

Session 6: Deep Generative Models: VAEs and GANs

Matthew Willetts and Alexander Camuto

Outline

Unsupervised learning

Outline

Unsupervised learning

Autoencoders

Outline

Unsupervised learning

Autoencoders

Generative Adversarial Networks

Unsupervised learning

Unsupervised learning

Generic goal of unsupervised learning is to **find underlying structure** in data.

Specific goals include:

- clustering: group similar observations together;
- reducing the dimensionality for visualization;
- building a better representation of data for a downstream supervised task;
- learning a likelihood function, e.g. to detect anomalies;
- generating new samples similar to past observations.

Unsupervised learning

For complex data (text, image, sound, ...), there is plenty of hidden latent structure we hope to capture:

- **Image data**: find low dimensional semantic representations, independent sources of variation;
- **Text data**: find fixed size, dense semantic representation of data.

Unsupervised learning

For complex data (text, image, sound, ...), there is plenty of hidden latent structure we hope to capture:

- **Image data**: find low dimensional semantic representations, independent sources of variation;
- **Text data**: find fixed size, dense semantic representation of data.

Latent space might be used to help build more efficient human labeling interfaces.

=> Goal: reduce labeling cost via active learning.

Graal of unsupervised learning

A low dimension space which captures all the **variations** of data and **disentangles** the different latent factors underlying the data.

0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	7
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9

(a) Varying c_1 on InfoGAN (Digit type)

7	7	7	7	7	7	7	7	7	7
0	0	0	0	0	0	0	0	0	0
7	7	7	7	7	7	7	7	7	7
9	9	9	9	9	9	9	9	9	9
8	8	8	8	8	8	8	8	8	8

(b) Varying c_1 on regular GAN (No clear meaning)

1	1	1	1	1	1	1	1	1	1
8	8	8	8	8	8	8	8	8	8
3	3	3	3	3	3	3	3	3	3
9	9	9	9	9	9	9	9	9	9
5	5	5	5	5	5	5	5	5	5

(c) Varying c_2 from -2 to 2 on InfoGAN (Rotation)

1	1	1	1	1	1	1	1	1	1
8	8	8	8	8	8	8	8	8	8
3	3	3	3	3	3	3	3	3	3
9	9	9	9	9	9	9	9	9	9
5	5	5	5	5	5	5	5	5	5

(d) Varying c_3 from -2 to 2 on InfoGAN (Width)

Chen, Xi, et al. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. NIPS, 2016.

Self-supervised learning

find smart ways to **build supervision** without labels, exploiting domain knowledge and regularities

Self-supervised learning

find smart ways to **build supervision** without labels, exploiting domain knowledge and regularities

Use **text structure** to create supervision

- Word2Vec, Skip-thought vectors, language models

Self-supervised learning

find smart ways to **build supervision** without labels, exploiting domain knowledge and regularities

Use **text structure** to create supervision

- Word2Vec, Skip-thought vectors, language models

Can we do the same for other domains?

- **Image:** exploit spatial context of an object
- **Sound, video:** exploit temporal context

Self-supervised learning

find smart ways to **build supervision** without labels, exploiting domain knowledge and regularities

Use **text structure** to create supervision

- Word2Vec, Skip-thought vectors, language models

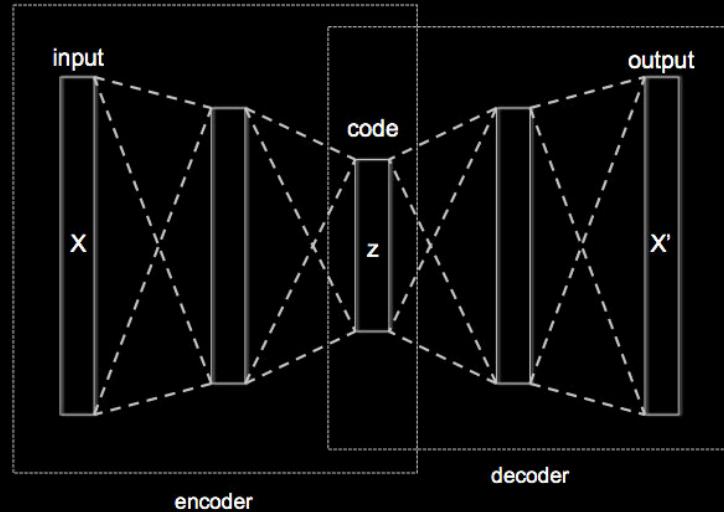
Can we do the same for other domains?

- **Image:** exploit spatial context of an object
- **Sound, video:** exploit temporal context

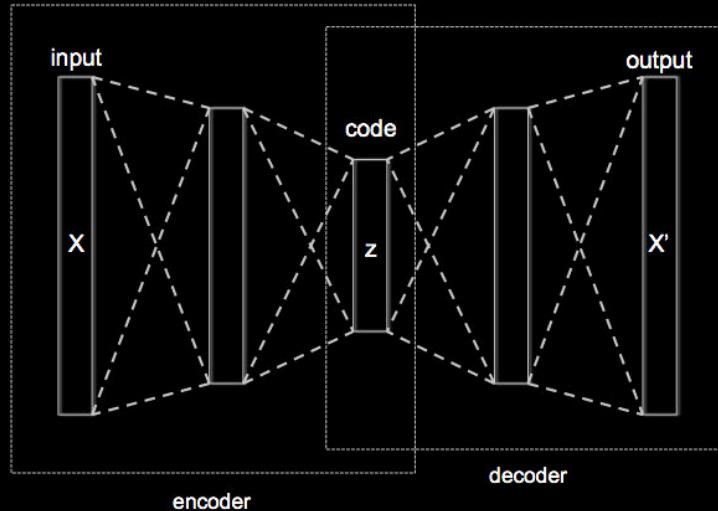
No direct **accuracy** measure: usually tested through a downstream task

Autoencoders

Autoencoder



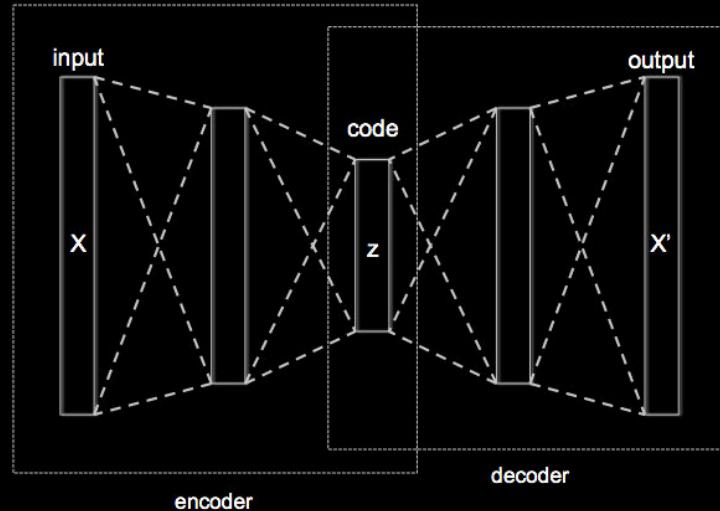
Autoencoder



Supervision : reconstruction loss of the input, usually:

$$l(x, f(x)) = \|f(x) - x\|_2^2$$

Autoencoder

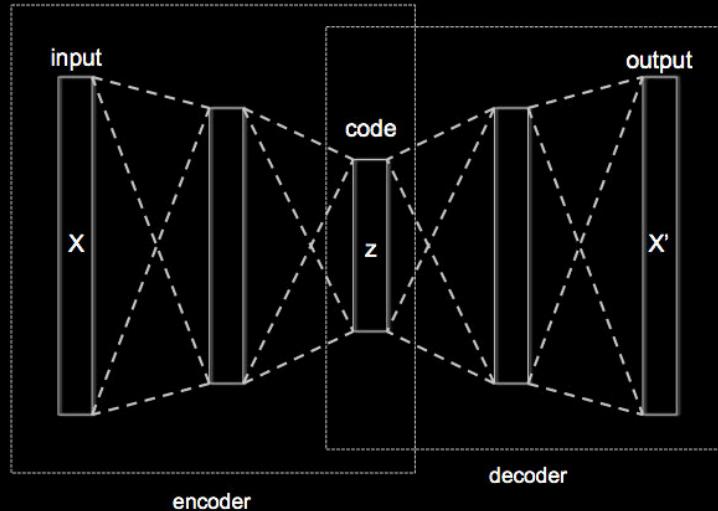


Supervision : reconstruction loss of the input, usually:

$$l(x, f(x)) = \|f(x) - x\|_2^2$$

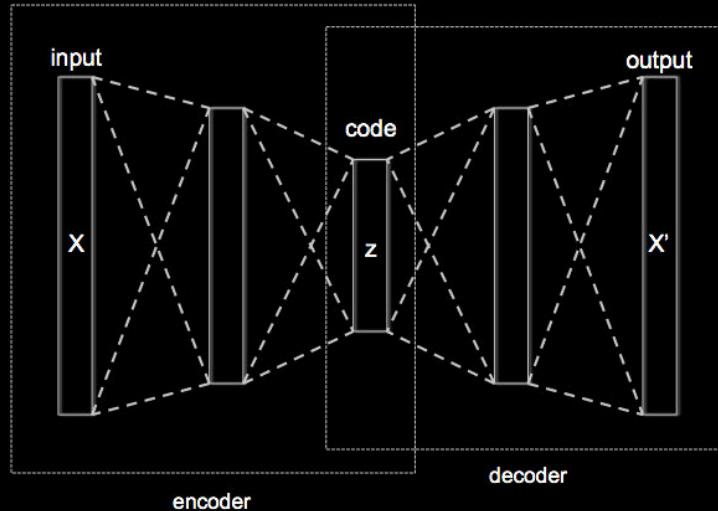
Binary crossentropy is also used

Autoencoder



Keeping the **latent code z** low-dimensional forces the network to learn a "smart" compression of the data, not just an identity function

Autoencoder



Keeping the **latent code z** low-dimensional forces the network to learn a "smart" compression of the data, not just an identity function

Encoder and decoder can have arbitrary architecture (CNNs, RNNs...)

