

Generative models

Florian Buettner

buettner.florian@siemens.com

Today's lecture: NNs as generative models

- Generative models – What and why
- Autoregressive models
- Variational autoencoders
- GANs

Unsupervised learning and generative models

- Just data, no labels
 - PCA, k-means,...
- Learn latent structure of data
- **Density estimation**
 - Training data from distribution p_{data}
 - Learn a distribution p_{model} that is similar to p_{data}

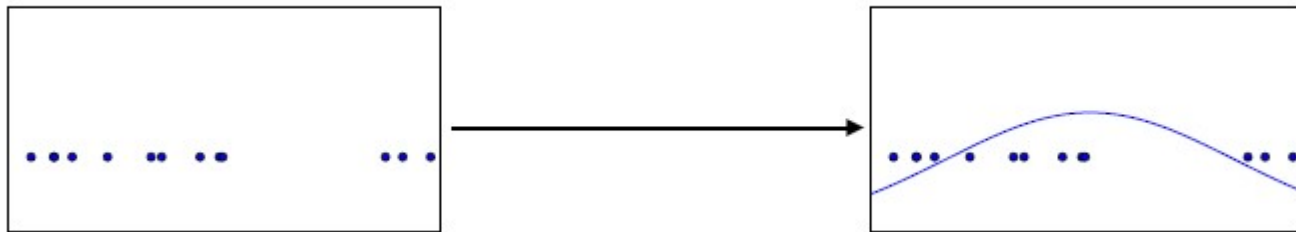


Figure: Goodfellow 2017 *arXiv:1701.00160*

Density modelling

- Simplest approach: learn everything about the data
 - Define explicit model and maximise overall likelihood
- Better: focus on what is useful!
 - pixel value Vs image content, n-gram Vs semantics
 - „not all bits are created equal“
 - Curse of dimensionality
- How can models be used for future tasks?
 - Access representations?
 - Get generative model for free

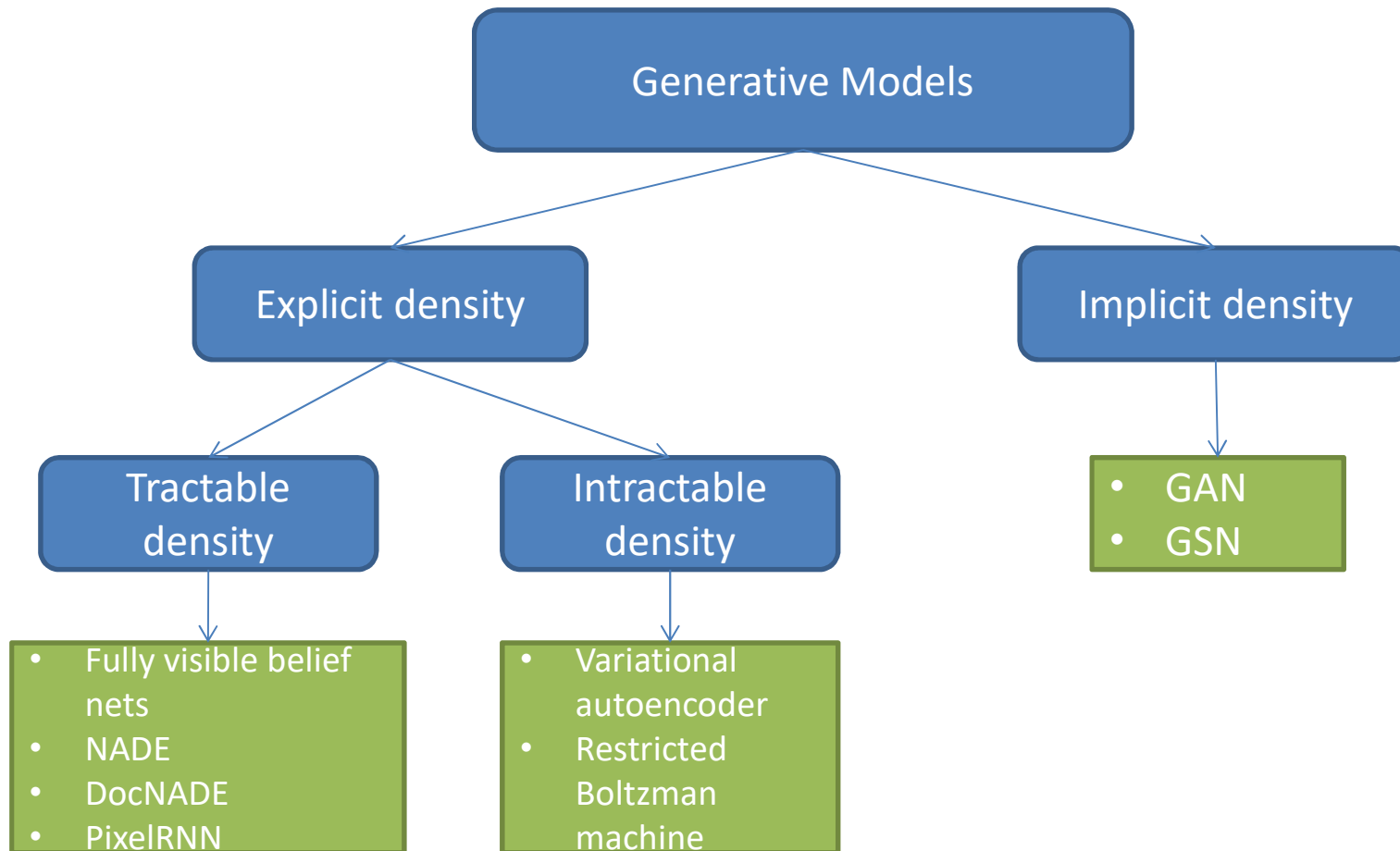
But why?

- Latent variables capturing data manifold as general feature
- Anomaly detection
- Domain transfer: Art, super-resolution, colorisation
- Simulation and planning (RL)
- Creating means understanding
 - What kind of patterns has the model learnt?

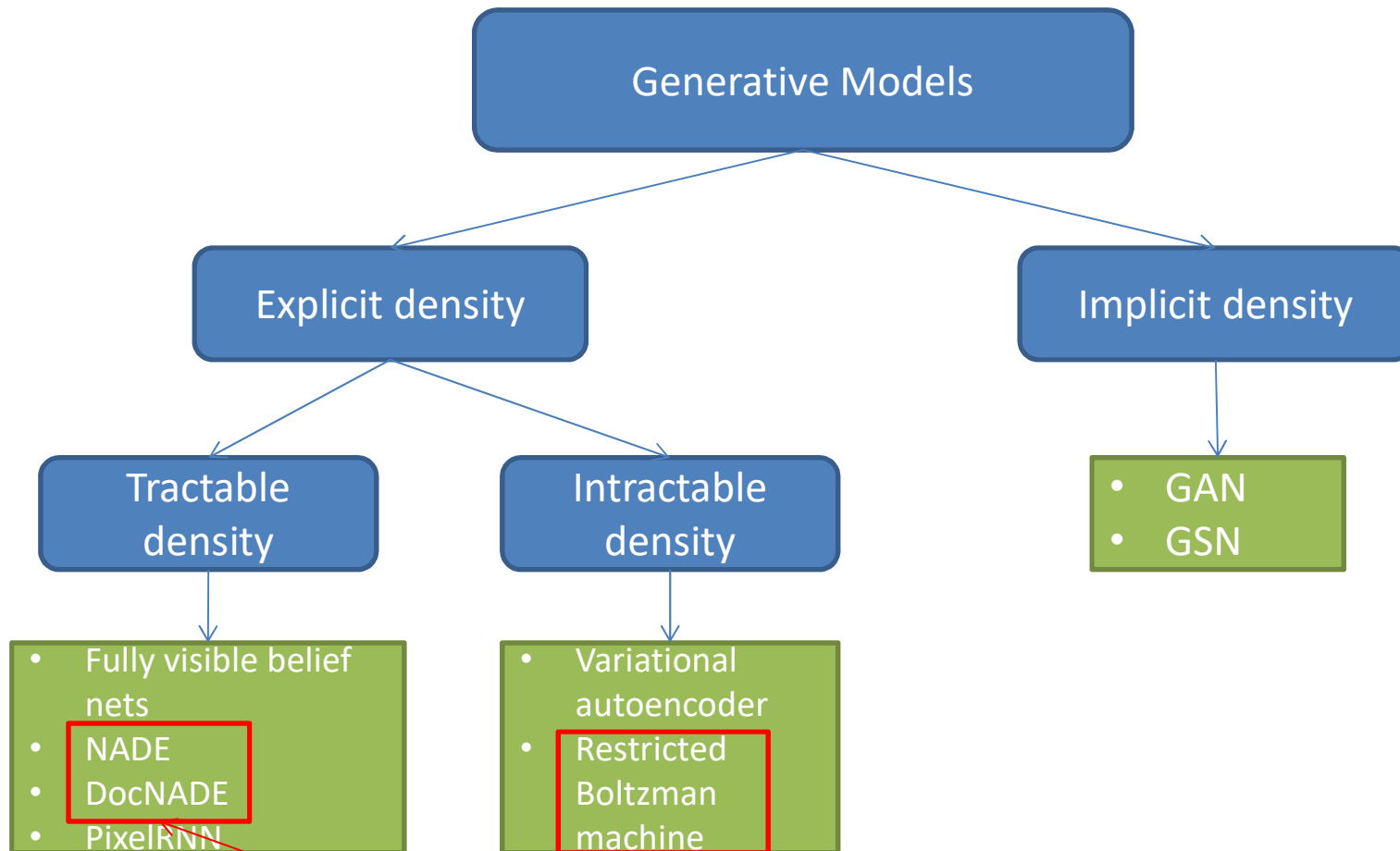


Figure: Zhu et al. *arXiv preprint*, 2017.

Types of generative models



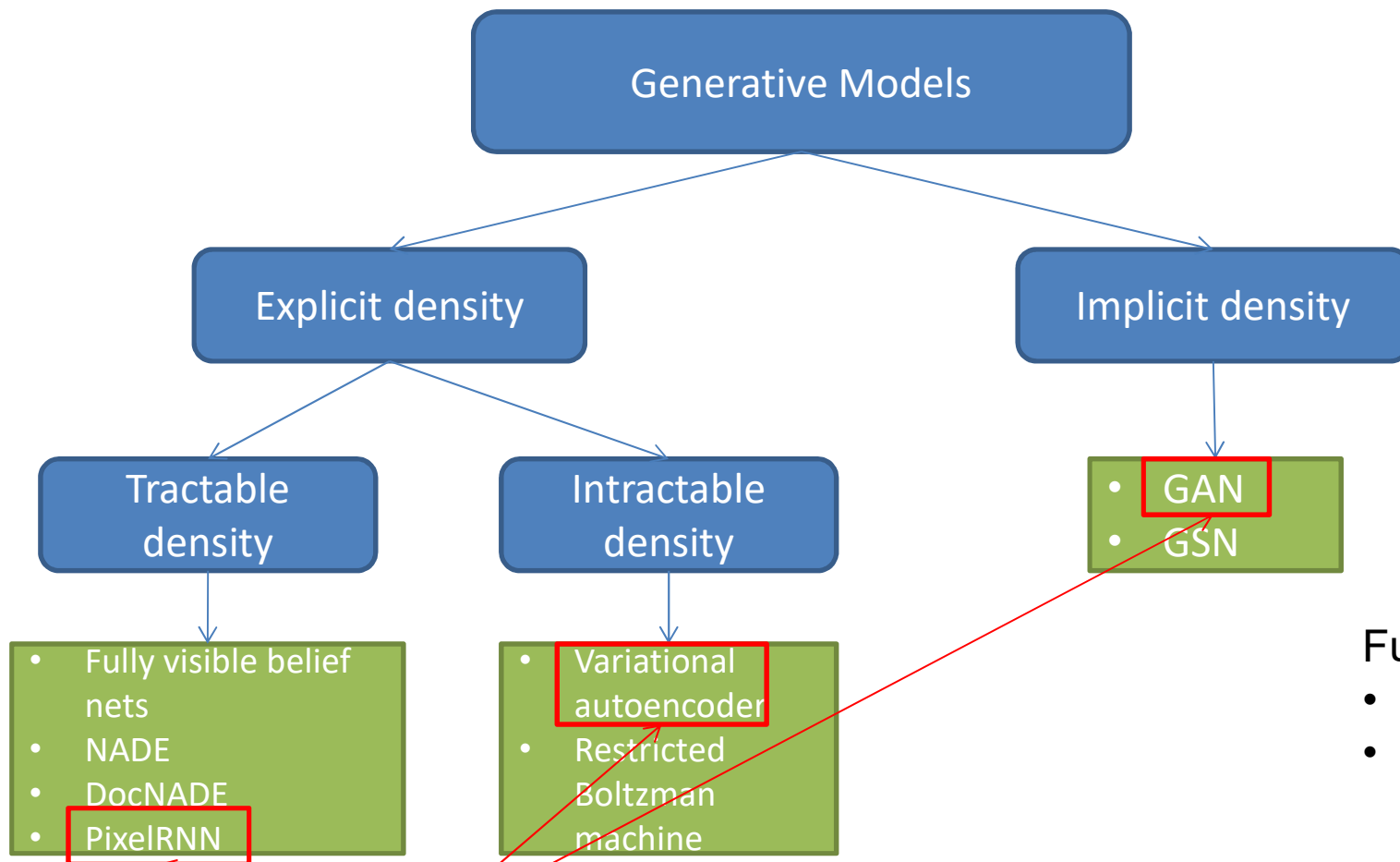
Types of generative models



Lecture 2 weeks ago

Figure adapted from Goodfellow 2017 *arXiv:1701.00160*

Types of generative models



Today's lecture

Further reading

- GSN paper
- Belief nets

Figure adapted from Goodfellow 2017, *arXiv:1701.00160*

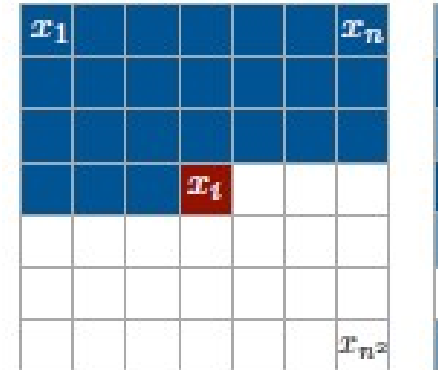
Autoregressive models

- Explicit density model
- Split high-dimensional input data into sequence
 - predict small piece of system (current state previous states)
 - no more curse of dimensionality

