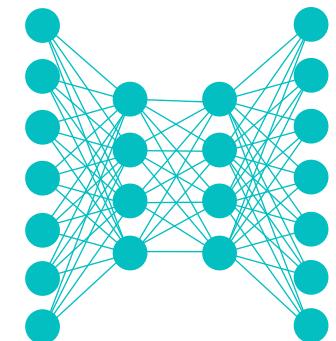


Lecture Notes for
Neural Networks
and Machine Learning



Generative
Networks



Logistics and Agenda

- Logistics
 - Welcome to online learning!
- Agenda
 - A historical perspective of generative Neural Networks
 - Variational Auto-Encoding
 - VAE in Keras Demo (next time)

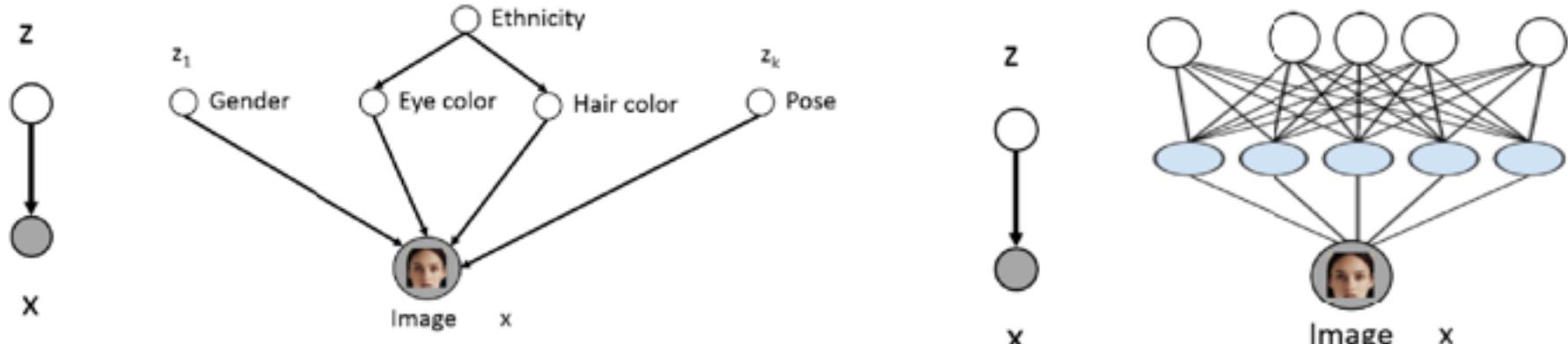


Last Time

- Multi-modal and Multi-task
- Town Hall had three options for what to work on
 - ChemBL
 - DeepFake
 - Biometrics
- Now: Generative Networks



Motivations: Latent Variables



$$p(\mathbf{x} | \mathbf{z})$$

Hard: \mathbf{z} is expertly chosen

$$p(\mathbf{x} | \mathbf{z})$$

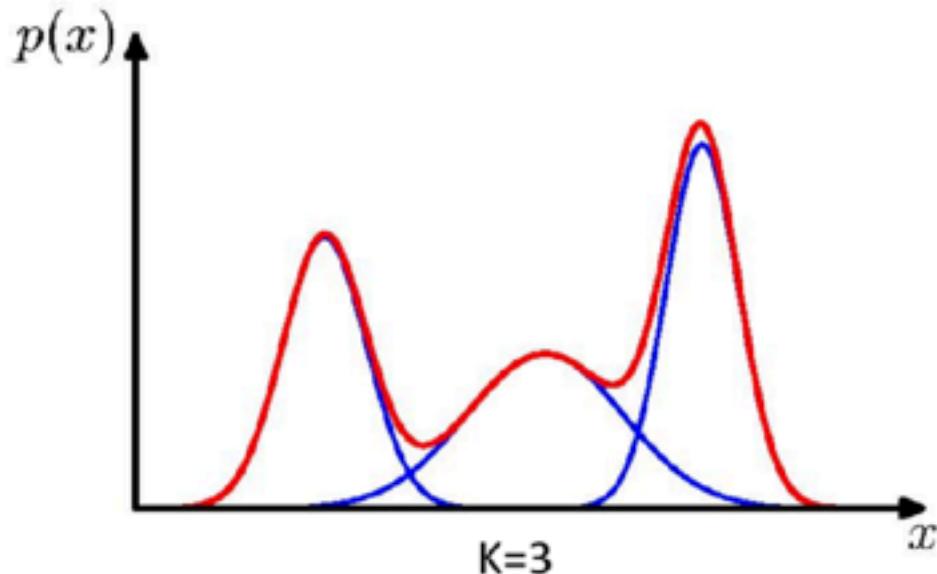
Not as Hard: \mathbf{z} is trained,
latent variables are uncontrolled