Sparse/Denoising Autoencoder

Adding a sparsity constraint on activations:

$$\|encoder(x)\|_1 \sim \rho, \rho = 0.05$$

Learns sparse features, easily interpretable

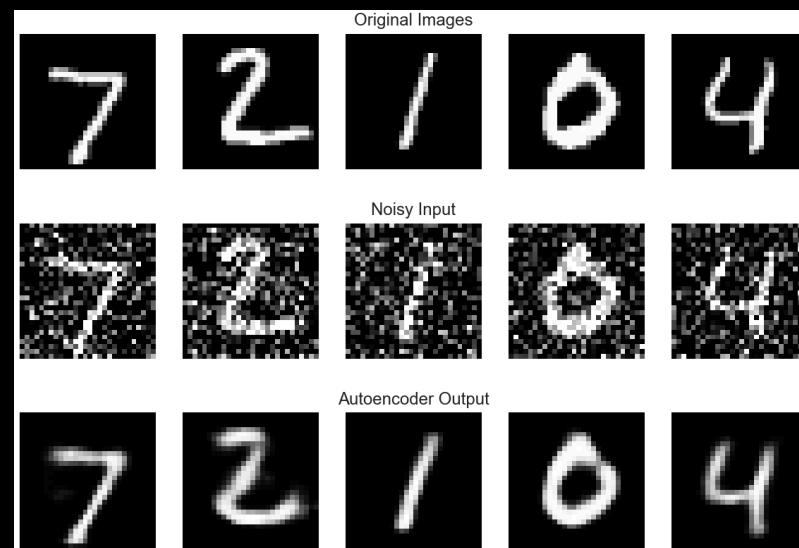
Sparse/Denoising Autoencoder

Adding a sparsity constraint on activations:

$$\|encoder(x)\|_1 \sim \rho, \rho = 0.05$$

Learns sparse features, easily interpretable

Denoising Autoencoder: train features for robustness to noise.



Uses and limitations

After **pre-training** use the latent code **z** as input to a classifier instead of **x**

Semi-supervised learning simultaneous learning of the latent code (on a large, unlabeled dataset) and the classifier (on a smaller, labeled dataset)

Uses and limitations

After **pre-training** use the latent code \mathbf{z} as input to a classifier instead of \mathbf{x}

Semi-supervised learning simultaneous learning of the latent code (on a large, unlabeled dataset) and the classifier (on a smaller, labeled dataset)

Other use: Use decoder $D(\mathbf{x})$ as a **Generative model**: generate samples from random noise

Uses and limitations

After **pre-training** use the latent code \mathbf{z} as input to a classifier instead of \mathbf{x}

Semi-supervised learning simultaneous learning of the latent code (on a large, unlabeled dataset) and the classifier (on a smaller, labeled dataset)

Other use: Use decoder $D(\mathbf{x})$ as a **Generative model**: generate samples from random noise

Limitations :

- Direct autoencoder fails to capture good representations for complex data such as images
- The generative model is usually of very poor quality (very blurry for images for instance)

Reality Check

For image features, **ImageNet pretraining** is still **much better** than unsupervised models

Pretraining Method	Supervision	Pretraining time	Classification	Detection	Segmentation
ImageNet [26]	1000 class labels	3 days	78.2%	56.8%	48.0%
Random Gaussian	initialization	< 1 minute	53.3%	43.4%	19.8%
Autoencoder	-	14 hours	53.8%	41.9%	25.2%
Agrawal <i>et al.</i> [1]	egomotion	10 hours	52.9%	41.8%	-
Doersh <i>et al.</i> [7]	context	4 weeks	55.3%	46.6%	-
Wang <i>et al.</i> [39]	motion	1 week	58.4%	44.0%	-
Ours	context	14 hours	56.5%	44.5%	29.7%

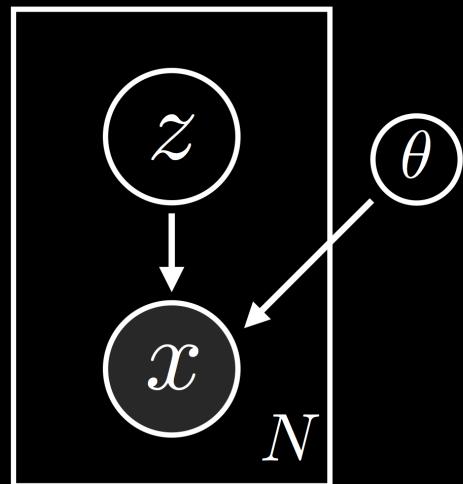
- Results shown after fine-tuned the network on **Pascal VOC dataset**
- The "ours" method is feature representation based on context inpainting

Variational Autoencoders

Variational Autoencoders (VAE).

Assume the data samples $\mathbf{x}^{(i)}$ are generated by the model:

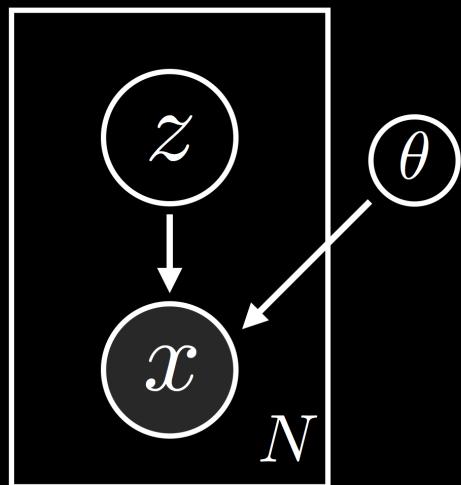
$$p_{\theta^*}(\mathbf{x}, \mathbf{z}) = p_{\theta^*}(\mathbf{z}) \cdot p_{\theta^*}(\mathbf{x}|\mathbf{z})$$



Variational Autoencoders (VAE).

Assume the data samples $\mathbf{x}^{(i)}$ are generated by the model:

$$p_{\theta^*}(\mathbf{x}, \mathbf{z}) = p_{\theta^*}(\mathbf{z}) \cdot p_{\theta^*}(\mathbf{x}|\mathbf{z})$$



- \mathbf{x} is an observed r.v. with values in \mathbb{R}^n ;
- \mathbf{z} is a latent r.v. with values in \mathbb{R}^d ;

Variational Autoencoders (VAE).

Variational Autoencoders (VAE).

- True continuous parameters θ^* are unknown;
- Estimate parameters θ from data $\mathbf{x}^{(i)}$ by maximizing the marginal likelihood (MLE):

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{z}) \cdot p_{\theta}(\mathbf{x}|\mathbf{z}) \, d\mathbf{z}$$

Variational Autoencoders (VAE).

- True continuous parameters θ^* are unknown;
- Estimate parameters θ from data $\mathbf{x}^{(i)}$ by maximizing the marginal likelihood (MLE):

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{z}) \cdot p_{\theta}(\mathbf{x}|\mathbf{z}) \, d\mathbf{z}$$

But this high dimensional integral cannot be estimated efficiently.

Variational Autoencoders (VAE).

- True continuous parameters θ^* are unknown;
- Estimate parameters θ from data $\mathbf{x}^{(i)}$ by maximizing the marginal likelihood (MLE):

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{z}) \cdot p_{\theta}(\mathbf{x}|\mathbf{z}) \, d\mathbf{z}$$

But this high dimensional integral cannot be estimated efficiently.

Variational Autoencoders (VAE) - Imp.

Imagine we try to estimate $p_\theta(\mathbf{x})$ by marginalising out \mathbf{z} :

$$p_\theta(\mathbf{x}) = \int p_\theta(\mathbf{z}) \cdot p_\theta(\mathbf{x}|\mathbf{z}) \, d\mathbf{z} = \mathbb{E}_{p_\theta(\mathbf{z})} p_\theta(\mathbf{x}|\mathbf{z})$$

But instead introduce $\frac{q_\phi(\mathbf{z}|\mathbf{x})}{q_\phi(\mathbf{z})}$ into the integral:

$$p_\theta(\mathbf{x}) = \int p_\theta(\mathbf{z}) \cdot p_\theta(\mathbf{x}|\mathbf{z}) \cdot \frac{q_\phi(\mathbf{z}|\mathbf{x})}{q_\phi(\mathbf{z})} \, d\mathbf{z} = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z})}$$

This is importance sampling - by picking a q distribution with continuous parameters ϕ that is well-peaked around high values of the joint distribution we can get an estimate of $p_\theta(\mathbf{x})$ with fewer samples.