$$p(\mathbf{x}) = \prod_{i=1}^{n^2} p(x_i | x_1, \dots, x_{i-1})$$

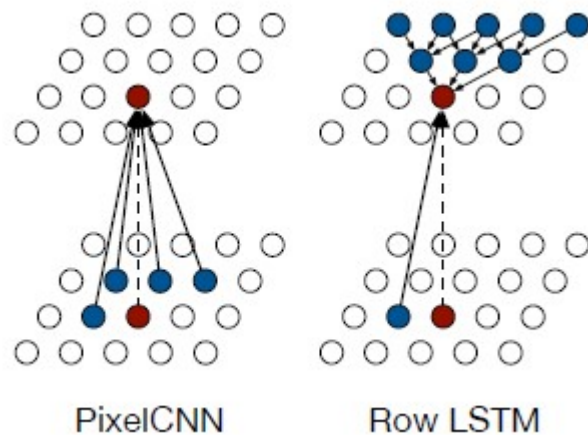
PixelRNN algorithm

- Use neural net to model distribution over pixel values
- Optimise weights by maximising likelihood of all images
- Need to choose order!
 - Start at a corner
 - Sequentially generate pixel values
- Use LSTM to model dependency on previous pixels



PixelCNN

- Same as before, but use CNN to model dependency on previous pixels
 - Masked convolutions
- Training more efficient thanks to possible parallelisations
 - Context for convs is known!
- Prediction is still slow



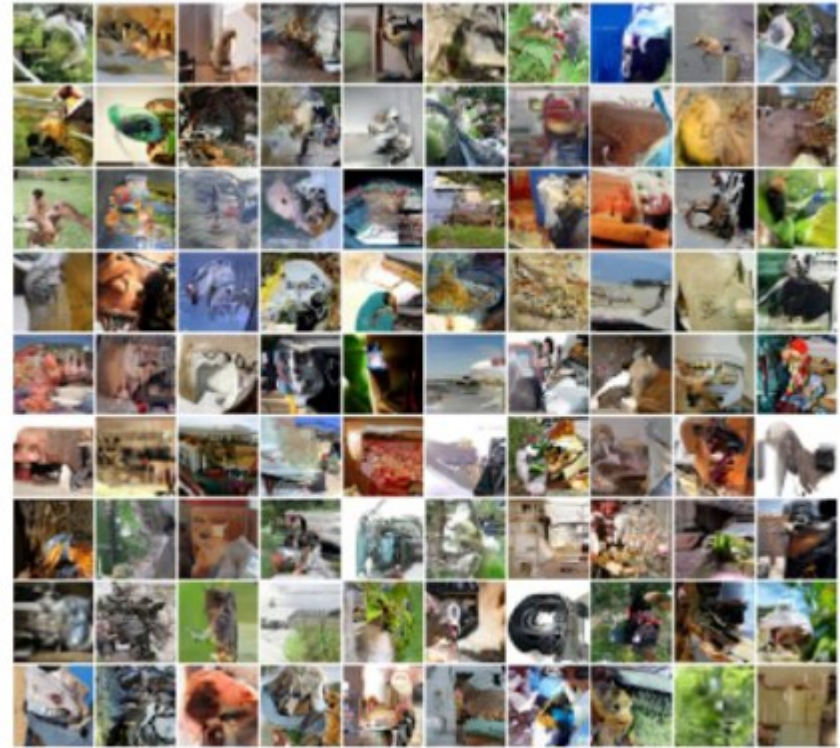
Further Reading

- Multi-scale RNN
- Conditional Image generation
- PixelCNN++
- Gated PixelCNN

Some results



32x32 CIFAR-10



32x32 ImageNet

Pros and Cons

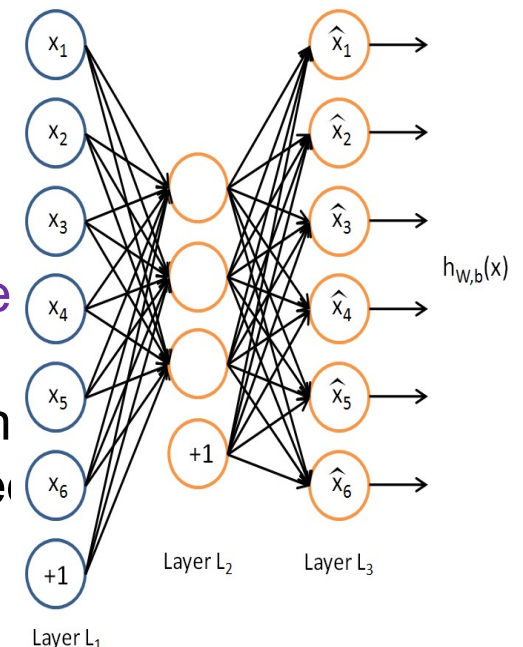
- Pros
 - Simple: just pick an order, no need to define prob distribution
 - easy to generate samples, like dreaming
- Cons
 - very expensive
 - as many predictions as pieces of data
 - parallelise during training but not testing
 - order dependance
 - where to start in an image?
 - how to deal with missing data?
 - Teacher forcing
 - difficult to generate long sequences
 - „blind representation“: not large structure of data that is actually interesting

Recap: Auto-encoders (AEs)

→ a *feed-forward neural network* trained to **reproduce its input** at the *output* layer

Key Facts about AEs:

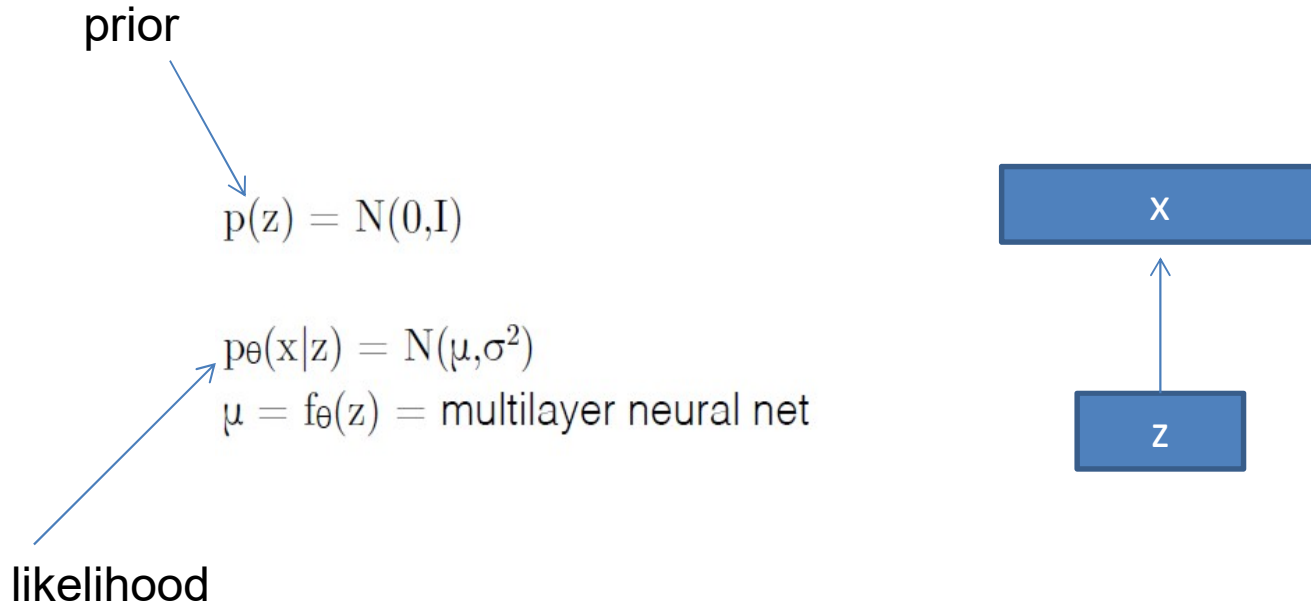
- unsupervised ML algorithm, similar to PCA
- neural network's target output is its input
- learn latent features/encoding of input (**no manual feature engineering**)
- represent both **linear** and **non-linear** transformation
- layered to form deep learning network, i.e., distributed representations
- tractable / easier optimization
- applications in **denoising** and **dimensionality** reduction (**dense representation**)
- powerful non-linear (i.e., non-linear encoding and decoding) generalization of PCA



<http://ufldl.stanford.edu/tutorial/unsupervised/Autoencoders/>

Variational autoencoders

- Probabilistic version of autoencoder
 - Place prior on latent space
 - Sample from prior and use decoder to generate new samples
 - AE as generative model



Inference

- With flexible neural net $f_{\theta}(z)$, the data distribution $p_{\theta}(x)$ can be almost arbitrarily complicated / multi-modal distribution
- But intractable posterior distribution $p(z|x)$
- Need approximate inference for learning

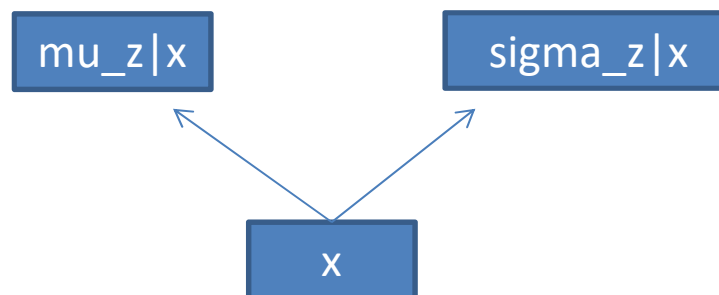
posterior $p_{\theta}(z|x) = p_{\theta}(x|z)p_{\theta}(z)/p_{\theta}(x)$

Data likelihood $p_{\theta}(x) = \int p_{\theta}(z)p_{\theta}(x|z)dz$

Variational inference with neural networks

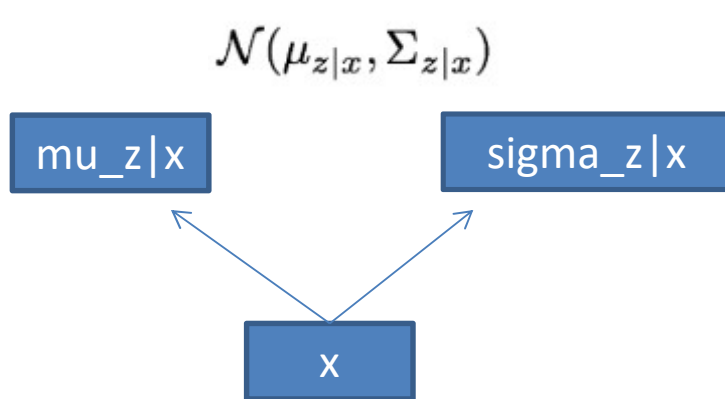
- Posterior $p(z|x)$ is not tractable
- Introduce parametric model $q_\phi(z|x)$ of true posterior
 - ϕ : variational parameters
 - parameterised by neural networks

$$q_\phi(z|x) = N(\mu, \sigma^2)$$
$$[\mu, \sigma^2] = f^{(z|x)}(x, \phi) = \text{multilayer neural net}$$