Want: $p(\mathbf{x}) \approx \sum_{\mathbf{z}} p(\mathbf{z}) p_{\theta}(\mathbf{x} | \mathbf{z})$



Motivation: Mixtures for Simplicity

Want: $p(\mathbf{x}) \approx \sum_{\mathbf{z}} p(\mathbf{z})p_{\theta}(\mathbf{x} | \mathbf{z})$



- Each latent variable mostly independent of other latent variables
- The sum of various mixtures can approximate most any distribution
- Good choice for conditional is Normal Distribution
- Can parameterize $p(x|z)$ to be a Neural Network

$$p_{\theta}(\mathbf{x} | \mathbf{z} = k) = \mathcal{N}(\mathbf{x}, \mu_k, \Sigma_k)$$

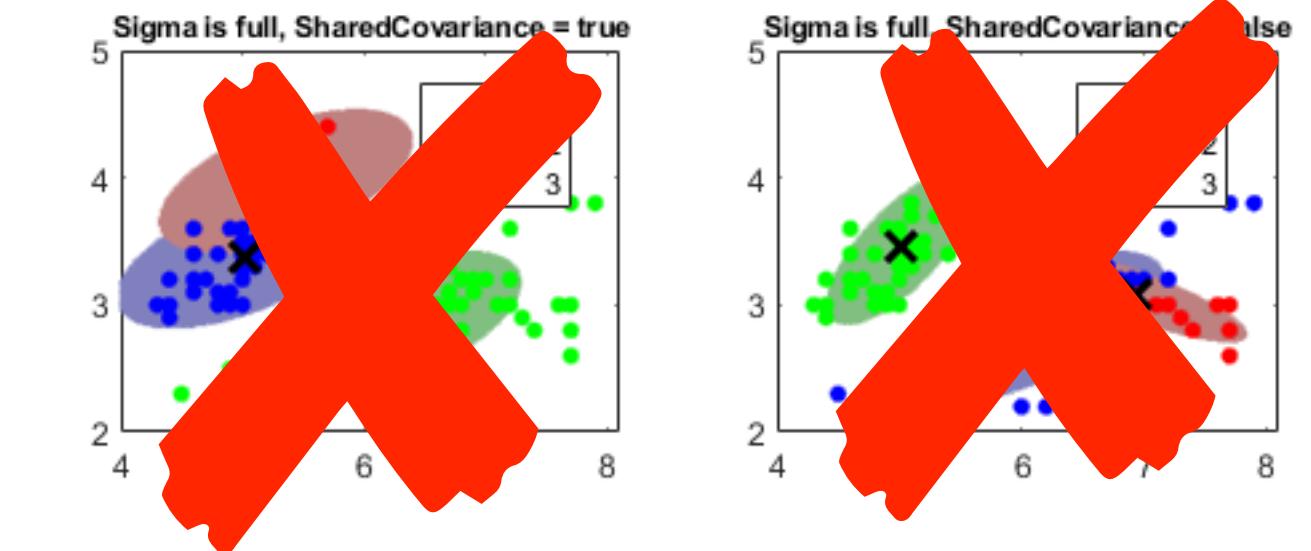
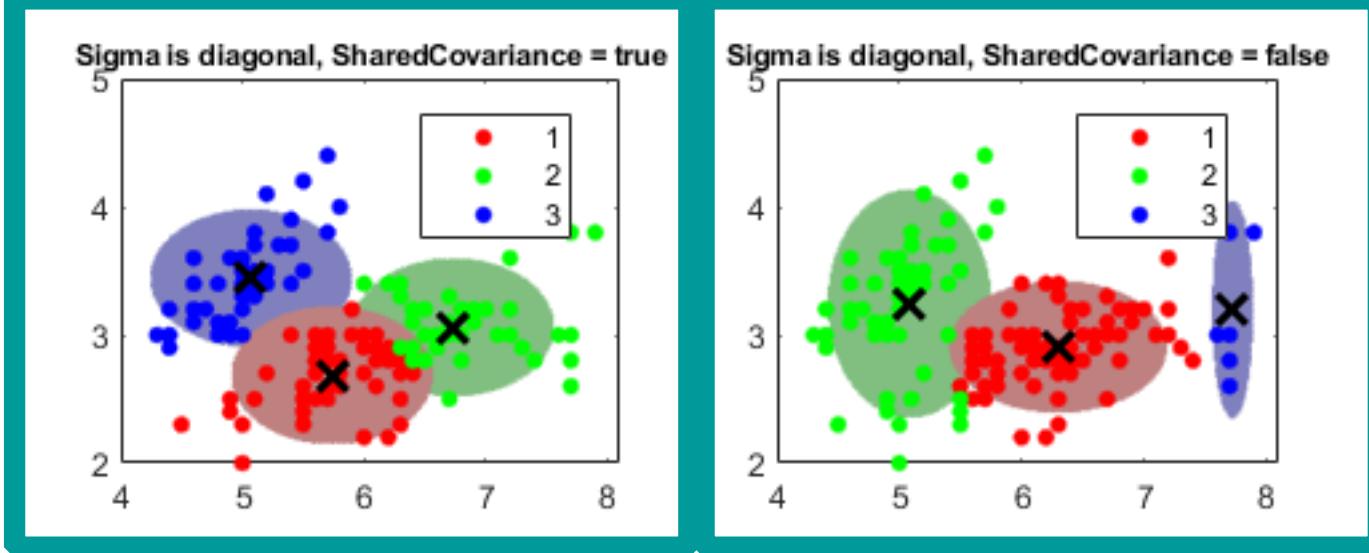
mean and covariance learned