Variational Autoencoders (VAE) - Imp.

$$p_\theta(\mathbf{x}) = \int p_\theta(\mathbf{z}) \cdot p_\theta(\mathbf{x}|\mathbf{z}) \cdot \frac{q_\phi(\mathbf{z}|\mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} d\mathbf{z} = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})}$$

Can we bound this? Yes, apply Jensen's Inequality:

$$\mathbb{E} f(x) \geq f(\mathbb{E} x)$$

where $f(\cdot) = -\log(\cdot)$

$$-\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \log \frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \geq -\log \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} = -\log p_\theta(\mathbf{x})$$

$$\log p_\theta(\mathbf{x}) \geq \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \log \frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})}$$

Variational Autoencoders (VAE) - Imp.

Expand this out

$$\log p_\theta(\mathbf{x}) \geq \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \log \frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \log p_\theta(\mathbf{x}|\mathbf{z}) + \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \log \frac{p_\theta(\mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})}$$

$$\log p_\theta(\mathbf{x}) \geq \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \log p_\theta(\mathbf{x}|\mathbf{z}) - D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}) || p_\theta(\mathbf{z}))$$

$$\log p_\theta(\mathbf{x}) \geq \mathcal{L}(x)$$

$$\mathcal{L}(x) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \log p_\theta(\mathbf{x}|\mathbf{z}) - D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}) || p_\theta(\mathbf{z}))$$

We call $\mathcal{L}(x)$ the Evidence Lower Bound, because it is. If we now take gradients wrt θ, ϕ we will be increasing the lower bound of the evidence of our model.

First though, another view:

Variational Autoencoders (VAE) - KL

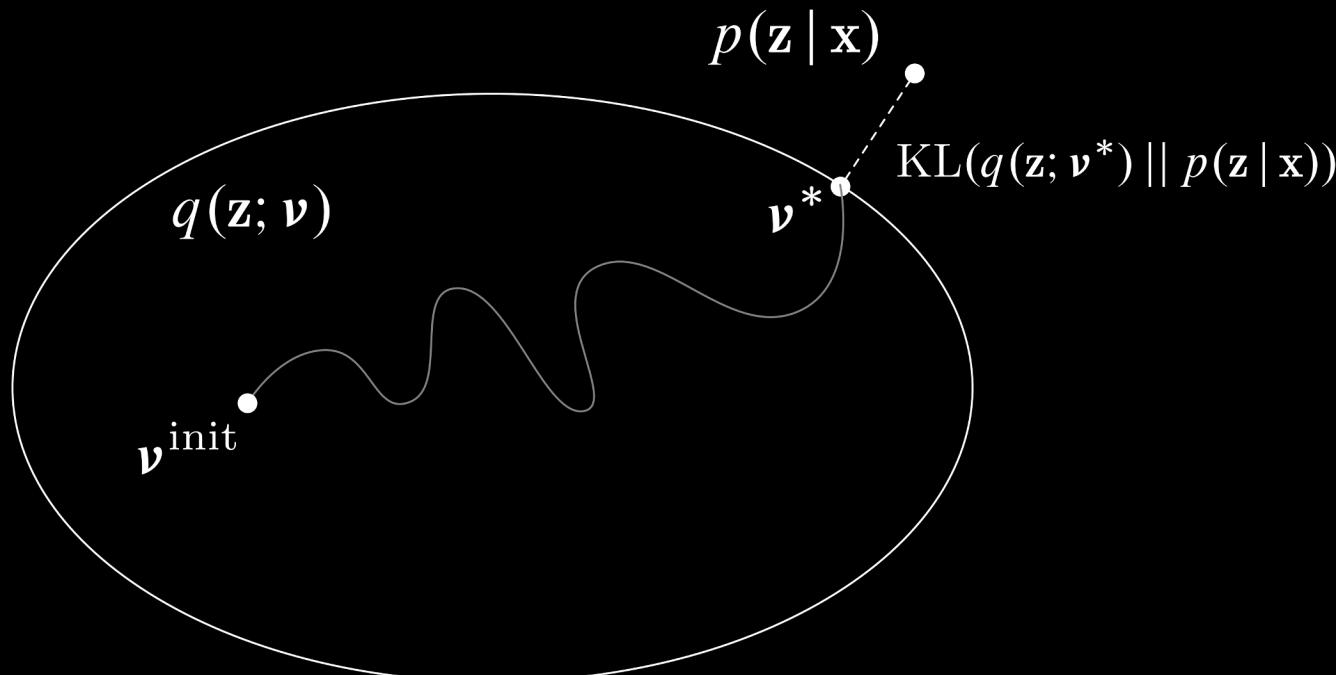
We wish we had the true posterior $p_\theta(\mathbf{z}|\mathbf{x})$ but that's really hard. We'd have to divide the joint by the marginal, we don't have that.

Instead, introduce a family of approximate (variational) posteriors on the latent variable: $q_\phi(\mathbf{z}|\mathbf{x})$ with continuous parameters ϕ :

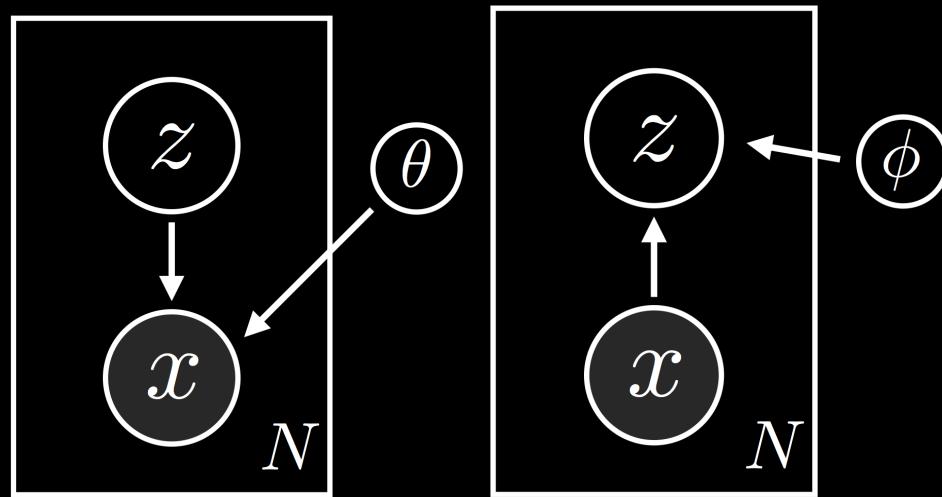
Variational Autoencoders (VAE) - KL

We wish we had the true posterior $p_\theta(\mathbf{z}|\mathbf{x})$ but that's really hard. We'd have to divide the joint by the marginal, we don't have that.

Instead, introduce a family of approximate (variational) posteriors on the latent variable: $q_\phi(\mathbf{z}|\mathbf{x})$ with continuous parameters ϕ :



Variational Autoencoders (VAE) - KL



Variational Autoencoders (VAE) - KL

We want to find the nearest member of this family $q_\phi(\mathbf{z}|\mathbf{x})$ to the true posterior in KL terms. We want to minimise:

$$\begin{aligned}\text{D}_{KL}[q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x})] &= - \int q_\phi(\mathbf{z}|\mathbf{x}) \log \frac{p_\theta(\mathbf{z}|\mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} d\mathbf{z} \\ &= - \int q_\phi(\mathbf{z}|\mathbf{x}) \log \frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} d\mathbf{z} + \int q_\phi(\mathbf{z}|\mathbf{x}) \log \frac{p_\theta(\mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} d\mathbf{z} \\ &= -\mathcal{L}(x) + p_\theta(\mathbf{x})\end{aligned}$$

Now we have found the gap from applying Jensen's inequality - rearrange:

$$p_\theta(\mathbf{x}) = \mathcal{L}(x) + \text{D}_{KL}[q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x})]$$

LHS: no dep. on $q_\phi(\mathbf{z}|\mathbf{x})$, so $\max \mathcal{L}(x) \rightarrow \min \text{D}_{KL}[q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x})]$ wrt ϕ .

Variational Autoencoders (VAE).

- Assume prior distribution for latent variable \mathbf{z} :

$$p_{\theta}(\mathbf{z}) = \mathcal{N}(0, 1)$$

Variational Autoencoders (VAE).

- Assume prior distribution for latent variable \mathbf{z} :

$$p_\theta(\mathbf{z}) = \mathcal{N}(0, 1)$$

- Parametrize $p_\theta(\mathbf{x}|\mathbf{z})$ by a neural network f_θ (decoder):

$$p_\theta(\mathbf{x}^{(i)}|\mathbf{z}) = \mathcal{N}(f_\theta(\mathbf{z}), 1)$$

Variational Autoencoders (VAE).