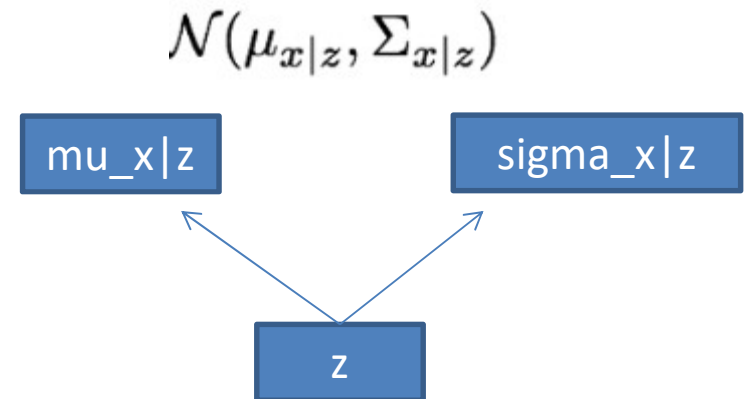


Encoder and decoder

- 2 NNs, encoder network $q_{\phi}(z|x)$ and decoder network $p_{\theta}(x|z)$



Encoder network q_{ϕ}



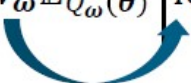
Decoder network p_{θ}

Recap: Variational inference

- Approximate posterior p with q -distribution
- Minimize KL divergence between q and p
 - Equivalent: maximise Evidence Lower Bound (ELBO) of data D

$$D_{KL}(Q_{\omega}(\boldsymbol{\theta}) || P(\boldsymbol{\theta}|D)) = \int d\boldsymbol{\theta}. Q_{\omega}(\boldsymbol{\theta}) \log \frac{Q_{\omega}(\boldsymbol{\theta})}{P(\boldsymbol{\theta}|D)}$$

$$\nabla_{\omega} \text{ELBO}(D) = \nabla_{\omega} \int Q_{\omega}(\boldsymbol{\theta}) \log P(D|\boldsymbol{\theta}) - Q_{\omega}(\boldsymbol{\theta}) \log \frac{Q_{\omega}(\boldsymbol{\theta})}{P(\boldsymbol{\theta})} d\boldsymbol{\theta}$$

$$= \nabla_{\omega} \mathbb{E}_{Q_{\omega}(\boldsymbol{\theta})} \left[\log P(d|\boldsymbol{\theta}) - \log \frac{Q_{\omega}(\boldsymbol{\theta})}{P(\boldsymbol{\theta})} \right]$$


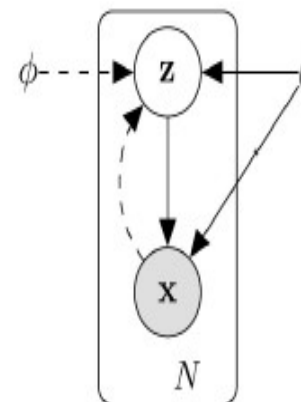
If we could move the derivative into the Expectation, we could approximate it by sampling!

Variational inference with neural networks

$$\underbrace{\mathbf{E}_z \left[\log p_{\theta}(x^{(i)} | z) \right] - D_{KL}(q_{\phi}(z | x^{(i)}) || p_{\theta}(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

$$\theta^*, \phi^* = \arg \max_{\theta, \phi} \sum_{i=1}^N \mathcal{L}(x^{(i)}, \theta, \phi)$$

Training: Maximize lower bound



- Jointly optimize w.r.t. ϕ and θ
- Simple SGD:
 - Sampling small minibatches of data
 - Sampling from approx. posterior
 - Use reparametrisation trick to approximate gradient of ELBO

Recap: reparameterisation trick

$$\nabla_{\omega} \mathbb{E}_{p(z|\omega)}[f(z)] = \nabla_{\omega} \int p(z|\omega) f(z) dz$$

Find a way to reparametrize $p(z|\omega)$ such that $z \sim g(\varepsilon, \omega)$ and we just sample ε

$$\nabla_{\omega} \int p(z|\omega) f(z) dz = \nabla_{\omega} \int p(\varepsilon) f(g(\varepsilon, \omega)) d\varepsilon$$

We move the derivative inside the integral:

$$= \int p(\varepsilon) \nabla_{\omega} f(g(\varepsilon, \omega)) d\varepsilon = \mathbb{E}_{p(\varepsilon)}[\nabla_{\omega} f(g(\varepsilon, \omega))]$$

warranted by the
dominated
convergence
theorem

Example:

$$z \sim \mathcal{N}(\mu, \sigma) \rightarrow z \sim \mu + \sigma \varepsilon; \varepsilon \sim \mathcal{N}(0, 1)$$

☑ ε is now independent of μ, σ and backpropagation works!

Stochastic Gradient VB

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}) &= \int q_{\boldsymbol{\phi}}(\mathbf{z}) [\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}) - \log q_{\boldsymbol{\phi}}(\mathbf{z})] d\mathbf{z} \\ &\simeq \frac{1}{L} \sum_{l=1}^L \left(\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}^{(l)}) - \log q_{\boldsymbol{\phi}}(\mathbf{z}^{(l)}) \right)\end{aligned}$$

where $\boldsymbol{\epsilon}^{(l)} \sim p(\boldsymbol{\epsilon})$ (samples from noise variable)

$$\mathbf{z}^{(l)} = g(\boldsymbol{\epsilon}^{(l)}, \boldsymbol{\phi})$$

(such that $\mathbf{z}^{(l)} \sim q_{\boldsymbol{\phi}}(\mathbf{z})$)

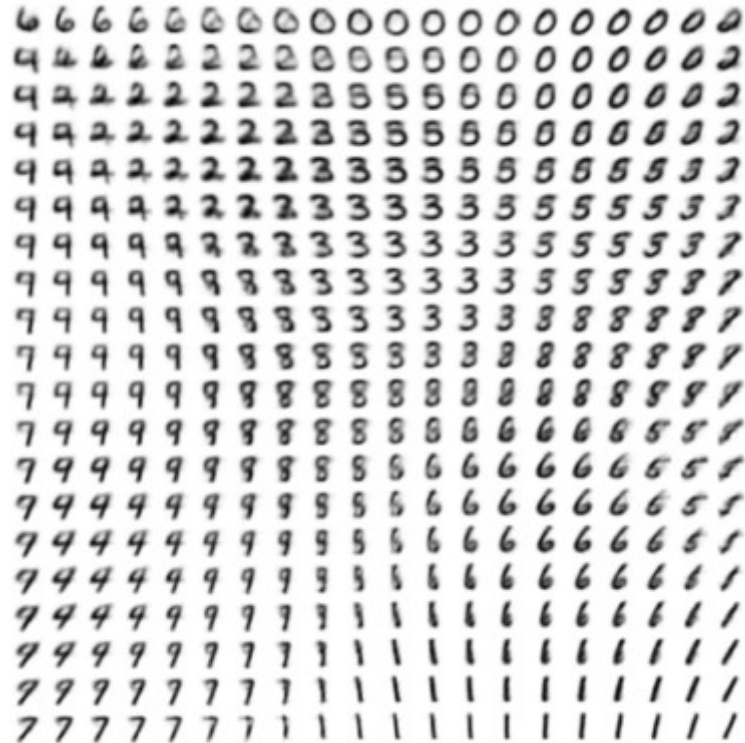
Stochastic VB in practice

- Draw mini-batch
- Sample from $p(\epsilon)$
- Compute gradients using backprop
- Update θ and ϕ

Some results



(a) Learned Frey Face manifold



(b) Learned MNIST manifold

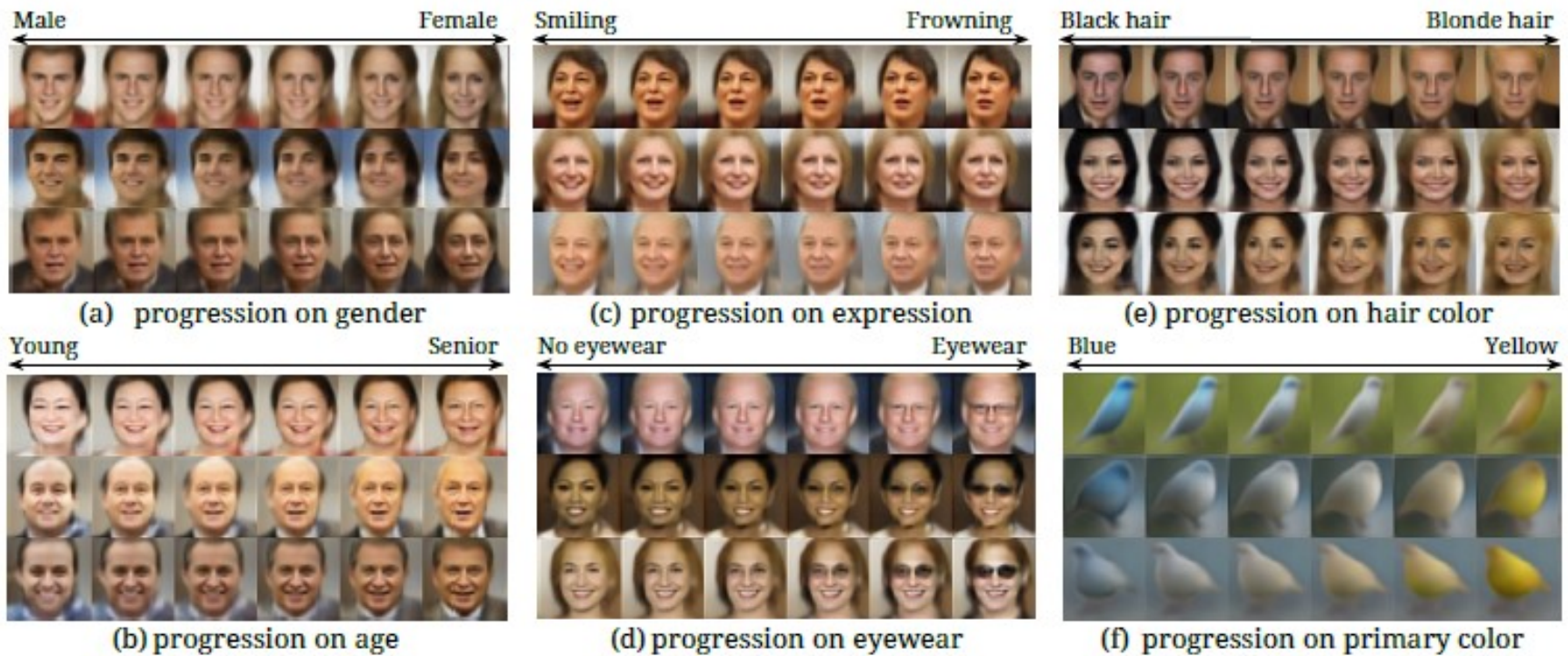
- Diagonal prior leads to independant z_i s
- Components are interpretable
- Representation is accesible (via $q(z|x)$)

Some improvements

- Improve encoder/decoder
 - Use convolutions
 - Conditional VAE
 - Replace all $P(X|z)$ with $P(X|z,Y)$
 - Replace all $Q(z|X)$ with $Q(z|X,Y)$
 - Hierarchical VAE
- Improve prior
 - Problem: mis-match prior and aggregate posterior
 - Bad samples for high density in prior and low density in posterior (hasn't seen samples!)
 - Bad reconstruction error



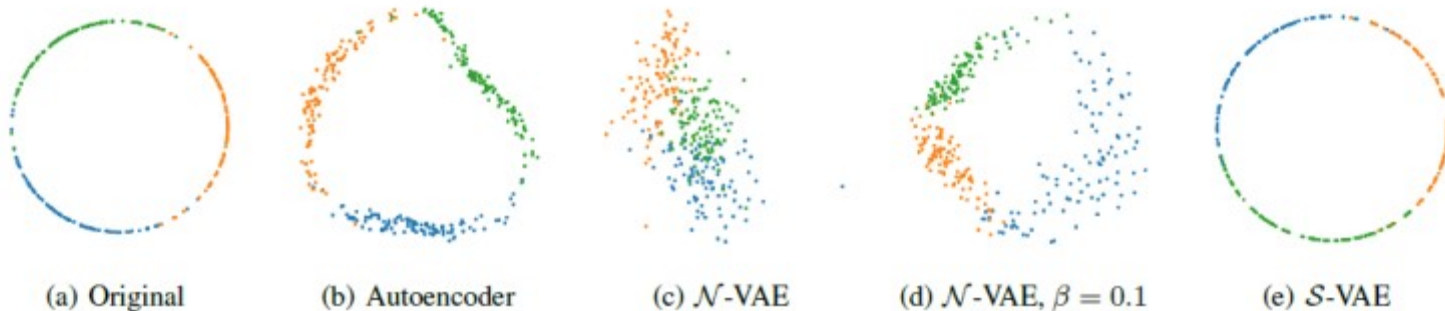
Results



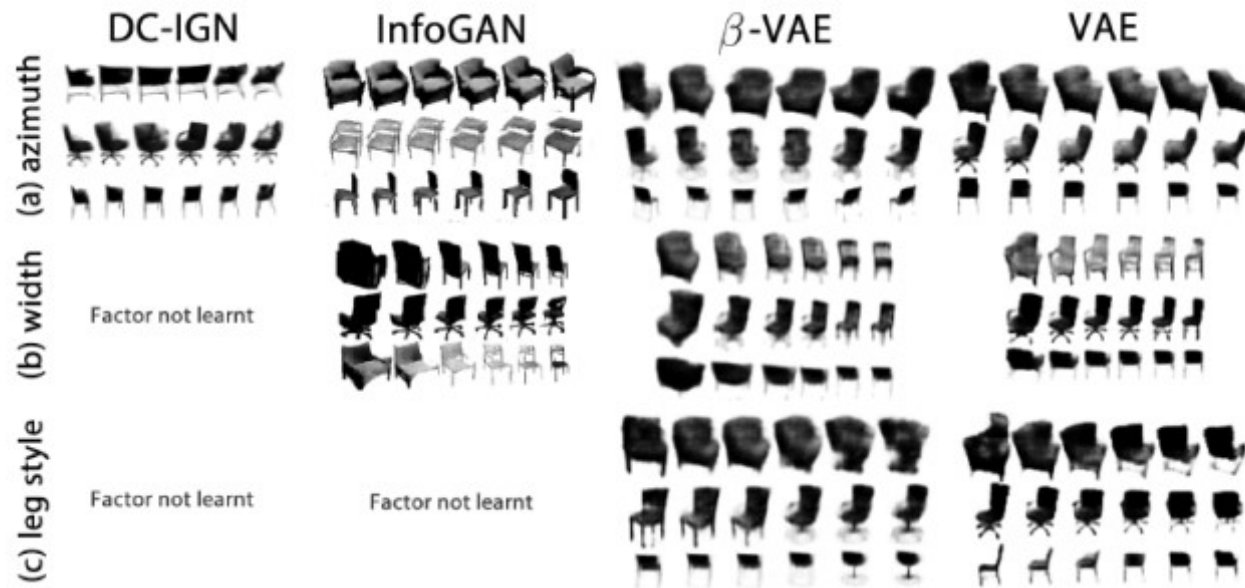
More improvements

- Beta-VAE
 - Encourage disentangled factors
 - introduce an adjustable hyperparameter that balances independence constraints with reconstruction accuracy