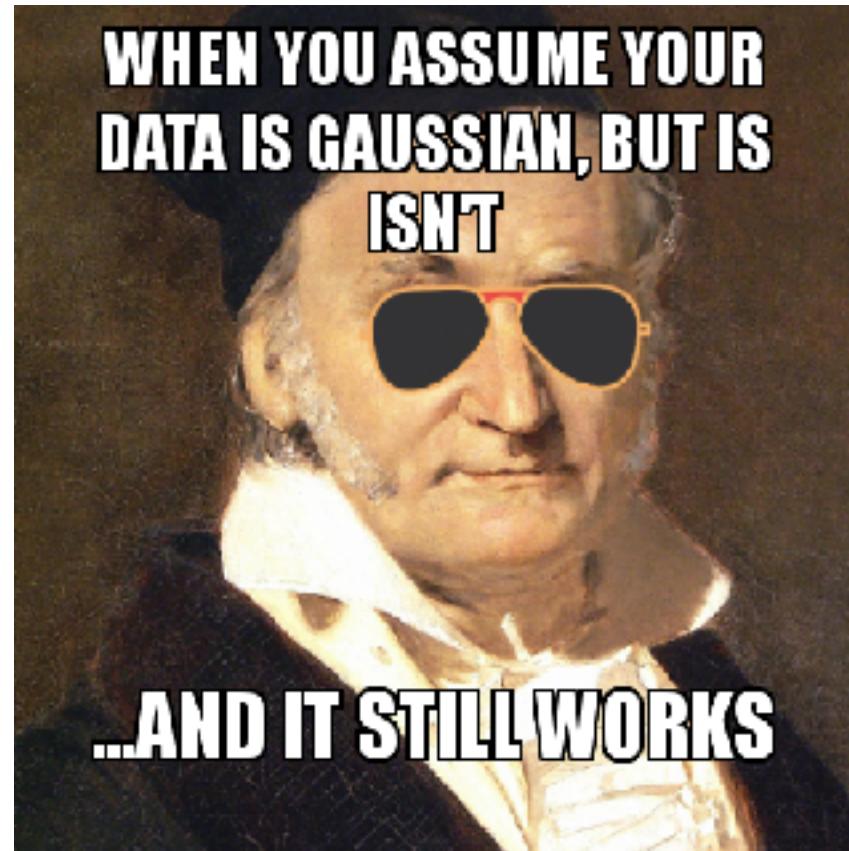
Motivation: Mixtures for Simplicity

$$= \mathcal{N}(\mathbf{x}, \mu_k, \Sigma_k)$$

mean and covariance learned



A History of Generative Networks



Aside: Sampling and expectation

$$\mathbf{E}_{s \leftarrow q(s|x)}[f(\cdot)] = \sum_{\forall i}^{\text{some function}} q(s|x^{(i)}) \cdot f(x^{(i)})^{\text{could be neural network}}$$

Expected value of f under conditional distribution, q
 s is latent variable, $x^{(i)}$ is an observation

$$\mathbf{E}_{s \leftarrow q(s|x)}[\log f(\cdot)] = \sum_{\forall i} q(s|x^{(i)}) \cdot \log(f(x^{(i)}))$$

If function is a probability, this is just the negative of cross entropy of distributions:

$$H(q, p) = - \sum_x q(x) \cdot \log(p(x))$$

Recall that KL divergence is a measure of difference in two distribution, and is just:

$$D(p\|q) = \sum_x p(x) \cdot \log\left(\frac{p(x)}{q(x)}\right) = \mathbf{E}_p \left[\log\left(\frac{p(x)}{q(x)}\right) \right]$$