- Assume prior distribution for latent variable \mathbf{z} :

$$p_\theta(\mathbf{z}) = \mathcal{N}(0, 1)$$

- Parametrize $p_\theta(\mathbf{x}|\mathbf{z})$ by a neural network f_θ (decoder):

$$p_\theta(\mathbf{x}^{(i)}|\mathbf{z}) = \mathcal{N}(f_\theta(\mathbf{z}), 1)$$

- Parametrize $q_\phi(\mathbf{z}|\mathbf{x}^{(i)})$ by a neural network with two heads μ_ϕ and σ_ϕ (encoder):

$$q_\phi(\mathbf{z}|\mathbf{x}^{(i)}) = \mathcal{N}(\mu_\phi(\mathbf{x}^{(i)}), \sigma_\phi(\mathbf{x}^{(i)}))$$

Variational Autoencoders (VAE).

- Assume prior distribution for latent variable \mathbf{z} :

$$p_\theta(\mathbf{z}) = \mathcal{N}(0, 1)$$

- Parametrize $p_\theta(\mathbf{x}|\mathbf{z})$ by a neural network f_θ (decoder):

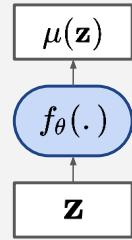
$$p_\theta(\mathbf{x}^{(i)}|\mathbf{z}) = \mathcal{N}(f_\theta(\mathbf{z}), 1)$$

- Parametrize $q_\phi(\mathbf{z}|\mathbf{x}^{(i)})$ by a neural network with two heads μ_ϕ and σ_ϕ (encoder):

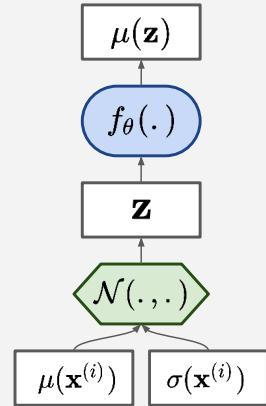
$$q_\phi(\mathbf{z}|\mathbf{x}^{(i)}) = \mathcal{N}(\mu_\phi(\mathbf{x}^{(i)}), \sigma_\phi(\mathbf{x}^{(i)}))$$

- Reparametrization trick:

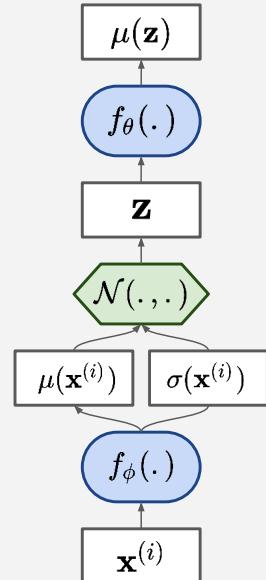
$$\mathbf{z} = \mu_\phi(\mathbf{x}^{(i)}) + \sigma_\phi(\mathbf{x}^{(i)}) \cdot \epsilon \quad \text{with} \quad \epsilon \sim \mathcal{N}(0, 1)$$



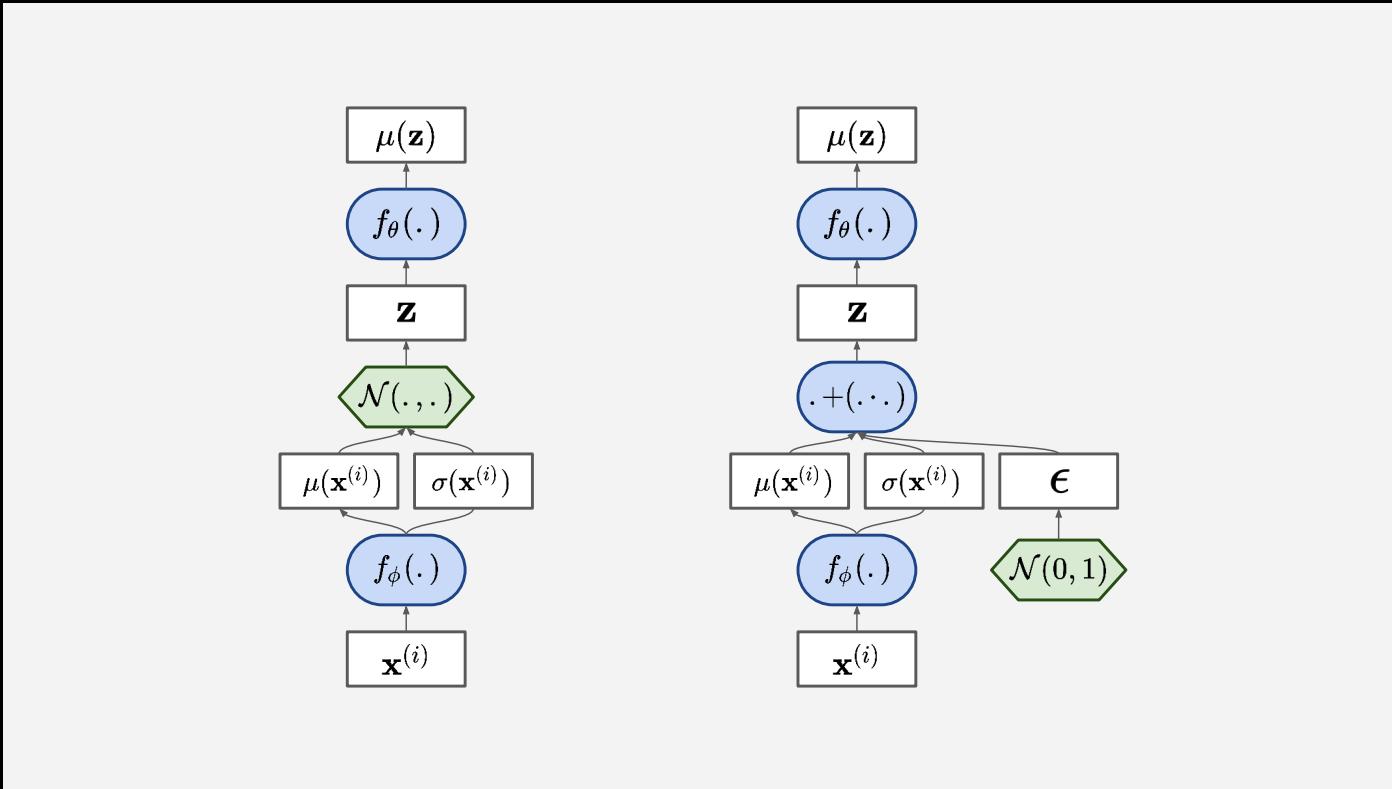
The **decoder** f_θ defines the likelihood of data.



The latent variable **z** is stochastic.



The **encoder** f_ϕ defines an approximate posterior on \mathbf{z} .



The reparametrization makes the objective differentiable wrt. θ & ϕ .

Conv/deconv VAEs



[conv/deconv VAE](#) trained by Alec Radford in 2015 on Labeled Faces in the Wild (LFW) dataset, 2h on single GTX 980

Variational Autoencoders (VAE).

Remarks

- Similar to Denoising AE but noise added to hidden layer;
- Motivated by a well-defined probabilistic model of the generative process;
- Quite easy to train in practice.

Variational Autoencoders (VAE).

Remarks

- Similar to Denoising AE but noise added to hidden layer;
- Motivated by a well-defined probabilistic model of the generative process;
- Quite easy to train in practice.

Limitations

- Is the continuous parametrization of posterior latent distribution too restrictive?
- Would a discrete latent variable make more sense?
- Gaussian parametrization of the decoder output results in blurry images.

Discrete latent variables VAE

Gumbel-Softmax / Concrete distribution VAEs

- Adapts the reparametrization trick for a discrete \mathbf{z} .
- Trains ok but no ground breaking applications so far.

Discrete latent variables VAE

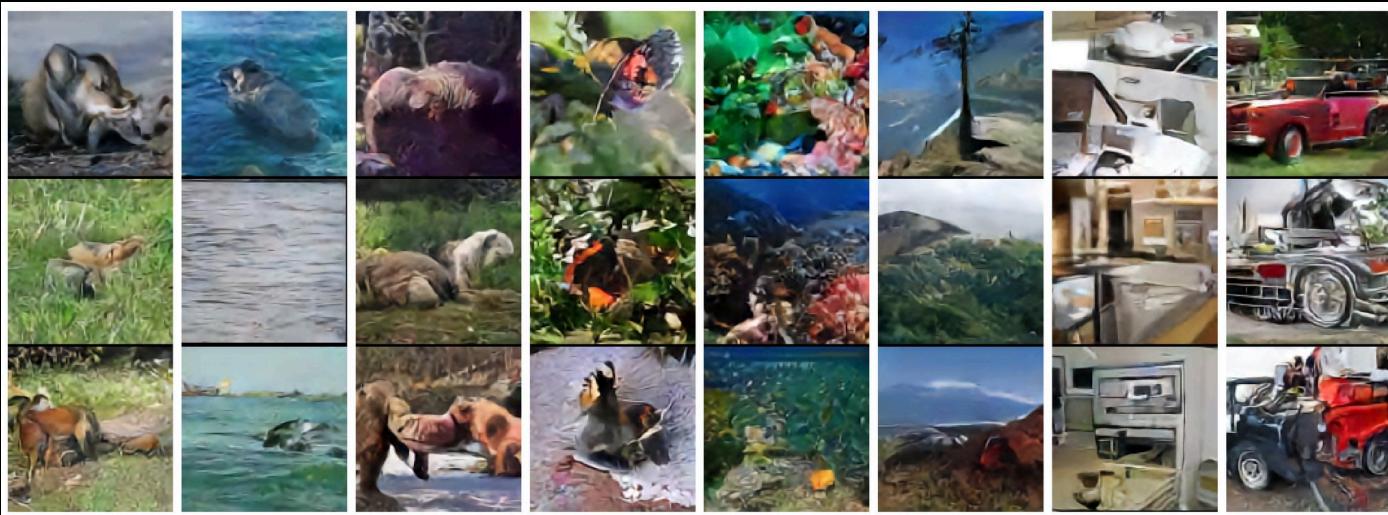
Gumbel-Softmax / Concrete distribution VAEs

- Adapts the reparametrization trick for a discrete \mathbf{z} .
- Trains ok but no ground breaking applications so far.

VQ-VAE

- \mathbf{z} is a vector indexed in a trainable embedding matrix.
- Select \mathbf{z} as embedding vector closest to encoder output.
- Approximate backprop via "gradient-copy" trick.
- Very expressive model, especially when combined with strong decoders and priors.

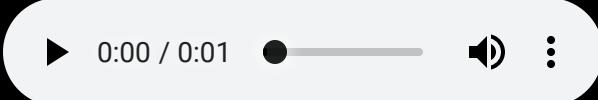
VQ-VAE imagenet results



VQ-VAE speech results

Speech synth demo: <https://avdnoord.github.io/homepage/vqvae/>

Example reconstruction:

- Original: 
- Reconstructed: 

Reconstruction conditionned on different speaker id:

- Original: 
- Reconstructed: 

Disentangled VAEs

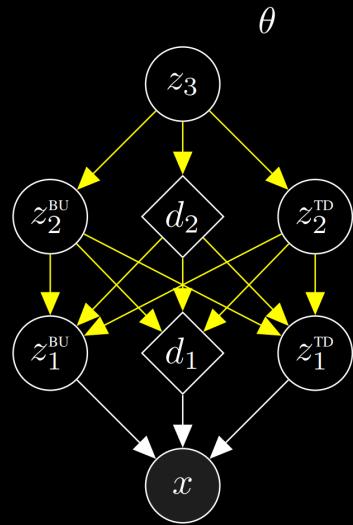


Disentangled VAEs

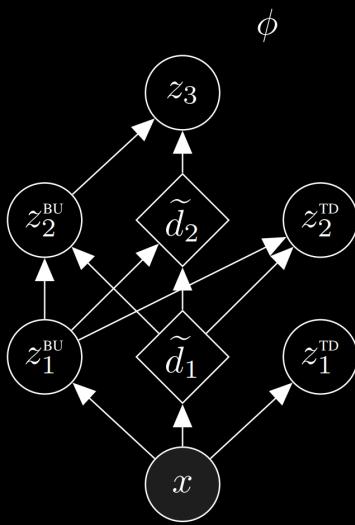


BIVA: Bidirectional-Inference VAE

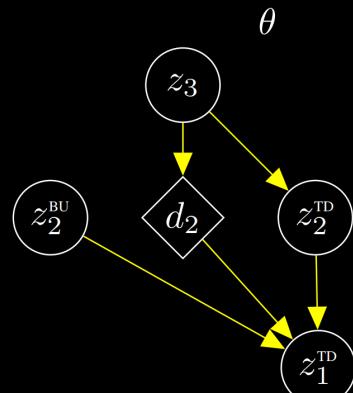
Lots of changes from 'Vanilla' VAE



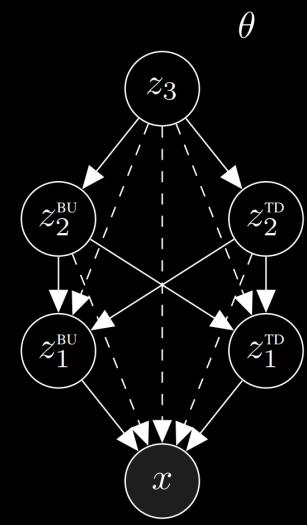
(a) Generative model



(b) BU inference



(c) TD inference



(d) Variable dependency

BIVA: Bidirectional-Inference VAE

Graphically:

- Hierarchy of \mathbf{z} layers:

$$\begin{aligned} q_{\phi}(\mathbf{z}|x) &= q_{\phi}(\mathbf{z}^{\text{TD}}, \mathbf{z}^{\text{BU}}|x) \\ &= q_{\phi}(z_L|x, z_{<L}^{\text{BU}}) \prod_{i=1}^{L-1} q_{\phi}(z_i^{\text{BU}}|x, z_{<i}^{\text{BU}}) q_{\phi,\theta}(z_i^{\text{TD}}|x, z_{\neq i}^{\text{BU}}, z_{>i}^{\text{TD}}) \end{aligned}$$

$$\begin{aligned} p_{\theta}(x, \mathbf{z}) &= p_{\theta}(x|\mathbf{z}) p_{\theta}(z_L) p_{\theta}(\mathbf{z}^{\text{BU}}, \mathbf{z}^{\text{TD}}) \\ &= p_{\theta}(x|\mathbf{z}) p_{\theta}(z_L) \prod_{i=1}^{L-1} p_{\theta}(z_i^{\text{BU}}|z_{>i}) p_{\theta}(z_i^{\text{TD}}|z_{>i}) \end{aligned}$$

$$\text{where } z_i = [z_i^{\text{BU}}, z_i^{\text{TD}}]$$

BIVA: Bidirectional-Inference VAE

Samples:

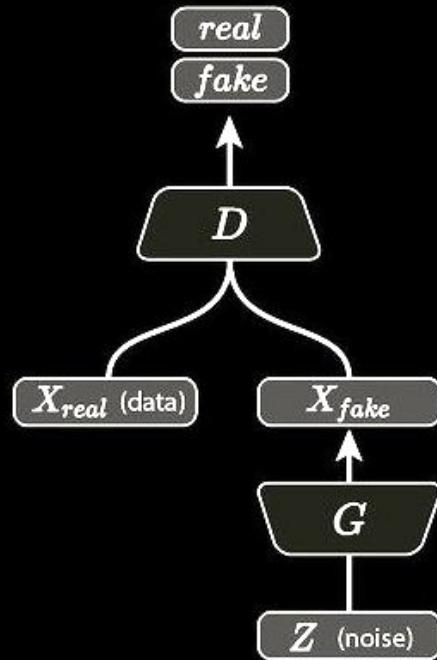


(c) $\sigma^2 = 0.5$

(d) $\sigma^2 = 1.0$

Generative Adversarial Networks

Generative Adversarial Networks



Alternate training of a **generative network** G and a **discriminative network** D

GANs

- D tries to find out which example are generated or real
- G tries to fool D into thinking its generated examples are real

GANs

- D tries to find out which example are generated or real
- G tries to fool D into thinking its generated examples are real

Sample real data $x \sim p_{data}$

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))]$$
$$D(x) \in [0, 1] \rightarrow 1$$

GANs

- D tries to find out which example are generated or real
- G tries to fool D into thinking its generated examples are real

Sample real data $x \sim p_{data}$

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))]$$
$$D(x) \in [0, 1] \rightarrow 1$$

Sample \mathbf{z} and generate fake data $G(\mathbf{z})$

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))]$$
$$D(G(z)) \in [0, 1] \rightarrow 0$$

GANs

- D tries to find out which example are generated or real
- G tries to fool D into thinking its generated examples are real

Sample real data $x \sim p_{data}$

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))]$$
$$D(x) \in [0, 1] \rightarrow 1$$

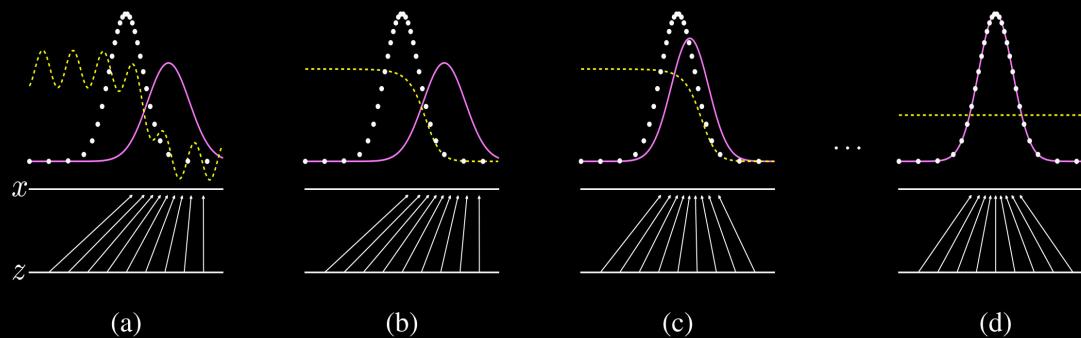
Sample \mathbf{z} and generate fake data $G(\mathbf{z})$

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))]$$
$$D(G(z)) \in [0, 1] \rightarrow 0$$

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))]$$
$$D(G(z)) \in [0, 1] \rightarrow 1$$

GANs

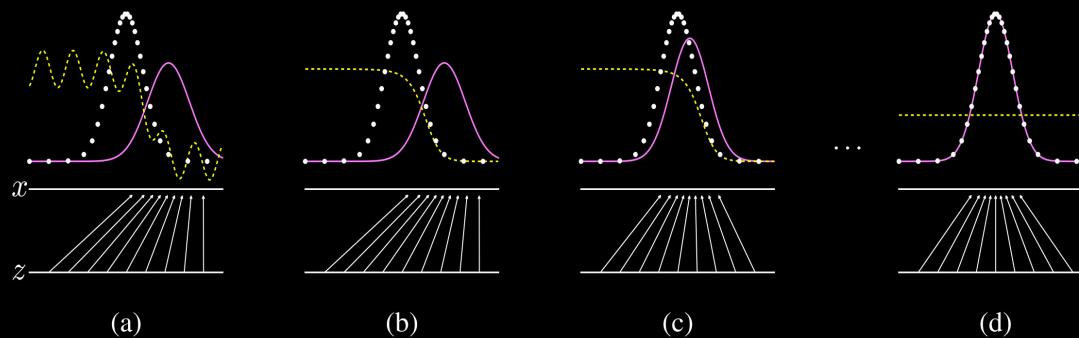
1D-example



- optimal: $D = \frac{1}{2}$, $G(\mathbf{z}) = p_{data}$

GANs

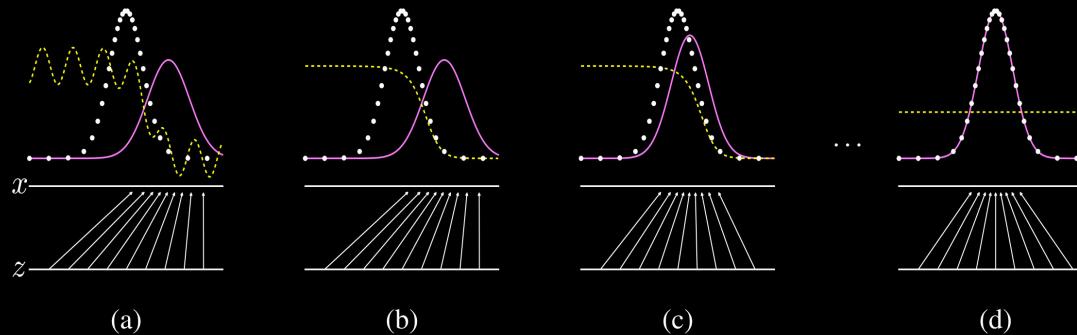
1D-example



- optimal: $D = \frac{1}{2}$, $G(\mathbf{z}) = p_{data}$
- G never "sees" training data, it is solely updated from gradients coming from D

GANs

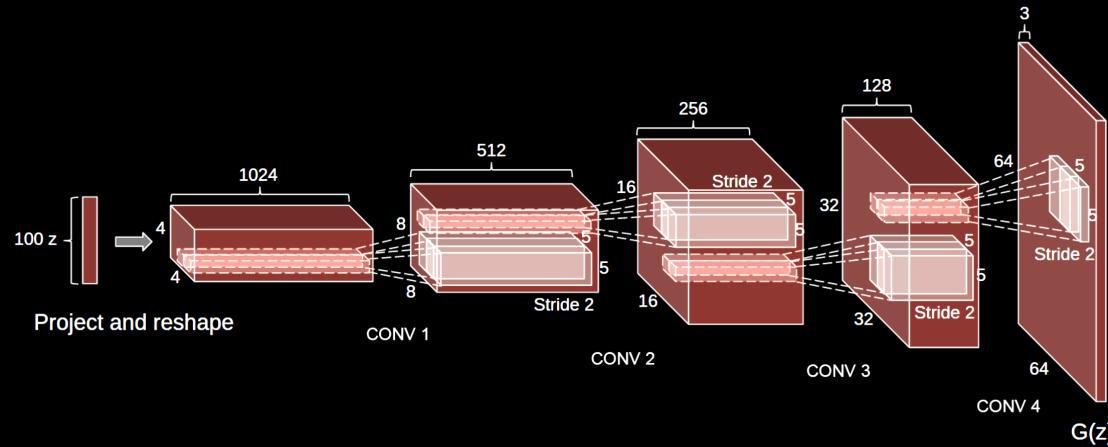
1D-example



- optimal: $D = \frac{1}{2}$, $G(\mathbf{z}) = p_{data}$
- G never "sees" training data, it is solely updated from gradients coming from D
- Naive Keras implementation:

```
d_loss = K.mean(-K.log(Dx) - K.log(1 - DGz))
g_loss = K.mean(K.log(1 - DGz))
```

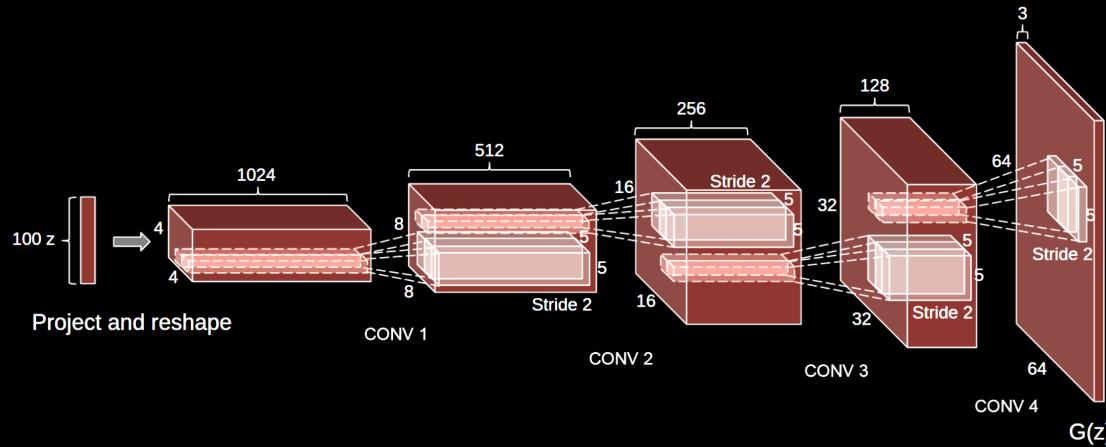
DC-GAN



- GANs training is unstable, and may suffer from **mode collapse**

Radford, Alec, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. 2015.

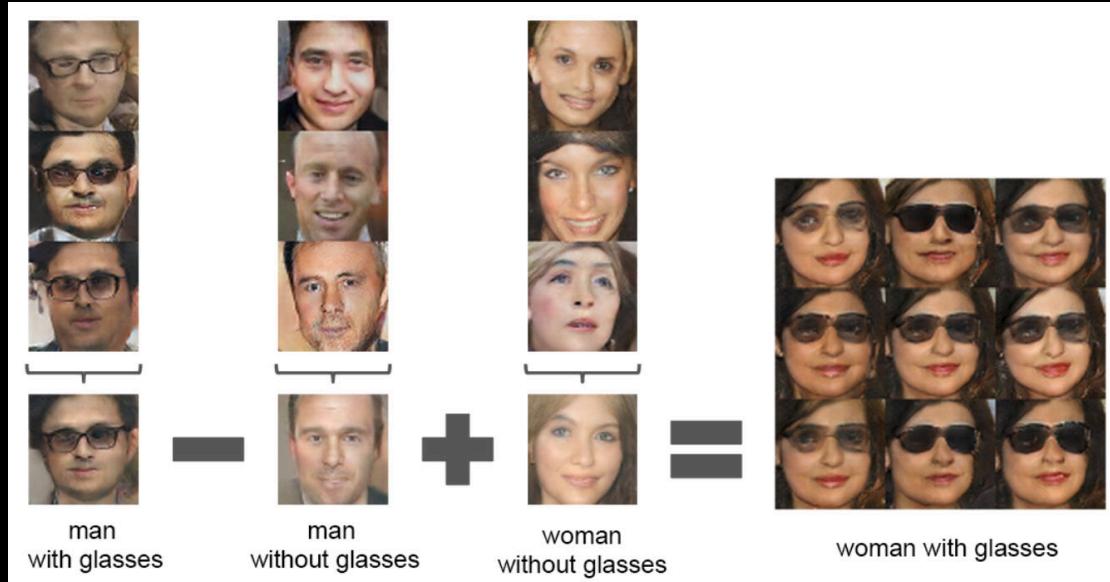
DC-GAN



- GANs training is unstable, and may suffer from **mode collapse**
- **Sensitive hyperparameters:** Use of batchnorm, strided convolutions, careful learning rates, several D updates per G updates...

Radford, Alec, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. 2015.

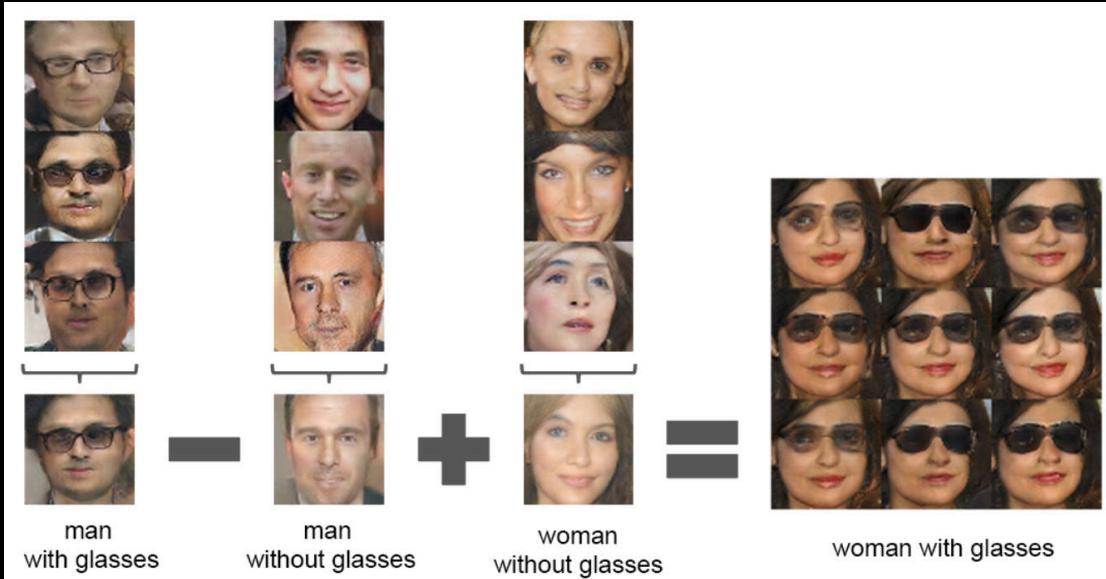
DC-GAN



- Generator generates less-blurry images than VAEs

Radford, Alec, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. 2015.

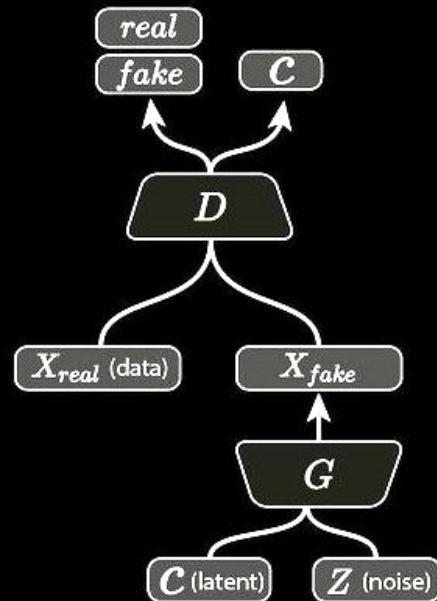
DC-GAN



- Generator generates less-blurry images than VAEs
- Latent space has some local linear properties (vector arithmetic like with Word2Vec)

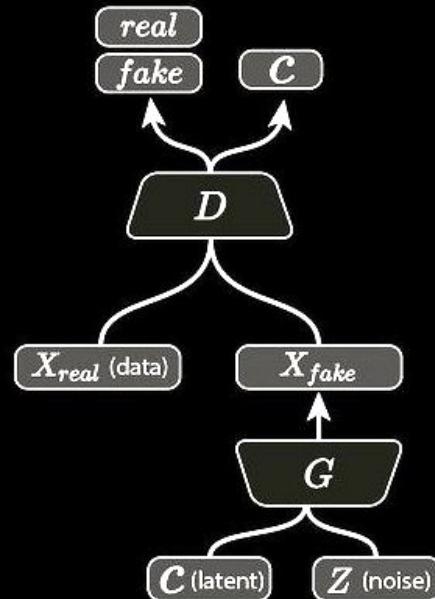
Radford, Alec, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. 2015.

Info GAN



Chen, Xi, et al. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. NIPS, 2016.

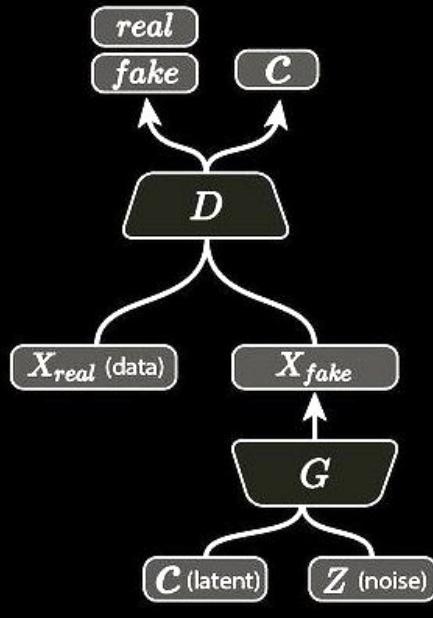
Info GAN



Separate noise **Z** and latent variables **C** and maximizes the mutual information between c the and the observation $G(z, c)$

Chen, Xi, et al. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. NIPS, 2016.

Info GAN



(a) Varying c_1 on InfoGAN (Digit type)

(b) Varying c_1 on regular GAN (No clear meaning)

(c) Varying c_2 from -2 to 2 on InfoGAN (Rotation)

(d) Varying c_3 from -2 to 2 on InfoGAN (Width)

Separate noise \mathbf{Z} and latent variables \mathbf{C} and maximizes the mutual information between \mathbf{c} the and the observation $G(z, c)$

Chen, Xi, et al. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. NIPS, 2016.

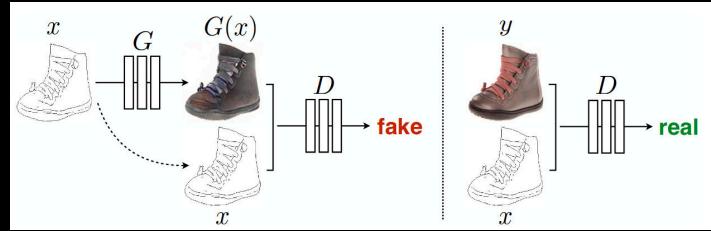
Progressive growing of GANs



All images are generated by walking through the latent space

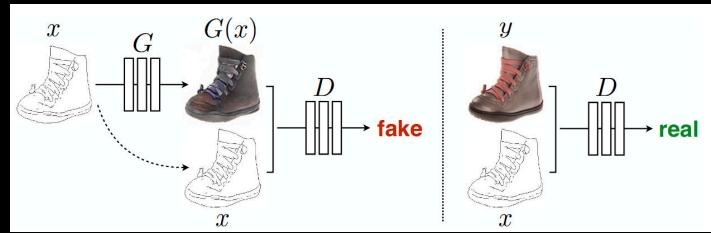
Karras, Tero, et al. Progressive growing of gans for improved quality, stability, and variation. 2017.

Pix2pix: Conditional GANs

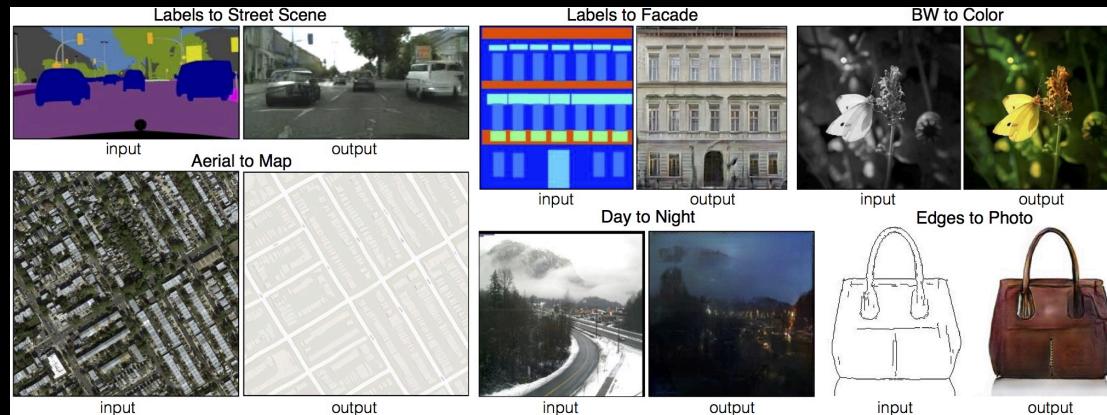


The generation no longer makes use of \mathbf{z} , rather is conditionned by an input \mathbf{x}

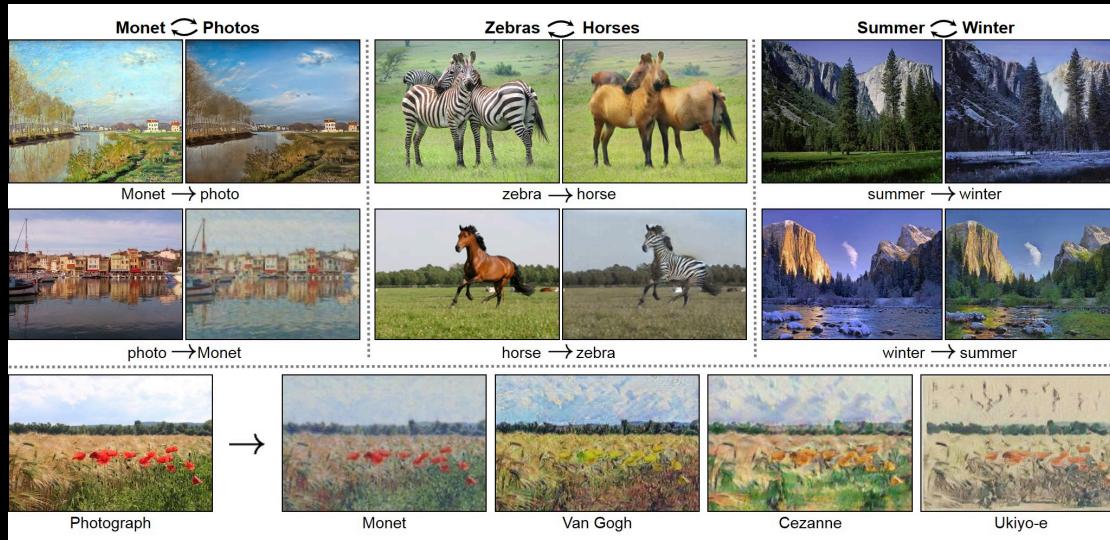
Pix2pix: Conditional GANs



The generation no longer makes use of \mathbf{z} , rather is conditionned by an input \mathbf{x}