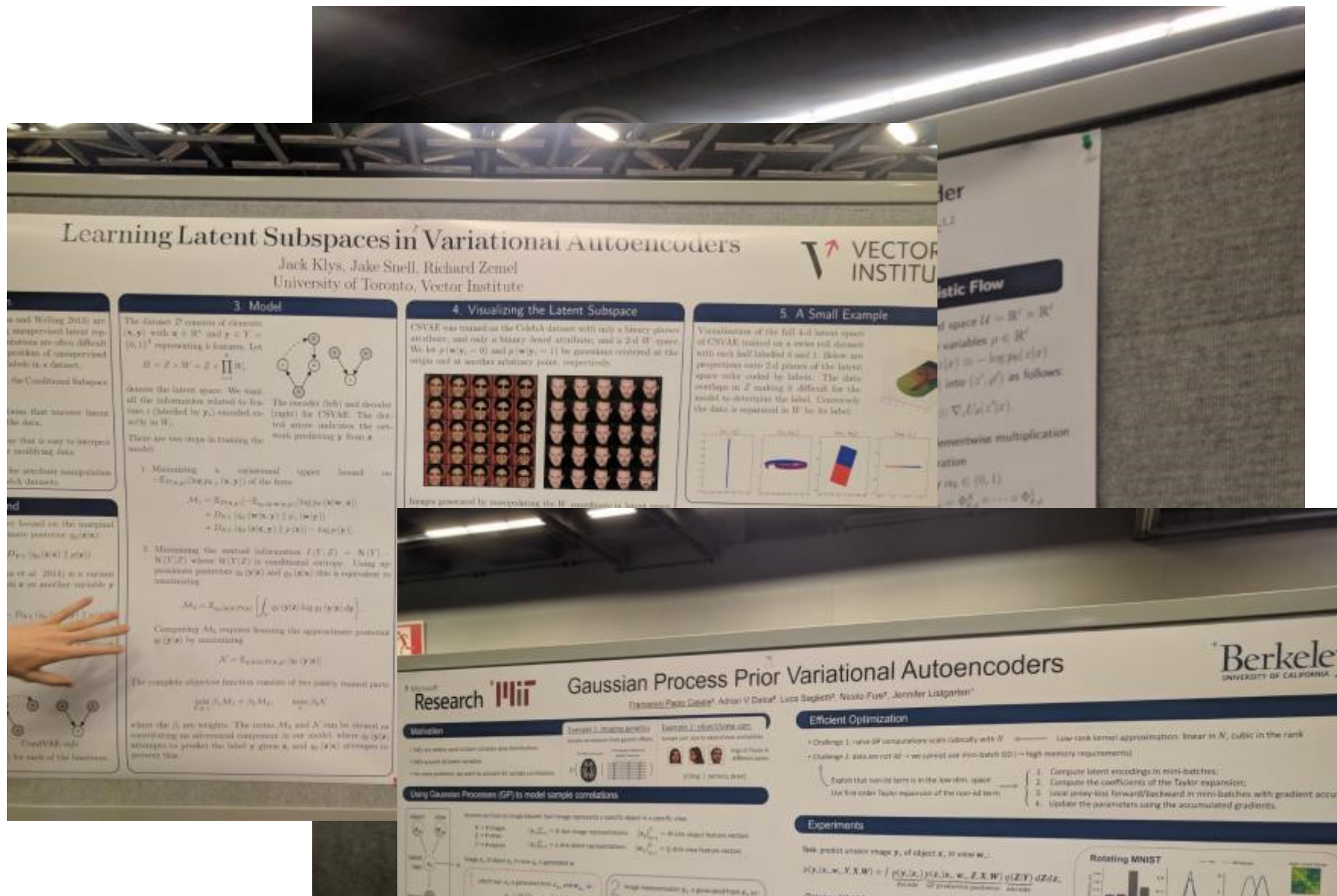
$$\mathcal{F}(\theta, \phi, \beta; \mathbf{x}, \mathbf{z}) \geq \mathcal{L}(\theta, \phi; \mathbf{x}, \mathbf{z}, \beta) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \beta D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$$



Beta-VAE and S-VAE



Some current research directions



Pros and Cons of VAEs

- Pros
 - Accesible representation
 - Robust and straight-forward to train
- Cons
 - Generated images blurrier than SOTA

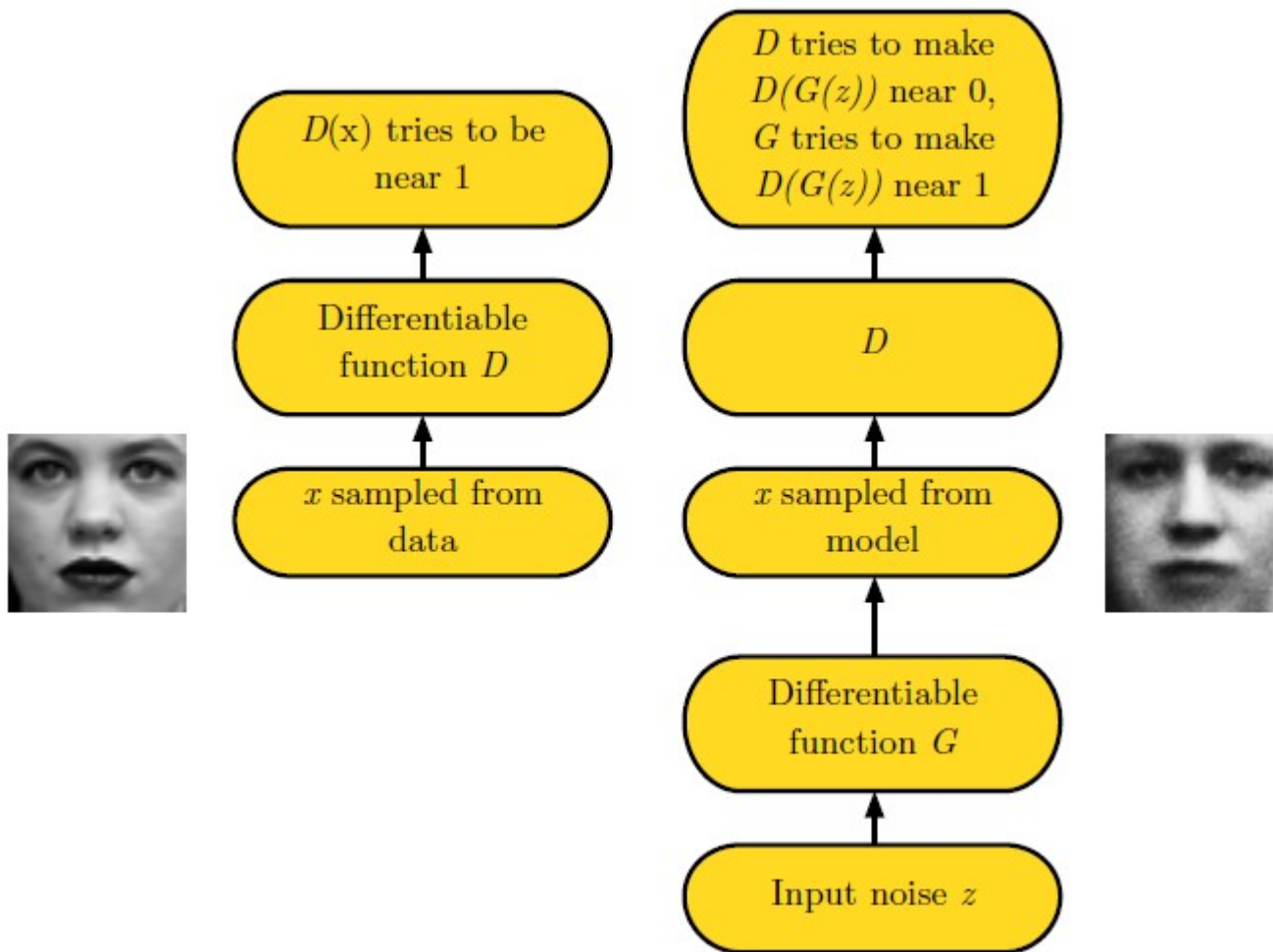
GANs

- We would like to sample from a complex, high-dimensional distribution
 - Problem: No direct way for this
 - Solution: Sample random noise and pass it through NN

A two-player game

- Set up a game between two players
- **Generator G**: generate samples that are intended to come from the same distribution as the training data
- **Discriminator D**: determines whether a determine whether a sampel is real or fake
 - use traditional supervised learning techniques
- Generator is trained to fool discriminator

Illustration of the game



Learning the parameters of G and D

- Cost functions of generator and discriminator (J_D and J_G) depend on both sets of parameters (θ_D and θ_G)
 - But: Each network has only access to its own parameters
 - Not optimisation, but game!
- Solution to game
 - Nash equilibrium: Tuple (θ_G, θ_D) that is a local optimum of J_D wrt θ_D and a local optimum of J_G wrt θ_G

Cost function of D

- Minimize cross-entropy
 - Train on 2 mini-batches
 - one coming from the dataset, where the label is 1 for all examples
 - one coming from the generator, where the label is 0 for all examples
- Co-operative view
 - Discriminator more like a teacher instructing the generator in how to improve

Zero-sum game

- Sum of all players' costs is always zero
 - $J_G = -J_D$
 - Also referred to as minmax
 - minimization in outer loop and maximization in inner loop

$$\theta^{(G)*} = \arg \min_{\theta^{(G)}} \max_{\theta^{(D)}} V \left(\theta^{(D)}, \theta^{(G)} \right)$$



Value function

$$V \left(\theta^{(D)}, \theta^{(G)} \right) = -J^{(D)} \left(\theta^{(D)}, \theta^{(G)} \right)$$

Heuristic non-saturating game

- Cost from zero-sum game does not perform well in practice
 - D minimizes a cross-entropy but G maximises same cross-entropy
 - When D rejects samples with high confidence, gradient of G vanishes
- Solution: flip target of cross-entropy for cost for G
 - Maximise log-prob of D being mistaken

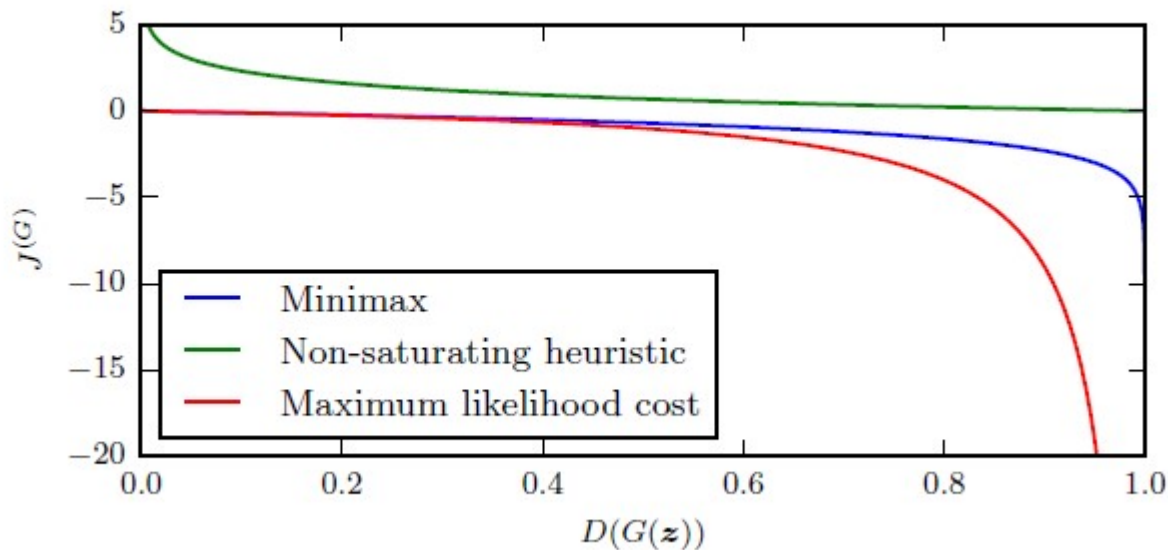
$$J^{(G)} = -\frac{1}{2} \mathbb{E}_z \log D(G(z))$$