Taxonomy of Generative Models

Taxonomy of Generative Models

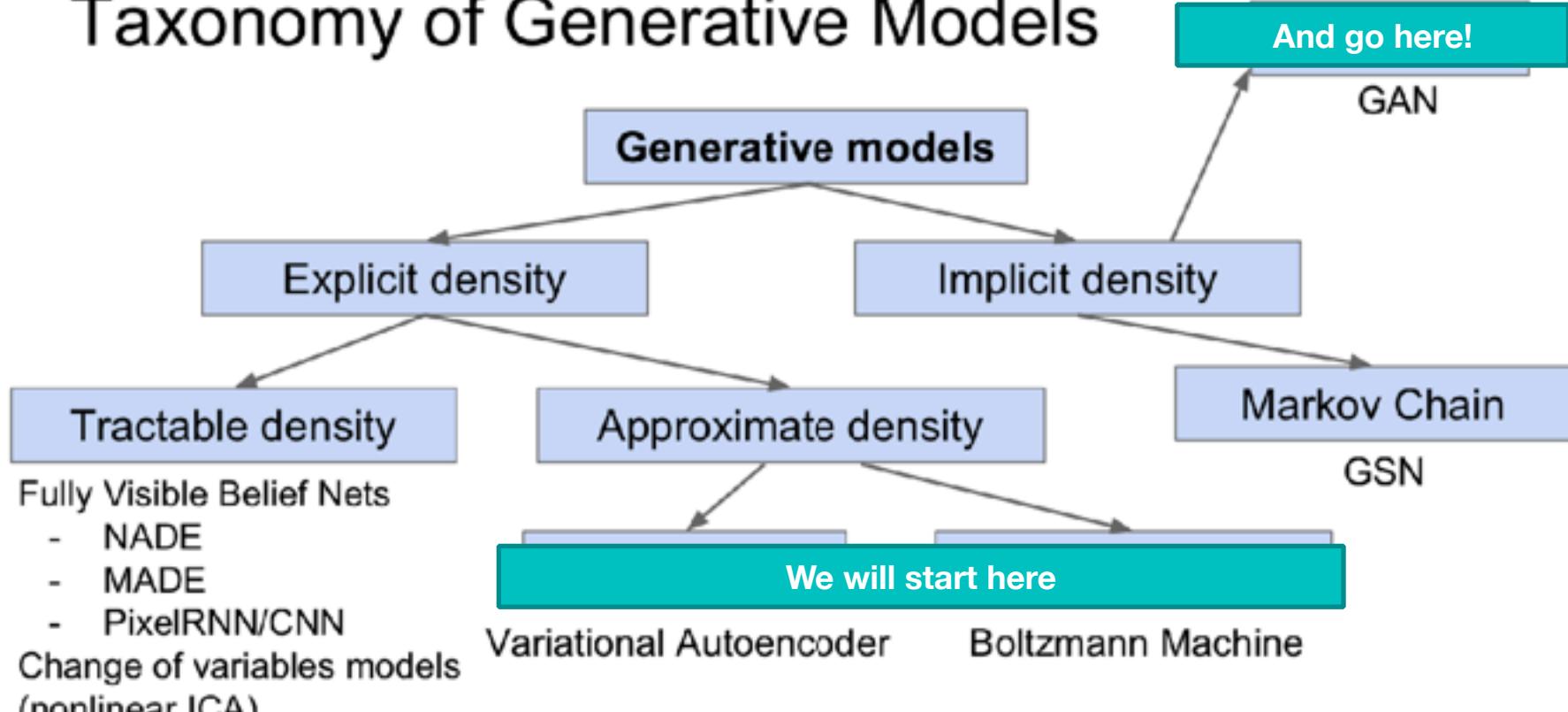
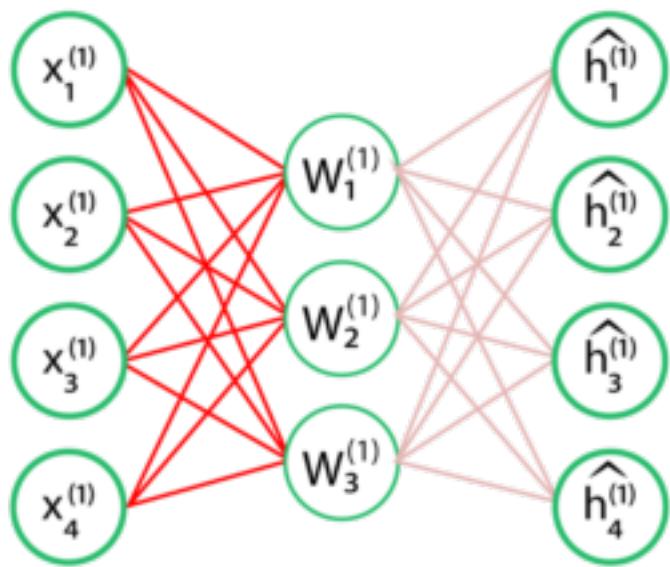
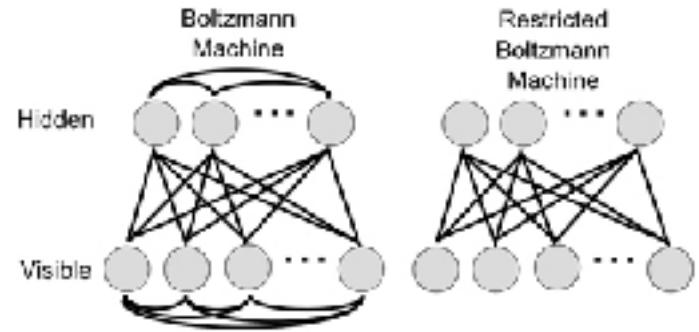


Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

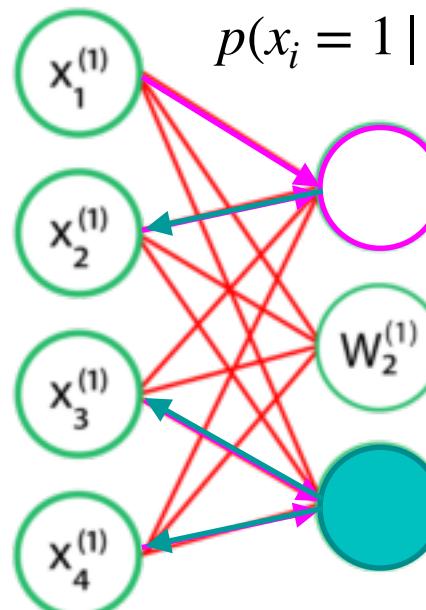


Deep Generative Models, 2006

- Restricted Boltzmann Machine
 - Forward pass (visible to latent)
 - Backward pass (latent to visible)
 - Similar to an auto encoder



AUTOENCODERS

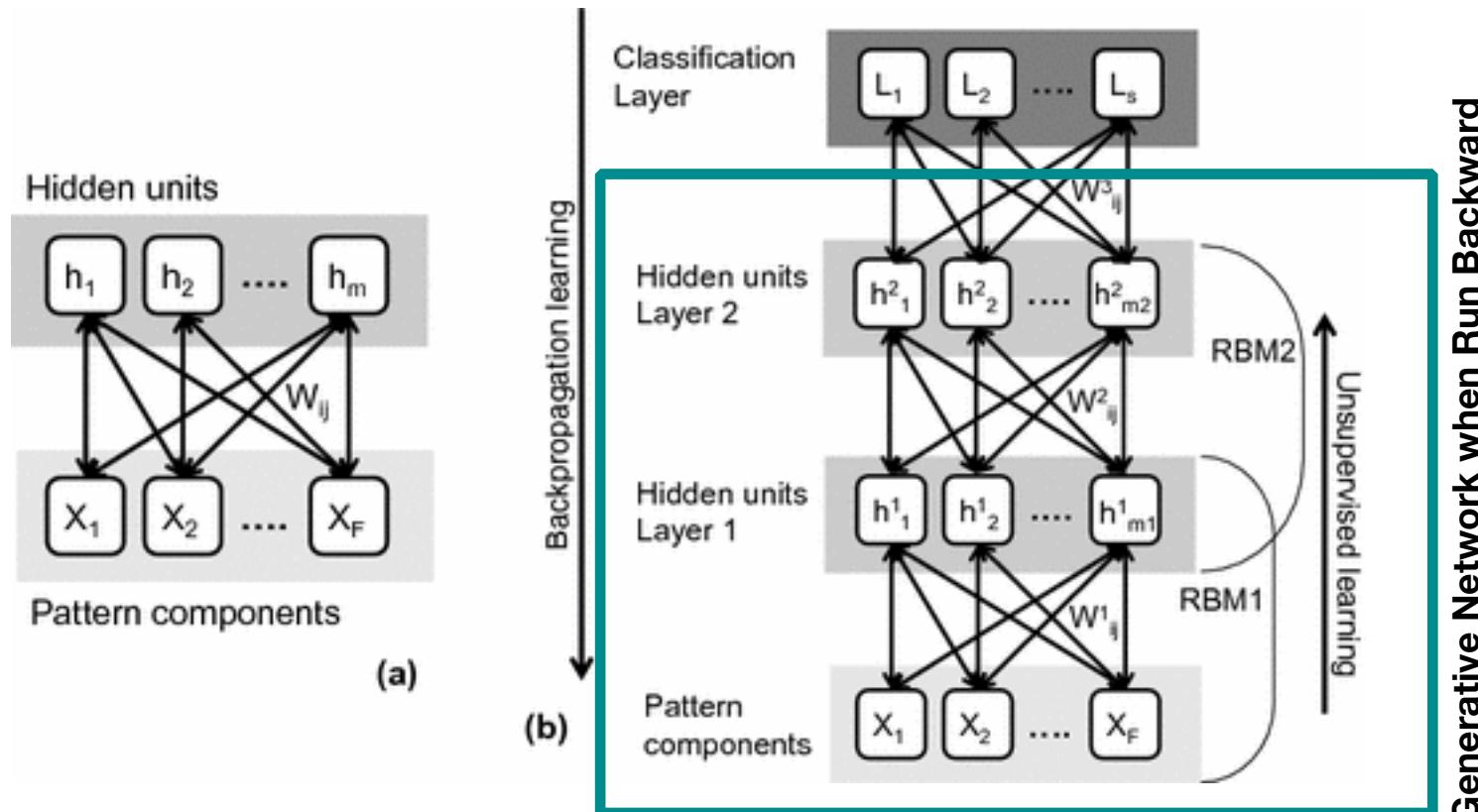


RBMs



Deep Generative Models, 2006

- Deep Belief Network
 - Many RBM blocks together!



https://link.springer.com/chapter/10.1007/978-3-319-59153-7_3

11



Deep Belief Networks, 2006

- Deep Belief Networks (DBNs)
- Iterative Layer Training (not feed forward)

Definitions for Restricted Boltzman Machine

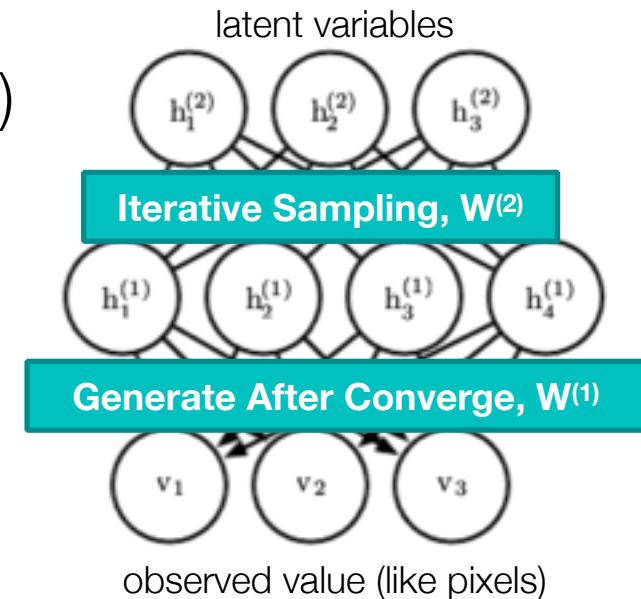
$$P(\mathbf{h}^{(l)}, \mathbf{h}^{(l-1)}) \propto \exp \left(\mathbf{b}^{(l)\top} \mathbf{h}^{(l)} + \mathbf{b}^{(l-1)\top} \mathbf{h}^{(l-1)} + \mathbf{h}^{(l-1)\top} \mathbf{W}^{(l)} \mathbf{h}^{(l)} \right), \quad (20.17)$$

Energy Function, relates each latent variable (binary)

$$P(h_i^{(k)} = 1 | \mathbf{h}^{(k+1)}) = \sigma \left(b_i^{(k)} + \mathbf{W}_{:,i}^{(k+1)\top} \mathbf{h}^{(k+1)} \right) \forall i, \forall k \in 1, \dots, l-2, \quad (20.18)$$

$$P(v_i = 1 | \mathbf{h}^{(1)}) = \sigma \left(b_i^{(0)} + \mathbf{W}_{:,i}^{(1)\top} \mathbf{h}^{(1)} \right) \forall i. \quad (20.19)$$

Sigmoids relate conditional distributions (like always)



$\mathbf{E}_{v \leftarrow p_{obs}} [\log p(v)]$ Maximize!! want samples from p to match observed v
 $p(v)$ is intractable

$\mathbf{E}_{v \leftarrow p} \left[\mathbf{E}_{h^{(1)} \leftarrow p(h^{(1)}|v)} [\log p(h^{(1)})] \right]$
optimize with contrastive divergence
i.e., approximation of EM, not Back Prop



Generative Models, 2009

- Deep Boltzmann Machine

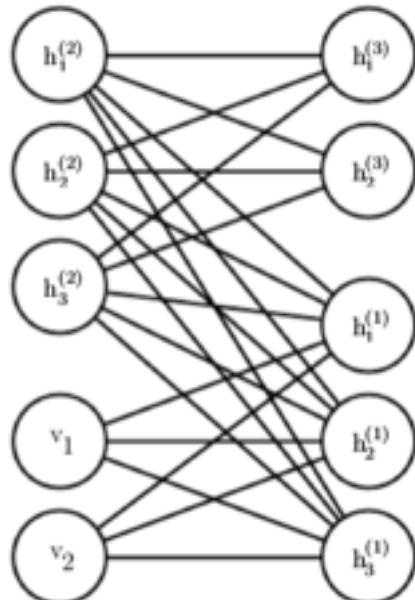
$$P(\mathbf{v}, \mathbf{h}^{(1)}, \mathbf{h}^{(2)}, \mathbf{h}^{(3)}) = \frac{1}{Z(\boldsymbol{\theta})} \exp(-E(\mathbf{v}, \mathbf{h}^{(1)}, \mathbf{h}^{(2)}, \mathbf{h}^{(3)}; \boldsymbol{\theta})). \quad (20.24)$$

To simplify our presentation, we omit the bias parameters below. The DBM energy function is then defined as follows:

$$E(\mathbf{v}, \mathbf{h}^{(1)}, \mathbf{h}^{(2)}, \mathbf{h}^{(3)}; \boldsymbol{\theta}) = -\mathbf{v}^\top \mathbf{W}^{(1)} \mathbf{h}^{(1)} - \mathbf{h}^{(1)\top} \mathbf{W}^{(2)} \mathbf{h}^{(2)} - \mathbf{h}^{(2)\top} \mathbf{W}^{(3)} \mathbf{h}^{(3)}. \quad (20.25)$$

We now develop the mean field approach for the example with two hidden layers. Let $Q(\mathbf{h}^{(1)}, \mathbf{h}^{(2)} | \mathbf{v})$ be the approximation of $P(\mathbf{h}^{(1)}, \mathbf{h}^{(2)} | \mathbf{v})$. The mean field assumption implies that

$$Q(\mathbf{h}^{(1)}, \mathbf{h}^{(2)} | \mathbf{v}) = \prod_j Q(h_j^{(1)} | \mathbf{v}) \prod_k Q(h_k^{(2)} | \mathbf{v}). \quad (20.29)$$



Not tractable: Can only optimize the Evidence lower bound, ELBO

The mean field approach is to minimize

Approximate via MCMC
like **Gibbs Sampling**

$$\text{KL}(Q \| P) = \sum_h Q(\mathbf{h}^{(1)}, \mathbf{h}^{(2)} | \mathbf{v}) \log \left(\frac{Q(\mathbf{h}^{(1)}, \mathbf{h}^{(2)} | \mathbf{v})}{P(\mathbf{h}^{(1)}, \mathbf{h}^{(2)} | \mathbf{v})} \right). \quad (20.30)$$



Image Generation from Samples

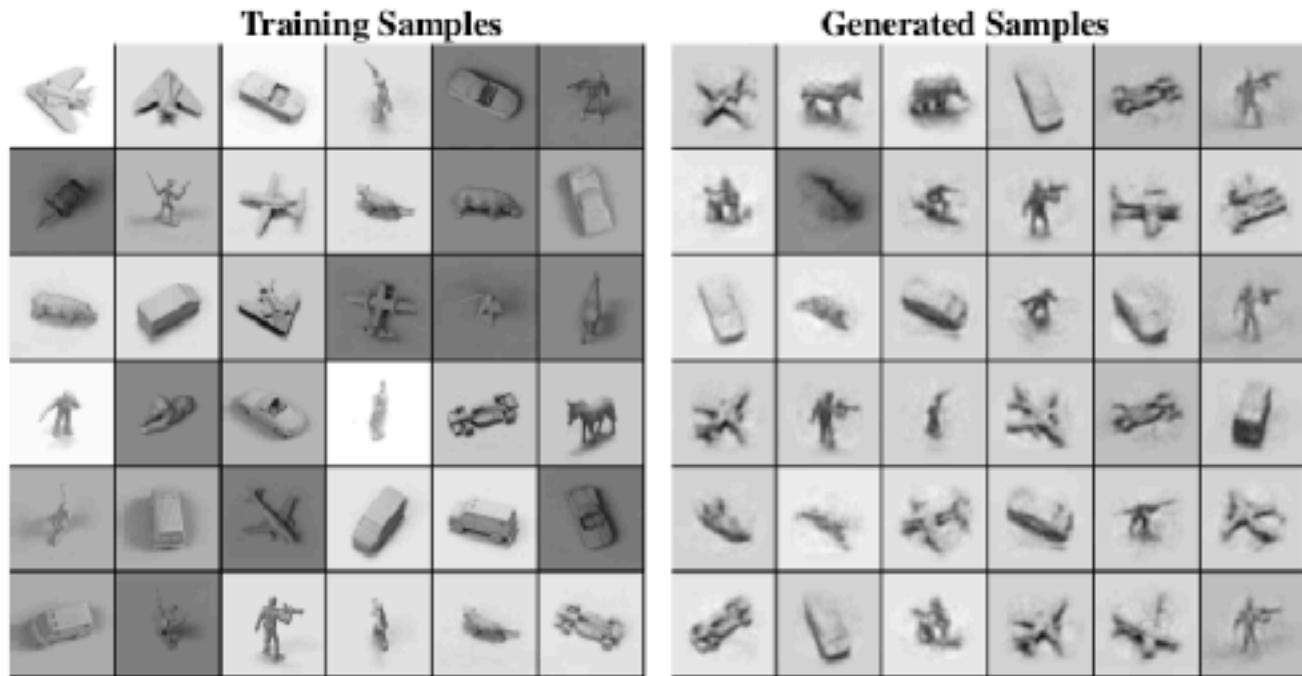
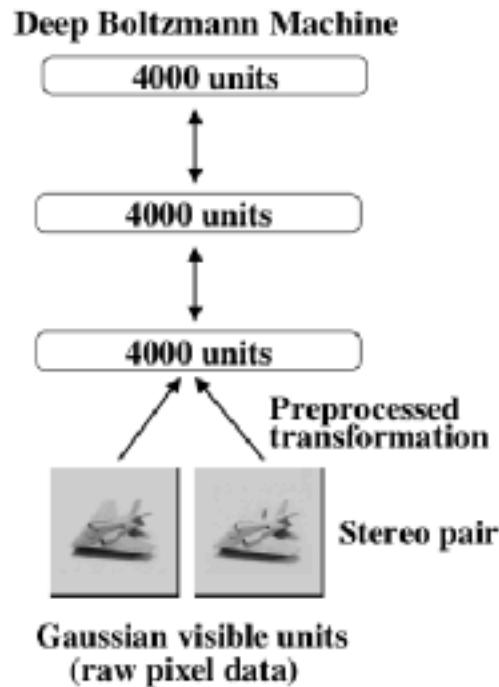


Figure 5: **Left:** The architecture of deep Boltzmann machine used for NORB. **Right:** Random samples from the training set, and samples generated from the deep Boltzmann machines by running the Gibbs sampler for 10,000 steps.



Contemporary Modeling

- DBNs and DBMs did not become very popular
 - Mathematics detracts from popular understanding
 - Often methods using sampling are not scalable
 - Cannot directly use Gradients (no Back Prop) 
- Evidence Lower Bound (ELBO) considered “good enough” because ... we can’t do better computationally
- Popular method for calculating generative networks with ELBO approximation:
 - Variational Auto Encoding
 - ◆ Guaranteed NOT to find global minimum
 - ◆ But scalable and will converge in finite time



Variational Auto Encoding

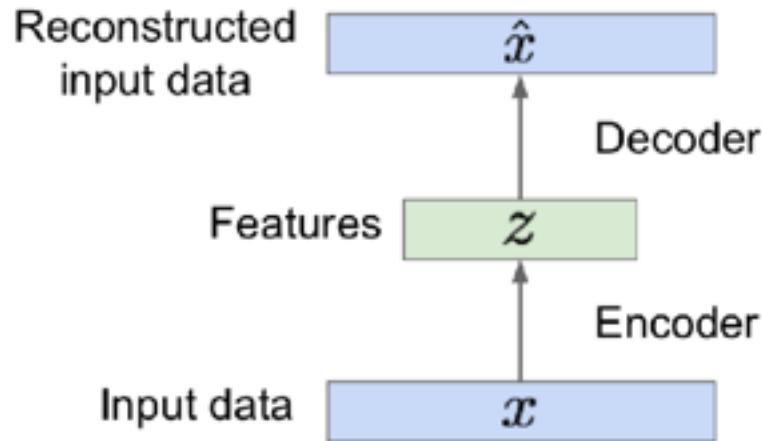
**“Mathematics is the
Khaleesi of sciences.”**



- Khal Friedrich Gauss



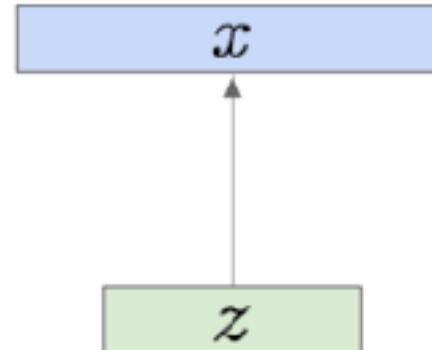
Aside: Auto Encoding



Once trained, is it possible to generate data?

Sample from
true conditional
 $p_{\theta^*}(x | z^{(i)})$

Sample from
true prior
 $p_{\theta^*}(z)$



- Does this work for simple auto encoding?
 - Nope.
- