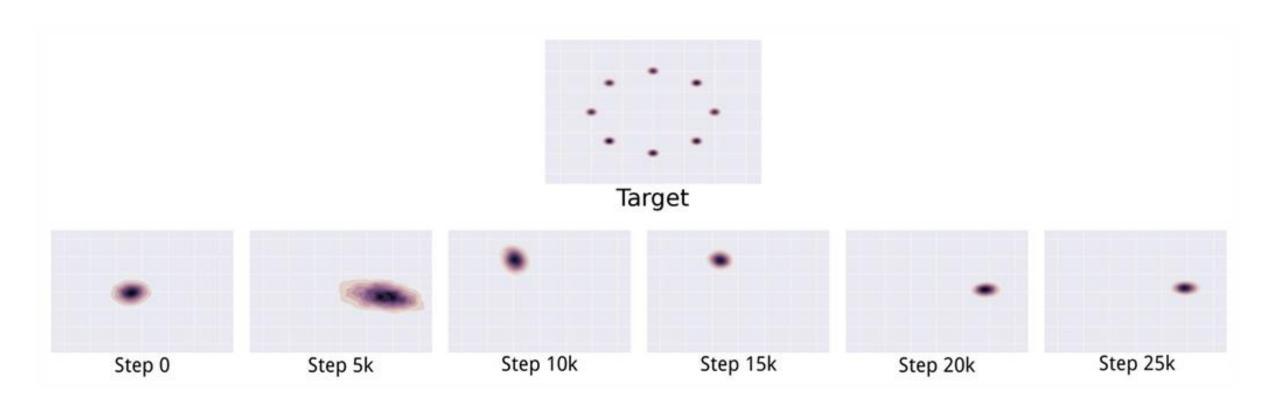
Potential Solutions (HACKS) https://github.com/soumith/ganhacks





Challenges with GANs: Mode Collapse

The generator maps all z values to the x that is mostly likely to fool the discriminator.



Some real examples



Reed, S., et al. Generating interpretable images with controllable structure. Technical report, 2016. 2, 2016.

Challenges with GANs: Mode Collapse

Possible Solutions:

There are a large variety of divergence measures for distributions:

• f-Divergences: (e.g. Jensen-Shannon, Kullback-Leibler)

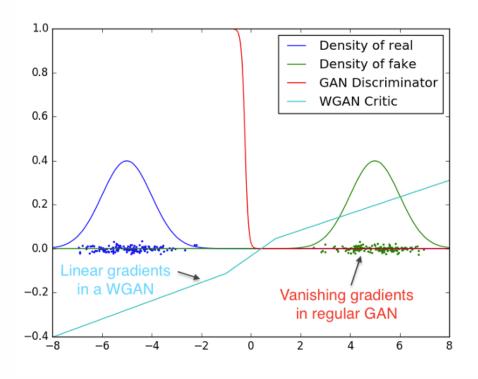
$$D_f (P || Q) = \int_{\chi} q(\mathbf{x}) f(\frac{p(\mathbf{x})}{q(\mathbf{x})}) d\mathbf{x}$$

- GANs [2], f-GANs [7], and more.
- Integral Probability Metrics: (e.g. Earth Movers Distance, Maximum Mean Discrepancy)

$$\gamma_F (P \mid\mid Q) = \sup_{f \in F} \left| \int f dP - \int f dQ \right|$$

 Wasserstein GANs [1], Fisher GANs [6], Sobolev GANs [5] and more.

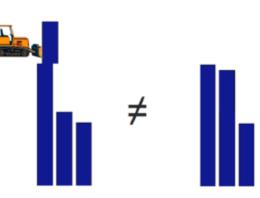
Wasserstein GAN



Wasserstein GANs

WGAN

- If our data are on a low-dimensional manifold of a high dimensional space the model's manifold and the true data manifold can have a negligible intersection in practice
- KL divergence is undefined or infinite
- The loss function and gradients may not be continuous and well behaved
- The Earth Mover's Distance is well defined:
 - Minimum transportation cost for making one pile of dirt (pdf/pmf) look like the other

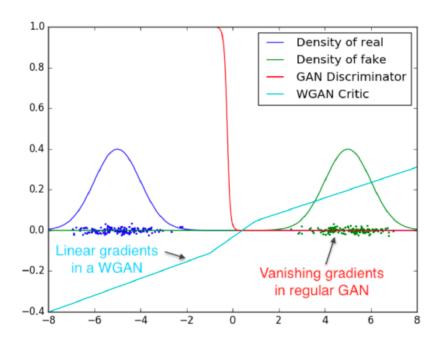


Wasserstein GANs

WGAN

$$J^{(D)}(\theta^{(D)}, \theta^{(G)}) = -\left[\mathbb{E}_{x \sim p_{data}} D(x) - \mathbb{E}_{z} D(G(z))\right]$$
$$J^{(G)}(\theta^{(D)}, \theta^{(G)}) = -\mathbb{E}_{z} D(G(z))$$

- Importantly, the discriminator is trained for many steps before the generator is updated
- Gradient-clipping is used in the discriminator to ensure D(x) has the Lipschitz continuity required by the theory
- The authors argue that this solves many training issues, including mode collapse



Conditional GANs

MNIST digits generated conditioned on their class label.

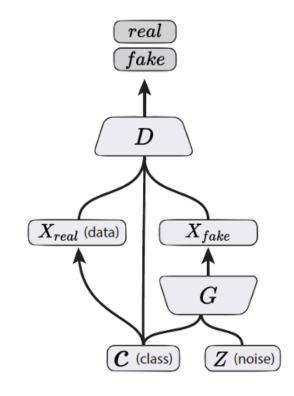
```
[1, 0, 0, 0, 0, 0, 0, 0, 0, 0]
                                                               00004
[0, 1, 0, 0, 0, 0, 0, 0, 0, 0]
[0, 0, 1, 0, 0, 0, 0, 0, 0, 0]
[0, 0, 0, 1, 0, 0, 0, 0, 0, 0]
[0, 0, 0, 0, 1, 0, 0, 0, 0, 0]
[0, 0, 0, 0, 0, 1, 0, 0, 0, 0]
[0, 0, 0, 0, 0, 0, 1, 0, 0, 0]
[0, 0, 0, 0, 0, 0, 0, 1, 0, 0]
[0, 0, 0, 0, 0, 0, 0, 0, 1, 0]
[0, 0, 0, 0, 0, 0, 0, 0, 0, 1]
```

Figure 2 in the original paper.

Conditional GANs

 Simple modification to the original GAN framework that conditions the model on additional information for better multi-modal learning.

• Lends to many practical applications of GANs when we have explicit *supervision* available.



Conditional GAN (Mirza & Osindero, 2014)

Image Credit: Figure 2 in Odena, A., Olah, C. and Shlens, J., 2016. Conditional image synthesis with auxiliary classifier GANs. arXiv preprint arXiv:1610.09585.

Image-to-Image Translation

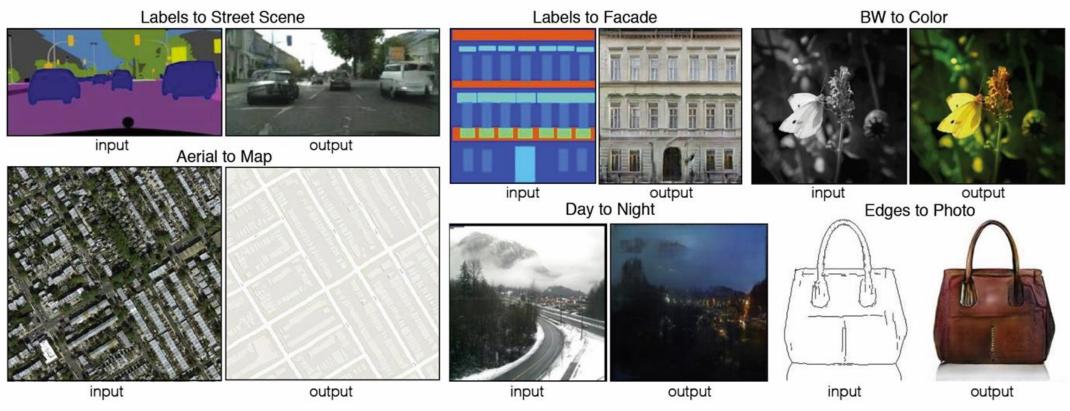


Figure 1 in the original paper.

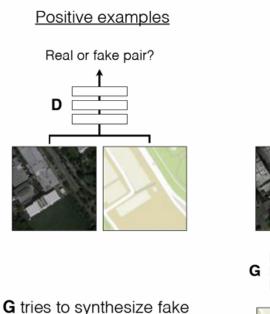
Link to an interactive demo of this paper

Image-to-Image Translation

Architecture: DCGAN-based architecture

 Training is conditioned on the images from the source domain.

 Conditional GANs provide an effective way to handle many complex domains without worrying about designing structured loss functions explicitly.



images that fool **D**

D tries to identify the fakes

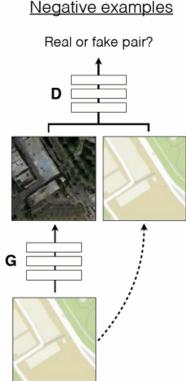


Figure 2 in the original paper.

Text-to-Image Synthesis

Motivation

Given a text description, generate images closely associated.

Uses a conditional GAN with the generator and discriminator being condition on "dense" text embedding.

this small bird has a pink breast and crown, and black primaries and secondaries.



the flower has petals that are bright pinkish purple with white stigma



this magnificent fellow is almost all black with a red crest, and white cheek patch.



this white and yellow flower have thin white petals and a round yellow stamen



Figure 1 in the original paper.

Text-to-Image Synthesis

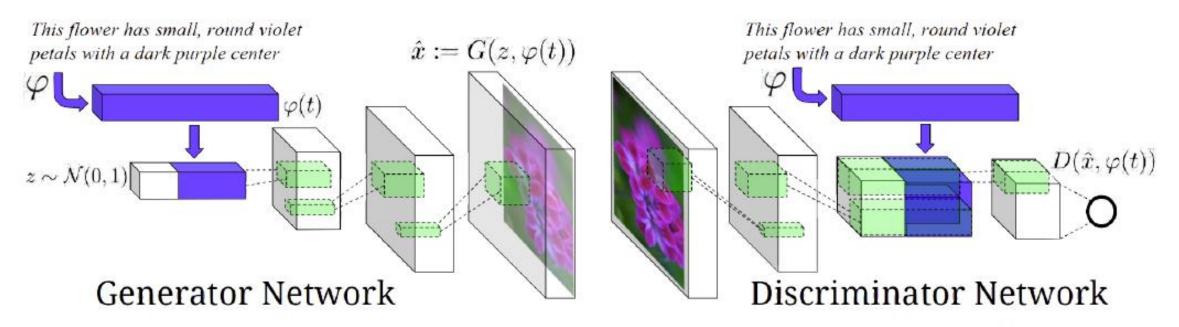


Figure 2 in the original paper.

Positive Example:

Real Image, Right Text

Negative Examples:

Real Image, Wrong Text Fake Image, Right Text

Face Aging with Conditional GANs

 Differentiating Feature: Uses an Identity Preservation Optimization using an auxiliary network to get a better approximation of the latent code (z*) for an input image.

Latent code is then conditioned on a discrete (one-hot) embedding of age

categories. Latent Vector Approximation Face Aging Input face x of age y_0 Identity Encoder Resulting face x_{target} Preserving of age "60+" Optimization | Optimized reconstruction Initial reconstruction Generator $\overline{x_0}$ of age y_0 \overline{x} of age y_0 "60+" Generator Generator y_0

Figure 1 in the original paper.

Face Aging with Conditional GANs

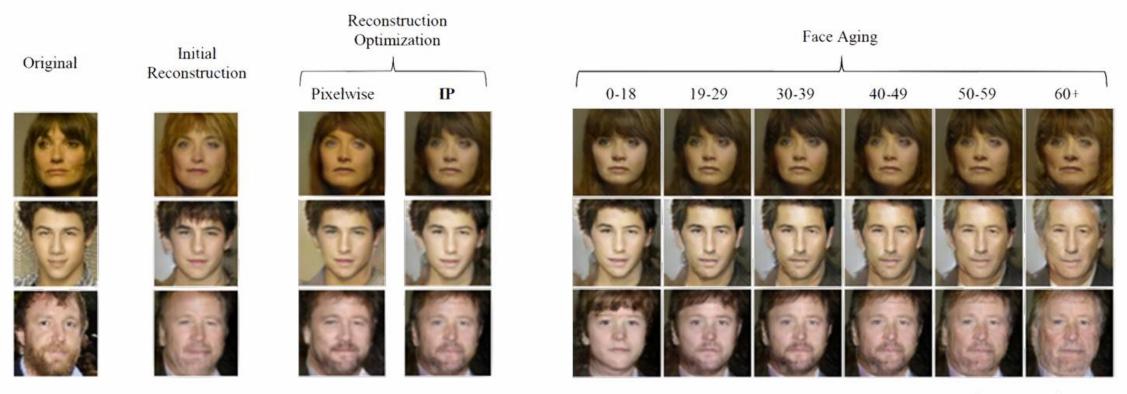
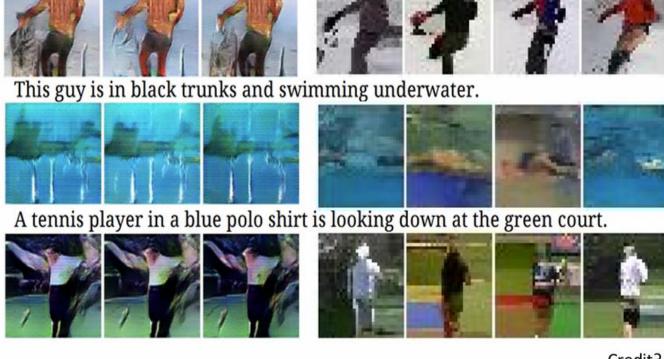


Figure 3 in the original paper.

Conditional GANs

Conditional Model Collapse

- Scenario observed when the Conditional GAN starts ignoring either the code (c) or the noise variables (z).
- This limits the diversity of images generated.



A man in a orange jacket with sunglasses and a hat ski down a hill.

Credit?

Additional Sources:

There are lots of excellent references on GANs:

- Sebastian Nowozins presentation at MLSS 2018: https://github.com/nowozin/mlss2018-madrid-gan
- NIPS 2016 tutorial on GANs by Ian Goodfellow: <u>https://arxiv.org/abs/1701.00160</u>
- A nice explanation of Wasserstein GANs by Alex Irpan: https://www.alexirpan.com/2017/02/22/wasserstein-gan.html