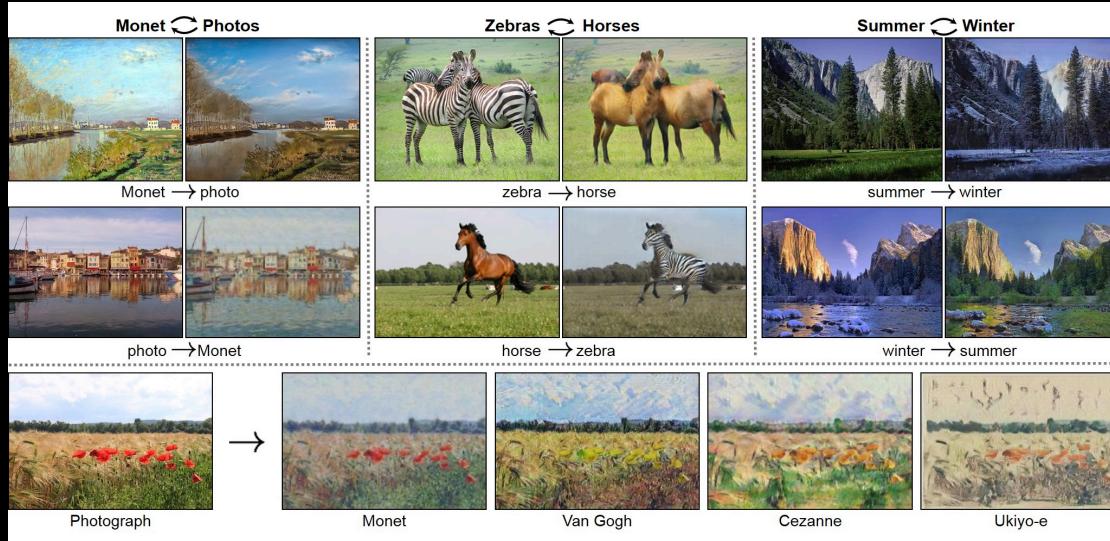


Cycle GANs



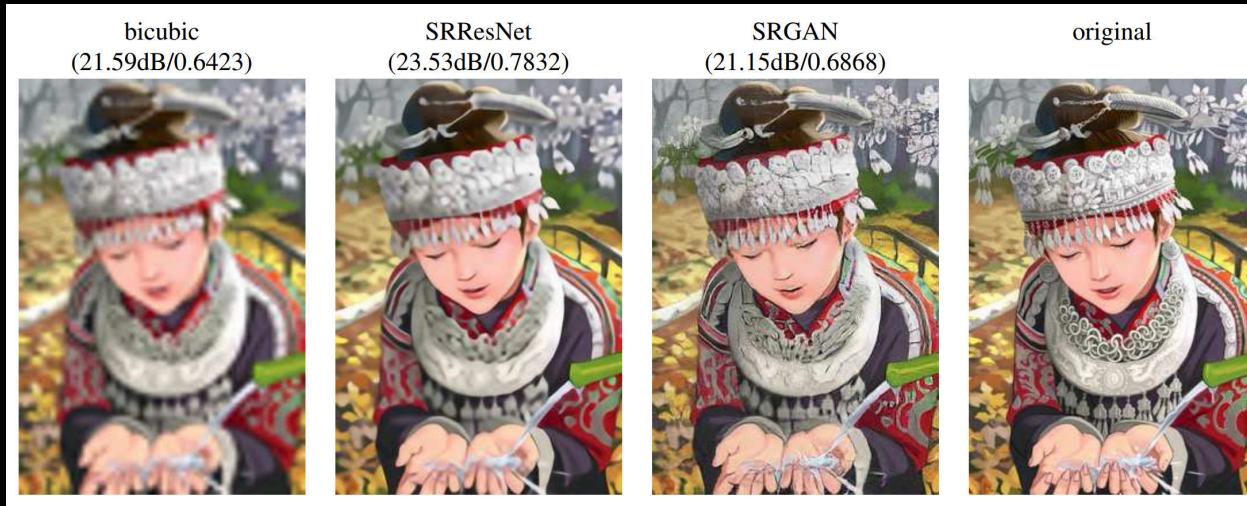
Jun-Yan Zhu et al. Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks, ICCV 2017

Cycle GANs



- No alignment between pairs needed, simply two different sets of images

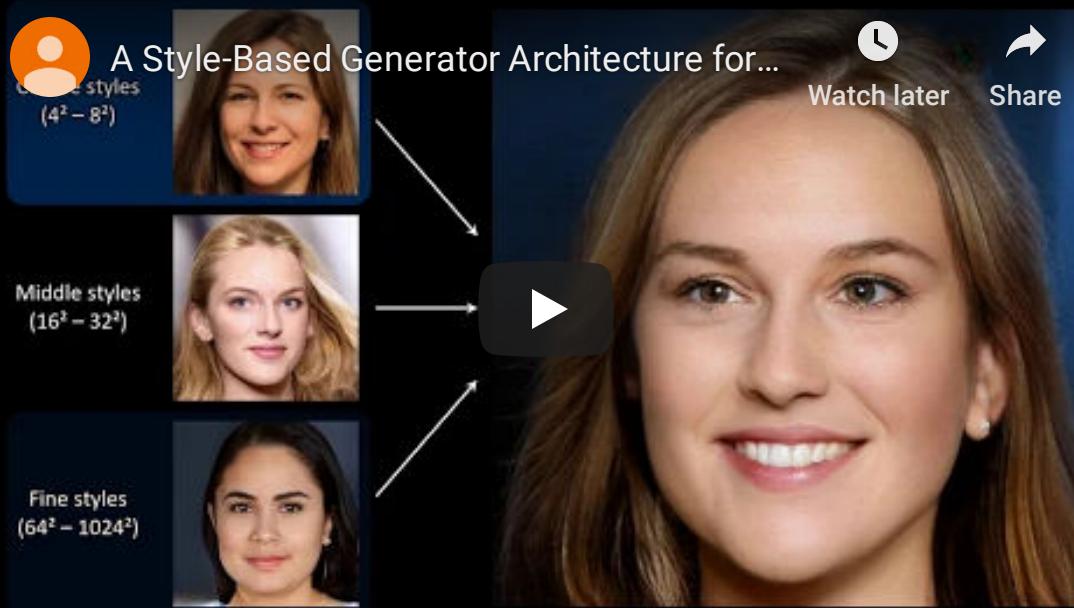
Super Resolution



"Perceptual" loss = combining pixel-wise loss mse-like loss with GAN loss

Ledig, Christian, et al. Photo-realistic single image super-resolution using a generative adversarial network.
CVPR 2016.

Style GANs



[A Style-Based Generator Architecture for Generative Adversarial Networks](#) by Tero Karras, Samuli Laine, Timo Aila, 2018 (preprint)

Takeaways

(Reconstruction) Autoencoders

- have no direct probabilistic interpretation;
- are not designed to generate useful samples;
- encoder defines a useful latent representation.

Takeaways

(Reconstruction) Autoencoders

- have no direct probabilistic interpretation;
- are not designed to generate useful samples;
- encoder defines a useful latent representation.

VAEs

- model explicitly (a lower bound of) the likelihood;
- high quality samples from high dimensional distributions;
- encoder defines a useful latent representation;
- optimization problem is often well-behaved.

Takeaways

GANs

- likelihood-free generative models;
- high quality samples from high dimensional distributions;
- discriminator not meant be used as encoder;
- optimization problem is trickier than for VAEs (open research).

Takeaways

GANs

- likelihood-free generative models;
- high quality samples from high dimensional distributions;
- discriminator not meant be used as encoder;
- optimization problem is trickier than for VAEs (open research).

There exists other kinds of generative models:

- auto-regressive models: PixelCNN, WaveNet, RNN language models...
- can be used as prior and decoder for VAEs, generators for GANs.
- flow-based models: Glow, WaveGlow...

Takeaways

Adversarial training is useful beyond generative models:

- domain adaptation;
- learning representations blind to sensitive attributes;
- defend against malicious inputs (adversarial examples);
- regularization by training on adversarial examples.

Takeaways

Adversarial training is useful beyond generative models:

- domain adaptation;
- learning representations blind to sensitive attributes;
- defend against malicious inputs (adversarial examples);
- regularization by training on adversarial examples.

Quality of samples from VAE and GAN depends a lot on the architectures of sub-networks.

Thanks to

Charles Ollion and Olivier Grisel, Paris-Saclay, for the slides

Chris Holmes and Gil McVean, as always