- Each player has strong gradient when he loses the game

Maximum likelihood game

- Minimizing the KL divergence between the data and the model
 - Equivalent (if D is optimal)

$$J^{(G)} = -\frac{1}{2} \mathbb{E}_z \exp(\sigma^{-1}(D(G(z))))$$



Putting it together

for number of training iterations **do**

for k steps **do**

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

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

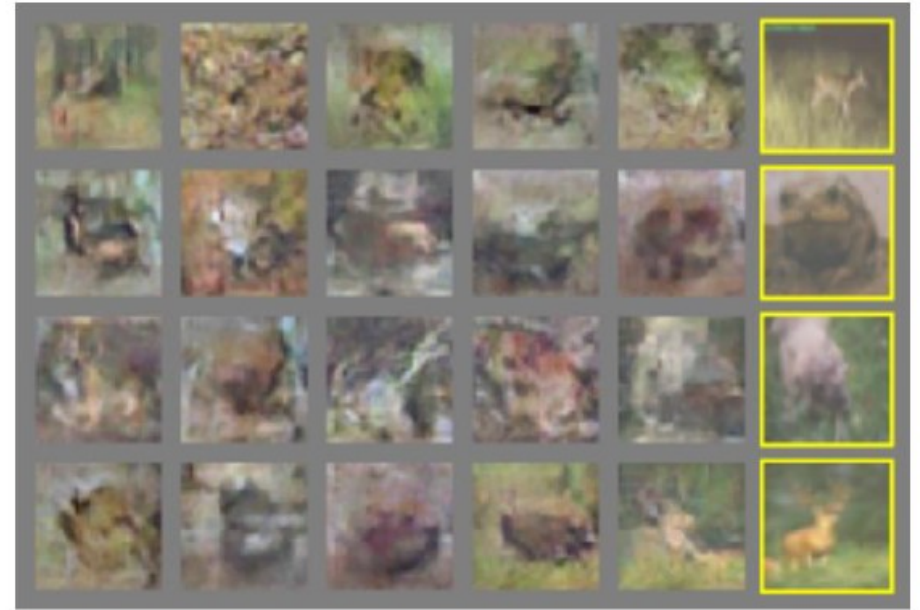
end for

- Sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.
- 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

Samples from generator

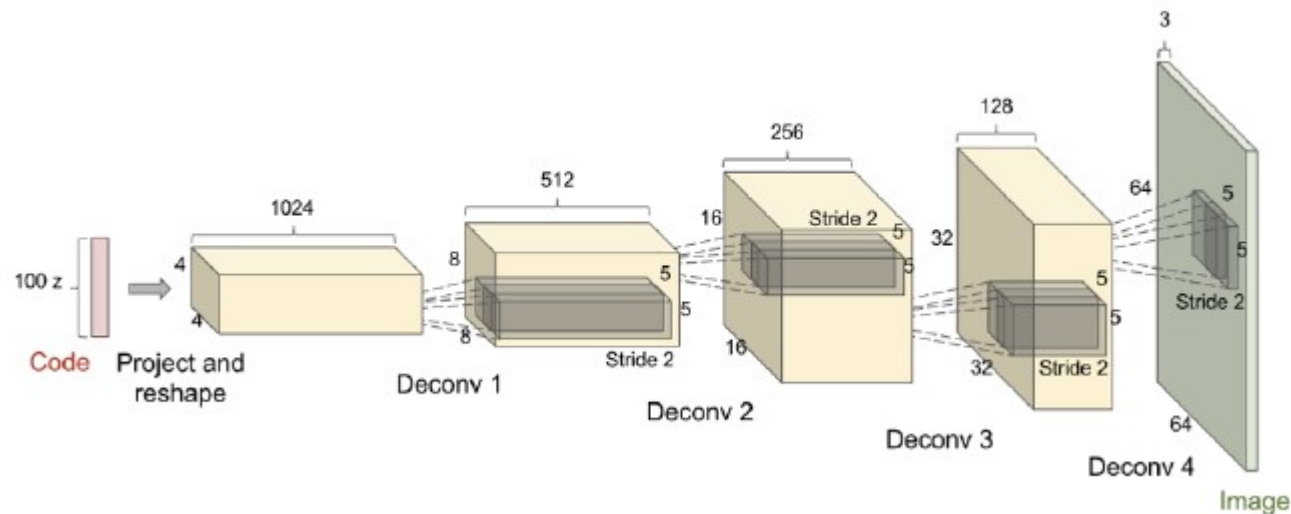


Nearest neighbour from training set

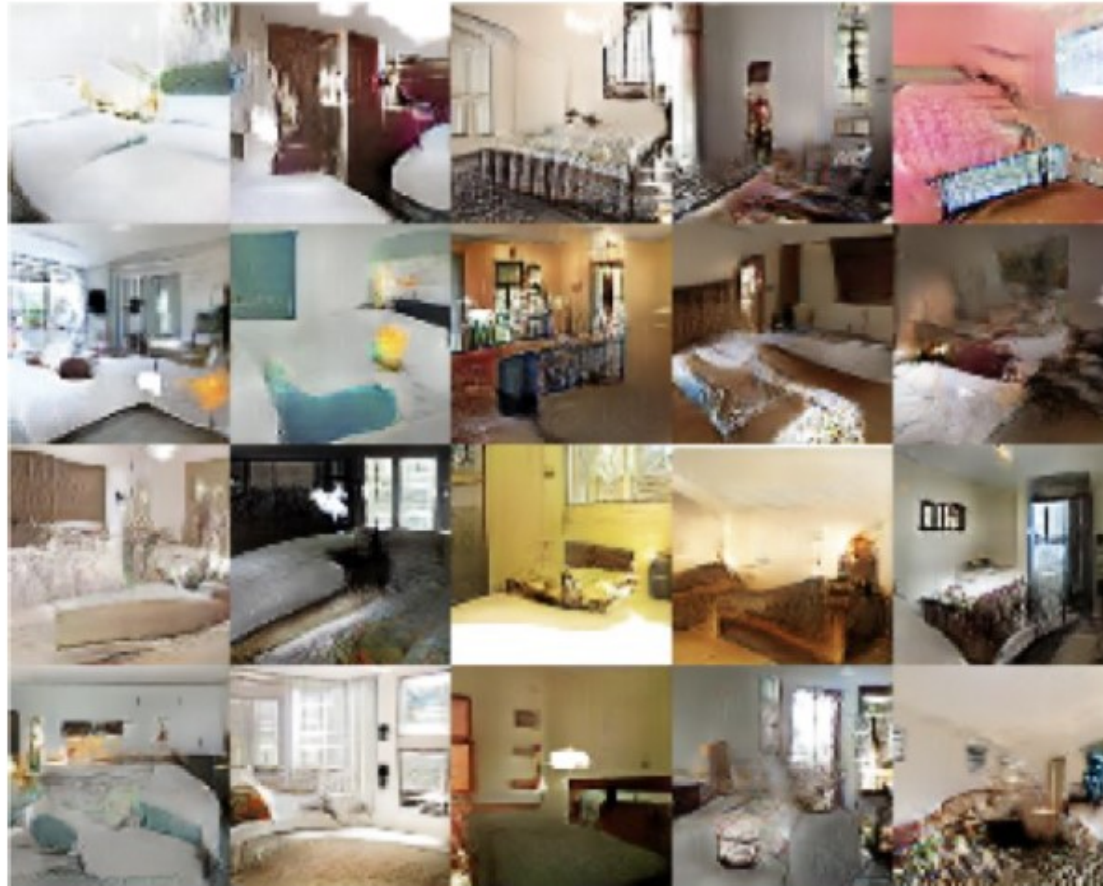
DGAN

Architecture guidelines for stable Deep Convolutional GANs

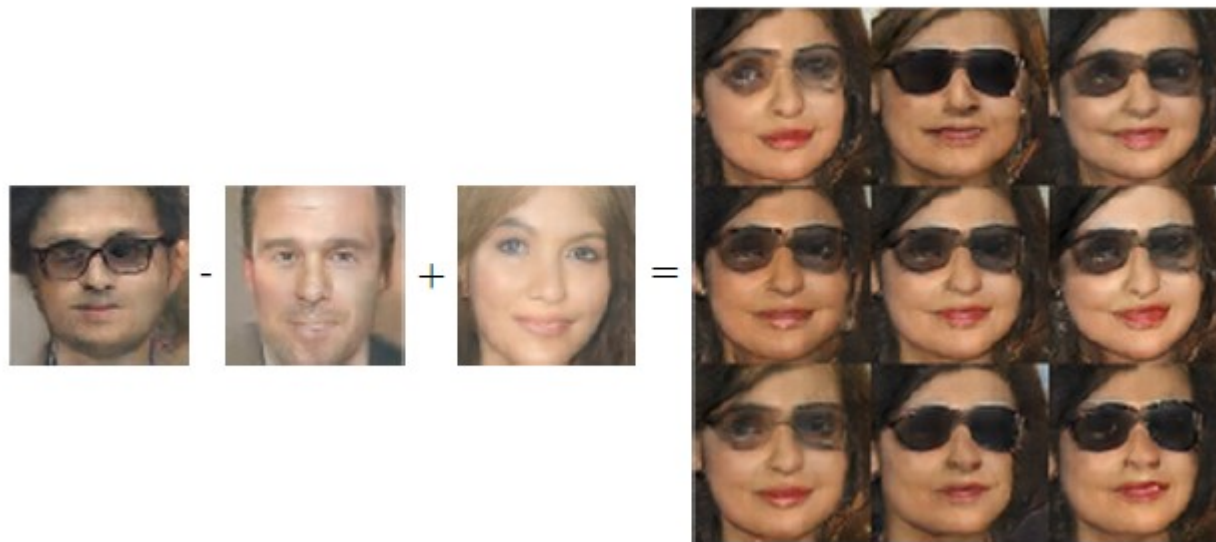
- Replace any pooling layers with strided convolutions (discriminator) and fractional-strided convolutions (generator).
- Use batchnorm in both the generator and the discriminator.
- Remove fully connected hidden layers for deeper architectures.
- Use ReLU activation in generator for all layers except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



Sample again



GAN maths

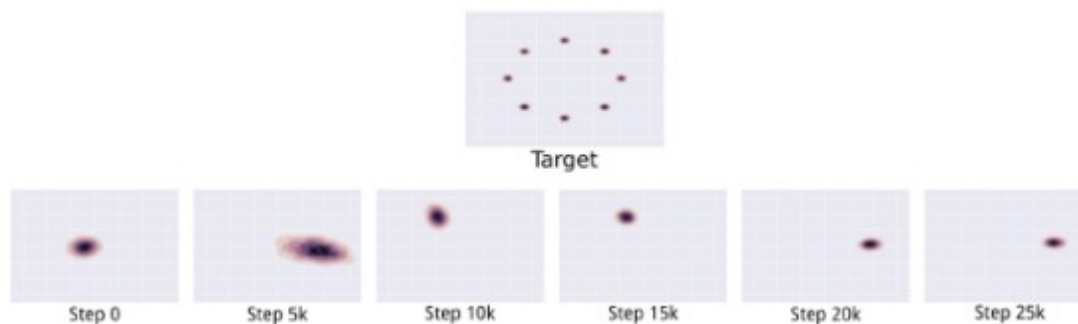


Mode collapse

- Most severe problem in terms of non-convergence
- Issue: maximin solution to the GAN game is different from the minimax solution
 - Generator asked to map every z value to the single x coordinate that discriminator believes is most likely to be real
 - Simultaneous gradient descent doesn't favour one over the other

$$\theta^{(G)*} = \arg \min_{\theta^{(G)}} \max_{\theta^{(D)}} V(\theta^{(D)}, \theta^{(G)})$$

\longleftrightarrow

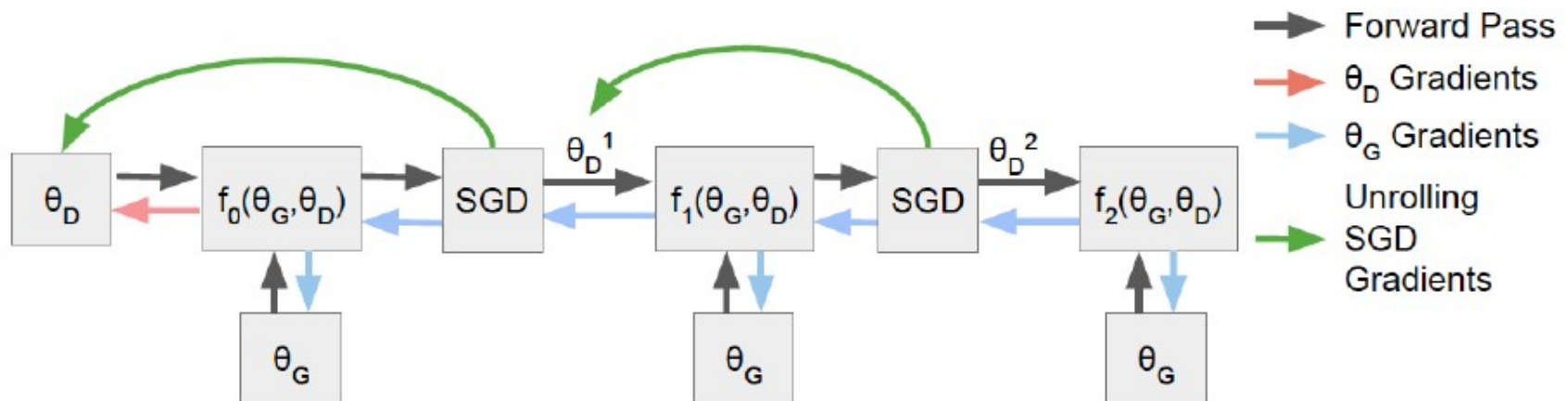


Solutions

- Minibatch features
 - Compare one example to batches of real/fake examples
 - D can detect if sample is unusually similar to other samples
- Unrolled GAN
 - back-propagate through the maximization operation

Unrolled GAN

- Consider several updates of the generator when updating the discriminator and vice versa
 - k steps in the discriminator
 - backpropagate all steps when computing the gradient on the generator

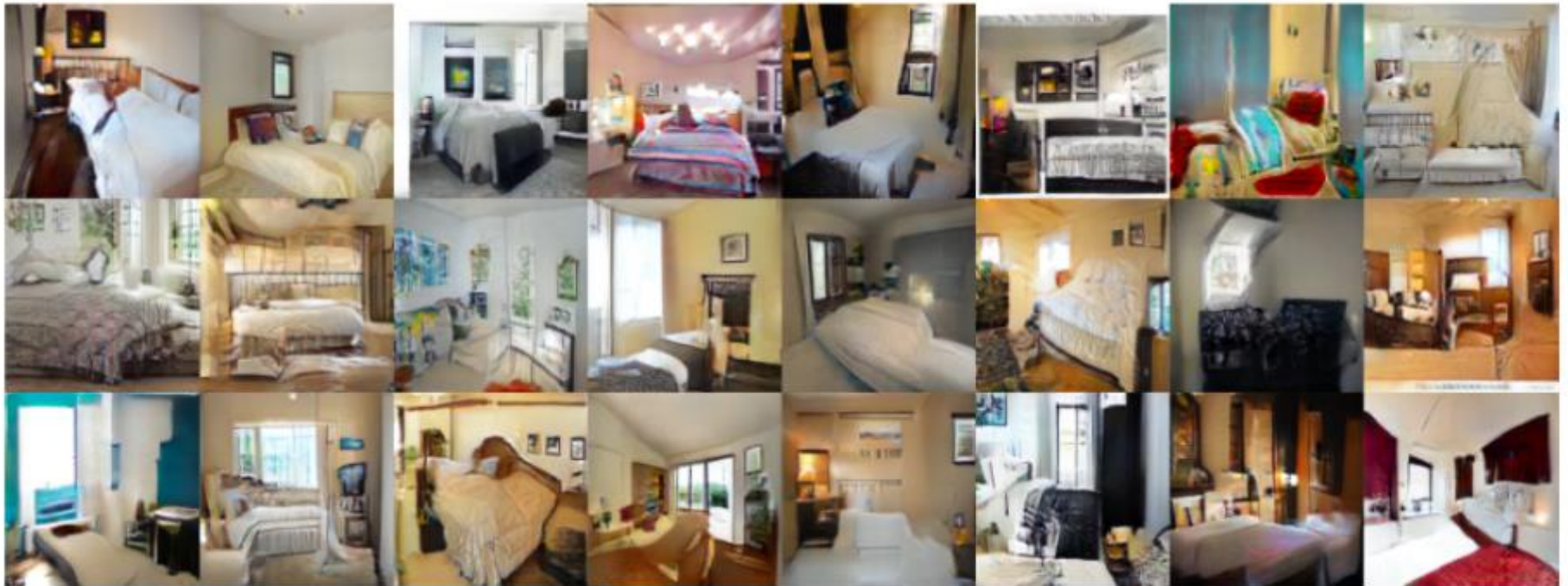


Better generators: LSGAN,

- Least squares GAN
 - X. Mao, Q. Li, H. Xie, R. Lau, Z. Wang, “Least squares generative adversarial networks” 2016
 - Still use a classifier but replace cross-entropy loss with Euclidean loss

Discriminator	
GAN	$\min_D E_{x \sim p_X} [-\log D(x)] + E_{z \sim p_Z} [-\log(1 - D(G(z)))]$
LSGAN	$\min_D E_{x \sim p_X} [(D(x) - 1)^2] + E_{z \sim p_Z} [D(G(z))^2]$
Generator	
GAN	$\min_G E_{z \sim p_Z} [-\log D(G(z))]$
LSGAN	$\min_G E_{z \sim p_Z} [(D(G(z)) - 1)^2]$

LSGAN

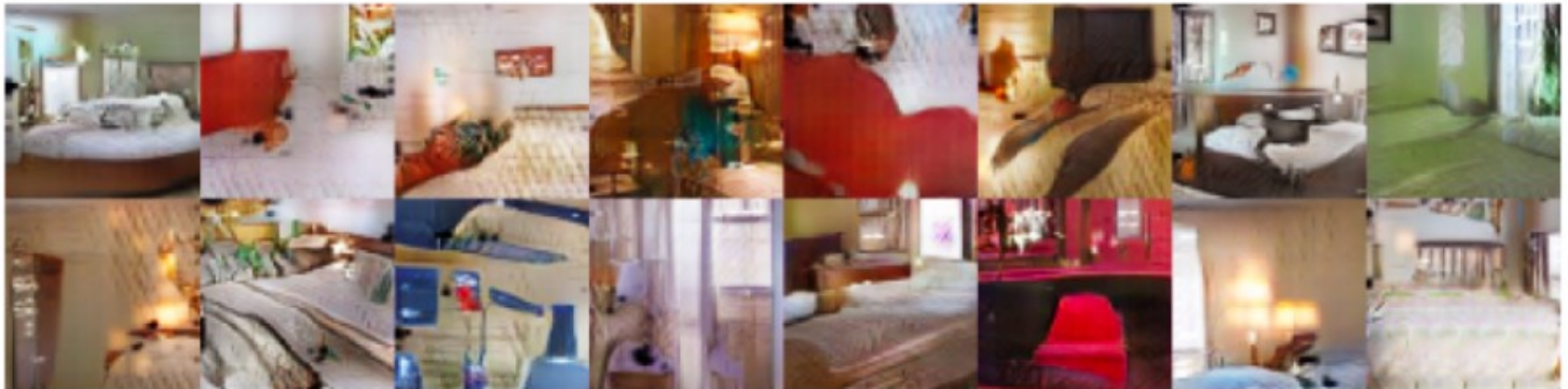


Wasserstein GAN

- M. Arjovsky, S. Chintala, L. Bottou “Wasserstein GAN” 2016
- Use critic instead of discriminator
 - Discriminator can output real number (same as before w/o sigmoid)

Discriminator	
GAN	$\max_D E_{x \sim p_X} [\log D(x)] + E_{z \sim p_Z} [\log(1 - D(G(z)))]$
WGAN	$\max_D E_{x \sim p_X} [D(x)] - E_{z \sim p_Z} [D(G(z))]$
Generator	
GAN	$\max_G E_{z \sim p_Z} [\log D(G(z))]$
WGAN	$\max_G E_{z \sim p_Z} [D(G(z))]$

WGAN



WGAN-GP

- I. Gulrajani, F. Ahmed, M. Arjovsky, V. Domoulin, A. Courville “Improved Training of Wasserstein GANs” 2017

$$\min_G \max_D E_{x \sim p_X} [D(x)] - E_{z \sim p_Z} [D(G(Z))] + \lambda E_{y \sim p_Y} [(\|\nabla_y D(y)\|_2 - 1)^2]$$

$$y = ux + (1 - u)G(z)$$

- y : imaginary samples

Optimal critic has unit gradient norm almost everywhere

DCGAN

LSGAN

WGAN (clipping)

WGAN-GP (ours)

Baseline (G : DCGAN, D : DCGAN)

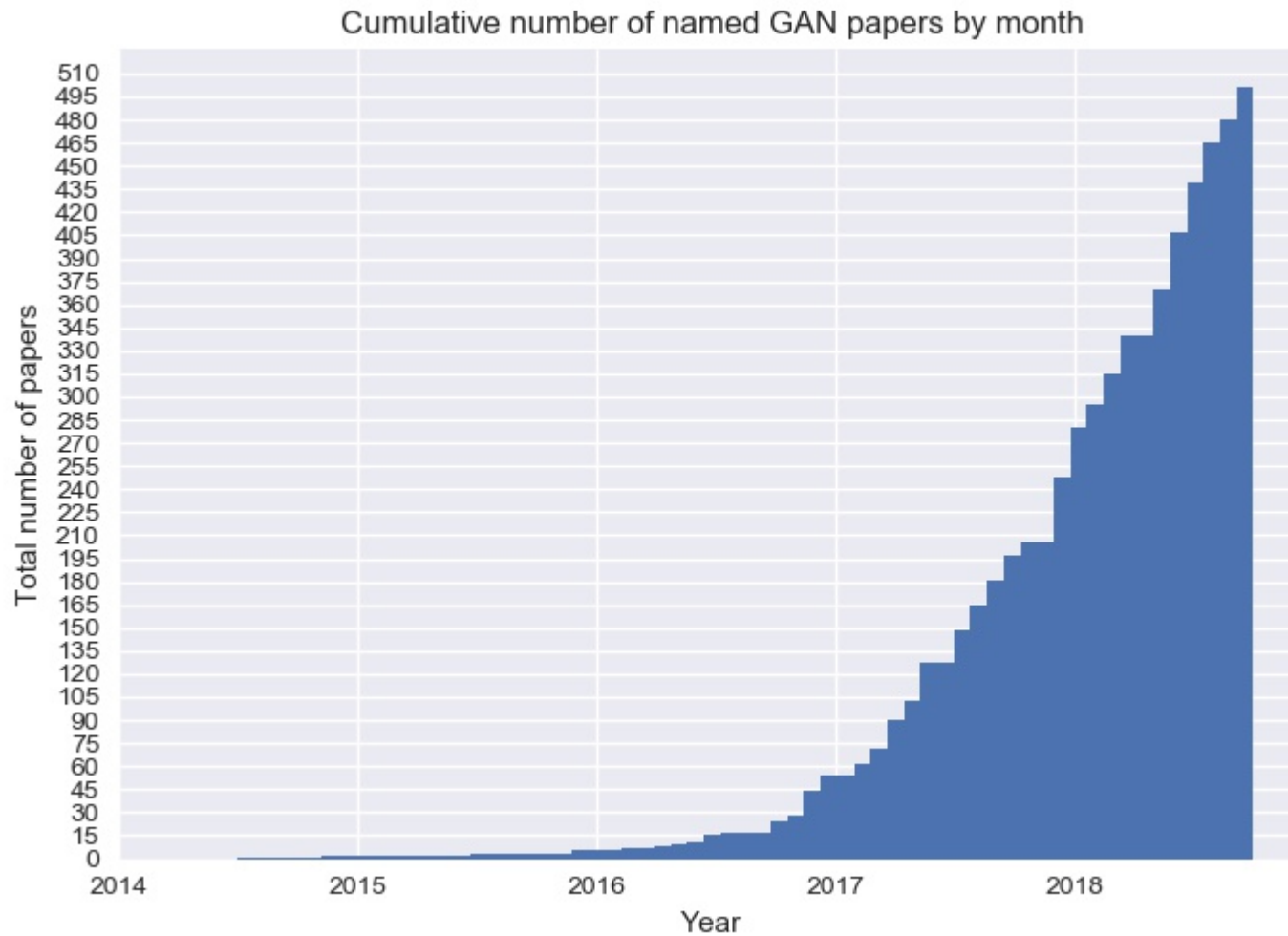


BEGAN and DRAGAN

- DRAGAN: Add gradient norm to standard GAN and evaluate around the data manifold
- BEGAN: use autoencoder as discriminator and optimize lower bound of the Wasserstein distance between auto-encoder loss distributions on real and fake data.

DRAGAN	$\mathcal{L}_D^{\text{DRAGAN}} = \mathcal{L}_D^{\text{GAN}} + \lambda \mathbb{E}_{\hat{x} \sim p_d + \mathcal{N}(0, c)} [(\ \nabla D(\hat{x})\ _2 - 1)^2]$	$\mathcal{L}_G^{\text{DRAGAN}} = \mathbb{E}_{\hat{x} \sim p_g} [\log(1 - D(\hat{x}))]$
BEGAN	$\mathcal{L}_D^{\text{BEGAN}} = \mathbb{E}_{x \sim p_d} [\ x - \text{AE}(x)\ _1] - k_t \mathbb{E}_{\hat{x} \sim p_g} [\ \hat{x} - \text{AE}(\hat{x})\ _1]$	$\mathcal{L}_G^{\text{BEGAN}} = \mathbb{E}_{\hat{x} \sim p_g} [\ \hat{x} - \text{AE}(\hat{x})\ _1]$

The GAN zoo



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

So....which one is best?

Are GANs Created Equal? A Large-Scale Study

Mario Lucic* Karol Kurach* Marcin Michalski Olivier Bousquet Sylvain Gelly
Google Brain

than others. We conduct a neutral, multi-faceted large-scale empirical study on state-of-the-art models and evaluation measures. We find that most models can reach similar scores with enough hyperparameter optimization and random restarts. This suggests that improvements can arise from a higher computational budget and tuning more than fundamental algorithmic changes. To overcome some limitations

Evaluating GANs

- Challenging to define appropriate metric
 - Maximum likelihood and other classical metrics not applicable
 - Subjective comparisons (visual quality) may be misleading
 - Inception score: discriminator has low entropy, while producing samples from all classes when passed through a classifier

$$\exp(\mathbb{E}_{x \sim G}[d_{KL}(p(y | x), p(y))])$$

- Fréchet Inception Distance: difference in embedding of true and fake data (assuming MVN in embedded space)
 - Strong negative correlation between visual quality and FID

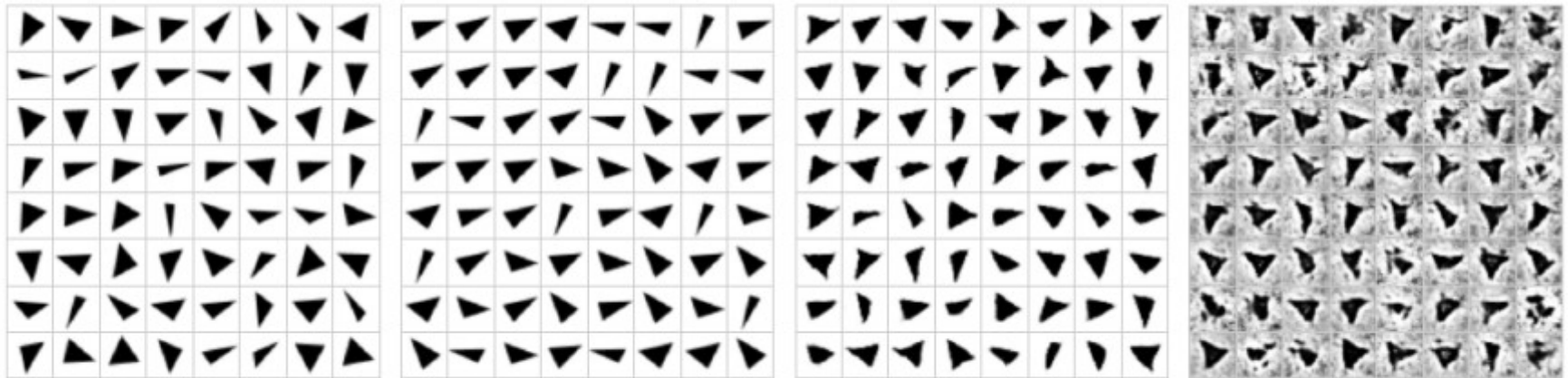
$$\text{FID}(x, g) = \|\mu_x - \mu_g\|_2^2 + \text{Tr}(\Sigma_x + \Sigma_g - 2(\Sigma_x \Sigma_g)^{\frac{1}{2}}),$$

- Precision/recall

Metrics

- Precision measures fraction of relevant retrieved instances among the retrieved instances
- Recall measures fraction of retrieved instances among relevant instances
- F1 score is harmonic average of precision and recall.
- IS captures precision: no penalization for not producing all modes of the data distribution
 - Only for not producing all classes
- FID captures both precision and recall

Precision-Recall for GANs



(a) High precision, high recall

(b) High precision, low recall

(c) Low precision, high recall

(d) Low precision, low recall

Fair assessment

- Compare state-of-the-art approaches

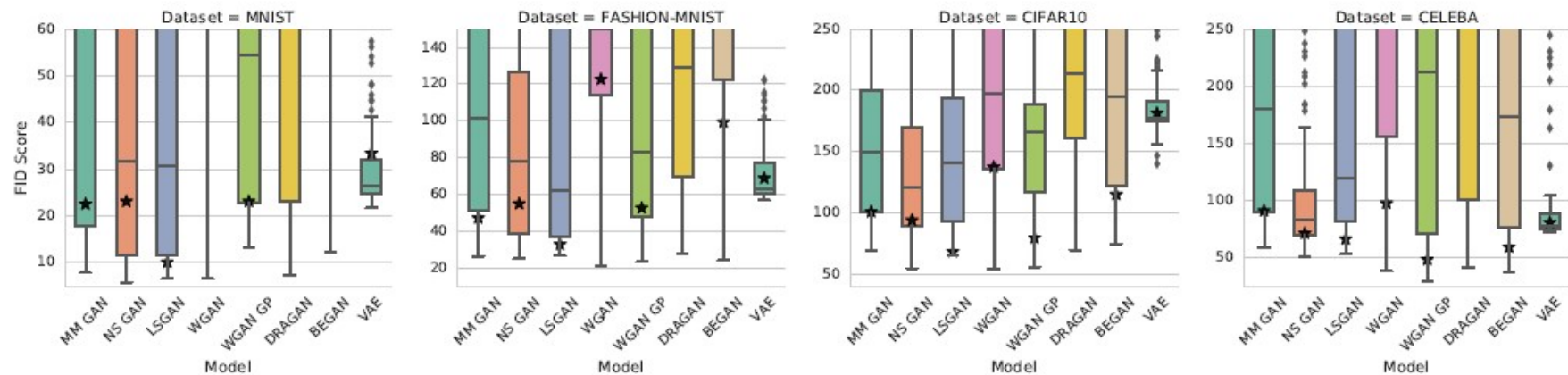
GAN	DISCRIMINATOR LOSS	GENERATOR LOSS
MM GAN	$\mathcal{L}_D^{\text{GAN}} = -\mathbb{E}_{x \sim p_d} [\log(D(x))] - \mathbb{E}_{\hat{x} \sim p_g} [\log(1 - D(\hat{x}))]$	$\mathcal{L}_G^{\text{GAN}} = \mathbb{E}_{\hat{x} \sim p_g} [\log(1 - D(\hat{x}))]$
NS GAN	$\mathcal{L}_D^{\text{NSGAN}} = -\mathbb{E}_{x \sim p_d} [\log(D(x))] - \mathbb{E}_{\hat{x} \sim p_g} [\log(1 - D(\hat{x}))]$	$\mathcal{L}_G^{\text{NSGAN}} = -\mathbb{E}_{\hat{x} \sim p_g} [\log(D(\hat{x}))]$
WGAN	$\mathcal{L}_D^{\text{WGAN}} = -\mathbb{E}_{x \sim p_d} [D(x)] + \mathbb{E}_{\hat{x} \sim p_g} [D(\hat{x})]$	$\mathcal{L}_G^{\text{WGAN}} = -\mathbb{E}_{\hat{x} \sim p_g} [D(\hat{x})]$
WGAN GP	$\mathcal{L}_D^{\text{WGANGP}} = \mathcal{L}_D^{\text{WGAN}} + \lambda \mathbb{E}_{\hat{x} \sim p_g} [(\nabla D(\alpha x + (1 - \alpha \hat{x})) _2 - 1)^2]$	$\mathcal{L}_G^{\text{WGANGP}} = -\mathbb{E}_{\hat{x} \sim p_g} [D(\hat{x})]$
LS GAN	$\mathcal{L}_D^{\text{LSGAN}} = -\mathbb{E}_{x \sim p_d} [(D(x) - 1)^2] + \mathbb{E}_{\hat{x} \sim p_g} [D(\hat{x})^2]$	$\mathcal{L}_G^{\text{LSGAN}} = -\mathbb{E}_{\hat{x} \sim p_g} [(D(\hat{x}) - 1)^2]$
DRAGAN	$\mathcal{L}_D^{\text{DRAGAN}} = \mathcal{L}_D^{\text{GAN}} + \lambda \mathbb{E}_{\hat{x} \sim p_d + \mathcal{N}(0, c)} [(\nabla D(\hat{x}) _2 - 1)^2]$	$\mathcal{L}_G^{\text{DRAGAN}} = \mathbb{E}_{\hat{x} \sim p_g} [\log(1 - D(\hat{x}))]$
BEGAN	$\mathcal{L}_D^{\text{BEGAN}} = \mathbb{E}_{x \sim p_d} [x - \text{AE}(x) _1] - k_t \mathbb{E}_{\hat{x} \sim p_g} [\hat{x} - \text{AE}(\hat{x}) _1]$	$\mathcal{L}_G^{\text{BEGAN}} = \mathbb{E}_{\hat{x} \sim p_g} [\hat{x} - \text{AE}(\hat{x}) _1]$

Fair comparisons

- Use same architecture
- Optimize hyperparameters on each dataset OR on one dataset only (infer for new datasets)
- Computational budget
 - Dependence on number of optimised hyperparameters

Are GANs created equal?

- Asterisk
 - Default hyperparameters

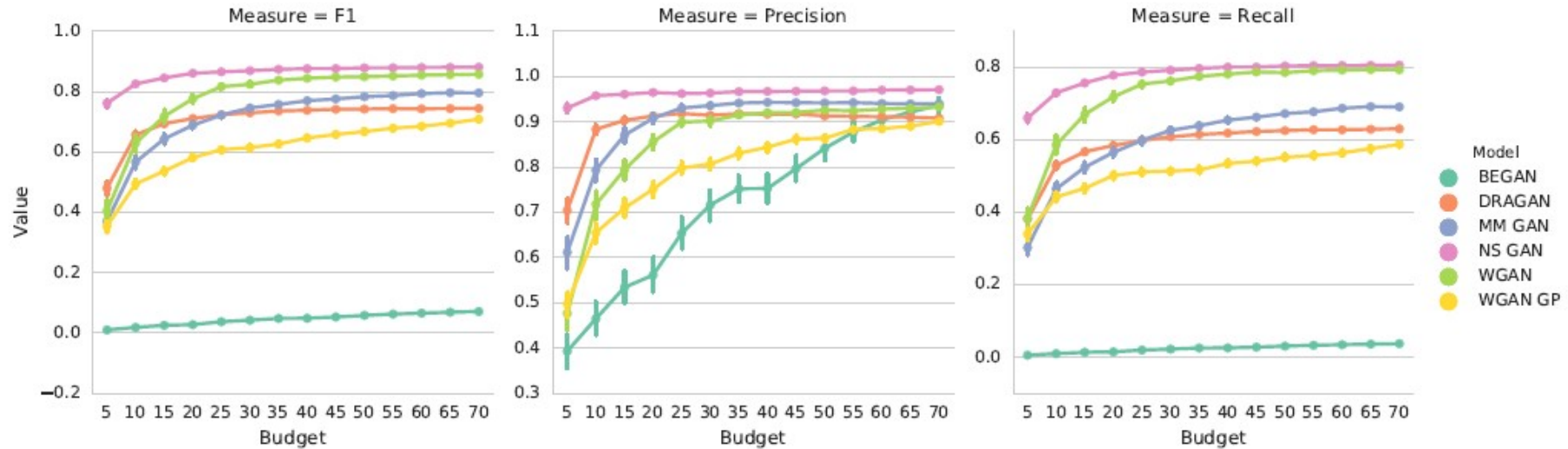


Large-scale hyperparameter optimisation

- No model strictly dominates the others
 - Strong dependence on dataset
- But: performance not SOTA
 - Larger networks would perform better
 - Authors report best FID (random seed optimisation!)

	MNIST	FASHION	CIFAR	CELEBA
MM GAN	9.8 ± 0.9	29.6 ± 1.6	72.7 ± 3.6	65.6 ± 4.2
NS GAN	6.8 ± 0.5	26.5 ± 1.6	58.5 ± 1.9	55.0 ± 3.3
LSGAN	$7.8 \pm 0.6^*$	30.7 ± 2.2	87.1 ± 47.5	$53.9 \pm 2.8^*$
WGAN	6.7 ± 0.4	21.5 ± 1.6	55.2 ± 2.3	41.3 ± 2.0
WGAN GP	20.3 ± 5.0	24.5 ± 2.1	55.8 ± 0.9	30.0 ± 1.0
DRAGAN	7.6 ± 0.4	27.7 ± 1.2	69.8 ± 2.0	42.3 ± 3.0
BEGAN	13.1 ± 1.0	22.9 ± 0.9	71.4 ± 1.6	38.9 ± 0.9
VAE	23.8 ± 0.6	58.7 ± 1.2	155.7 ± 11.6	85.7 ± 3.8

Budget matters



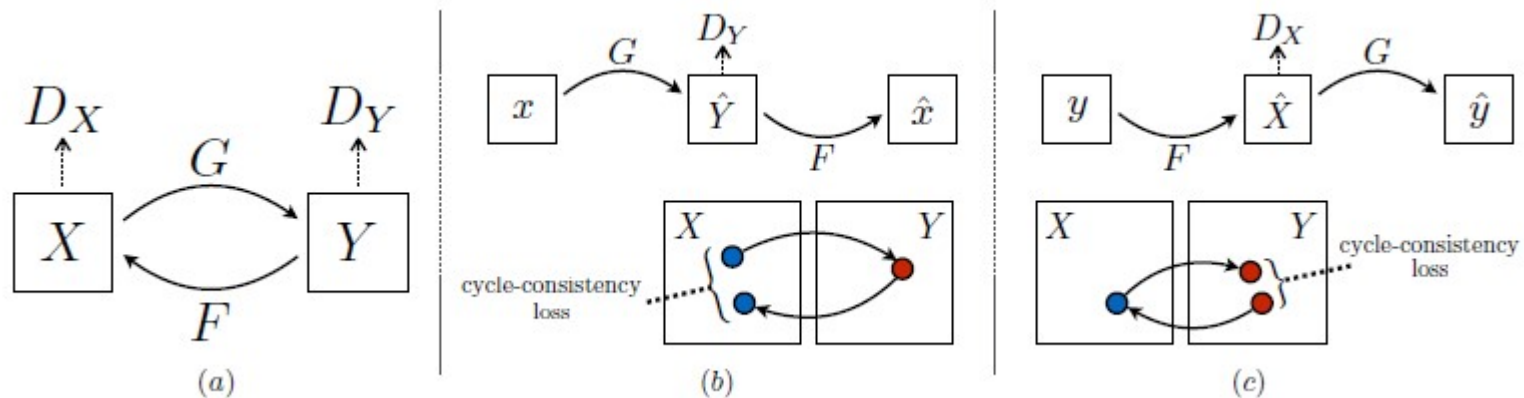
Combining VAEs and GANs

- VEEGAN
 - Combine likelihood-based and likelihood-free models
 - variational inference with synthetic likelihoods
- IntroVAE
 - minimize the divergence of the approximate posterior with the prior for real data while maximizing it for the generated samples
 - generator model attempts to mislead the inference model by minimizing the divergence of the generated samples
- Adversarial Autoencoder
- Adversarial Variational Bayes
- ALI/BiGAN
- AlphaGAN

Unpaired Image-to-Image Translation with CycleGAN

- Unpaired data is cheap
- How to use unpaired data for paired image-to-image translation?
- Idea:
 - Capture special characteristics of one image collection and translate into another image collection
 - Cycle consistency
 - Define additional mapping from generated space to data space
 - Translator $G : X \rightarrow Y$ and translator $F : Y \rightarrow X$
 - F and G inverses of each other

CycleGAN



$$\|F(G(x)) - x\|_1 \quad \|G(F(y)) - y\|_1$$

$$\begin{aligned} \mathcal{L}_{\text{cyc}}(G, F) = & \mathbb{E}_{x \sim p_{\text{data}}(x)} [\|F(G(x)) - x\|_1] \\ & + \mathbb{E}_{y \sim p_{\text{data}}(y)} [\|G(F(y)) - y\|_1]. \end{aligned}$$

$$\begin{aligned} \mathcal{L}_{\text{GAN}}(G, D_Y, X, Y) = & \mathbb{E}_{y \sim p_{\text{data}}(y)} [\log D_Y(y)] \\ & + \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log(1 - D_Y(G(x)))], \end{aligned} \quad (1)$$

$$\begin{aligned} \mathcal{L}(G, F, D_X, D_Y) = & \mathcal{L}_{\text{GAN}}(G, D_Y, X, Y) \\ & + \mathcal{L}_{\text{GAN}}(F, D_X, Y, X) \\ & + \lambda \mathcal{L}_{\text{cyc}}(G, F), \end{aligned}$$

Style transfer



Figure: Zhu et al. *arXiv preprint*, 2017.

More applications

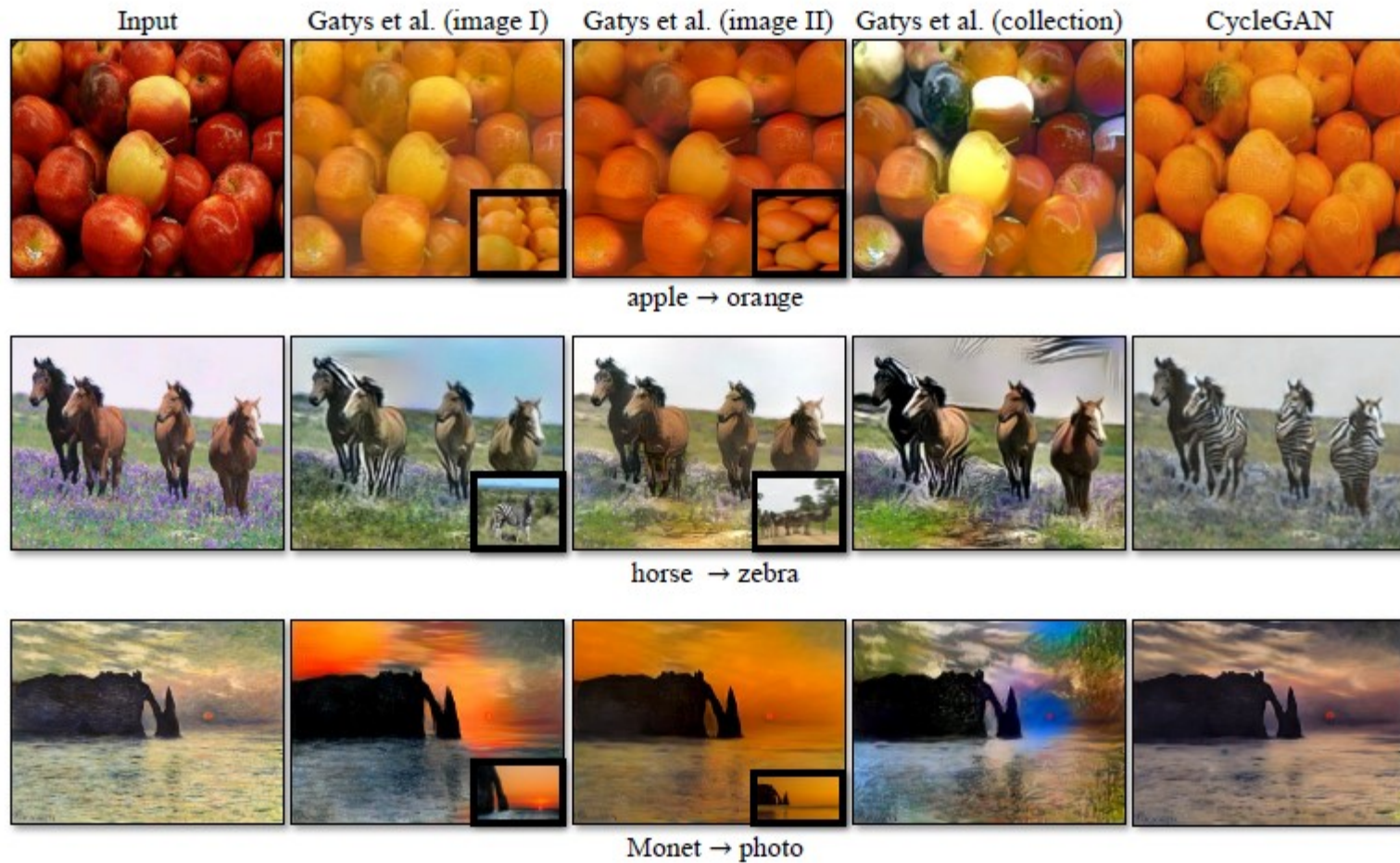


Figure: Zhu et al. *arXiv preprint*, 2017.

Failures

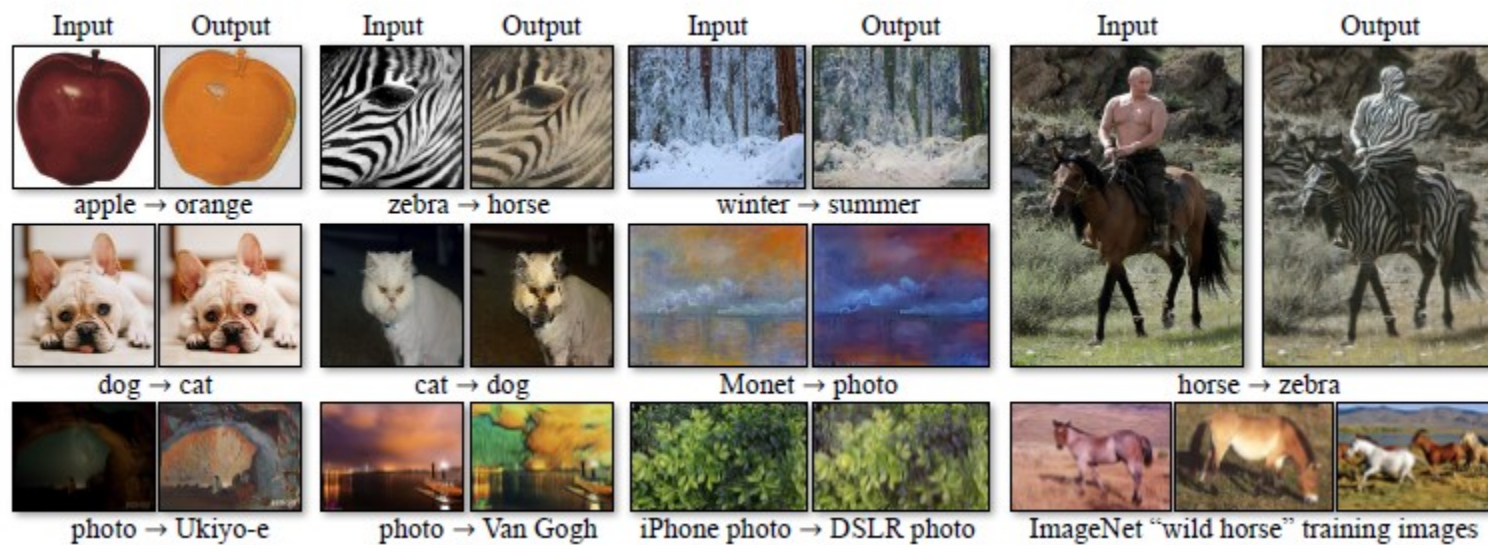


Figure: Zhu et al. *arXiv preprint*, 2017.

CycleGAN

- Excellent qualitative results on several tasks where paired training data does not exist, including collection style transfer, object transfiguration, season transfer, photo enhancement, etc.

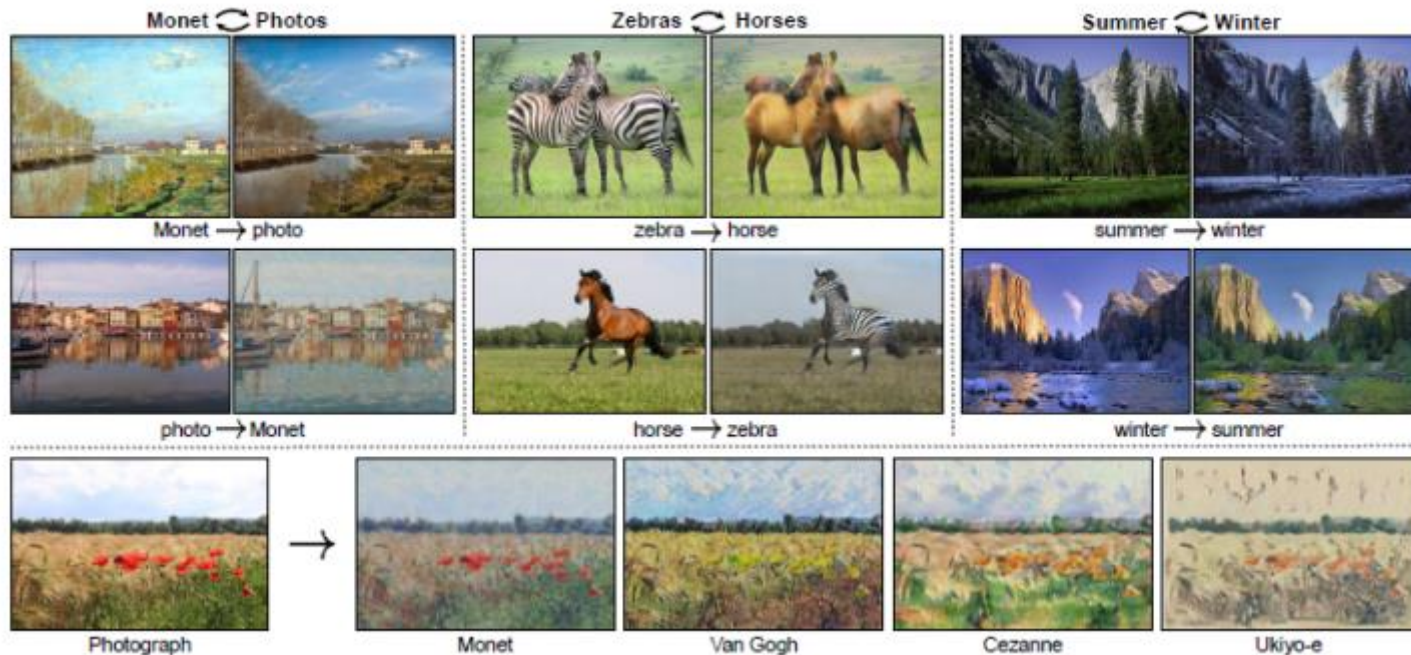


Figure: Zhu et al. *arXiv preprint*, 2017.

Conclusion

- Various types of models
 - Autoregressive models,
 - Explicit density model, optimizes exact likelihood, good samples.
Slow
 - VAEs
 - Optimises lower bound on likelihood. Useful representation and inference queries. Blurry samples.
 - GANs
 - Game-theoretic approach, best samples. Tricky and unstable to train
- Large variation between datasets, no one-size-fits-all model
- Lots of open research question and ongoing research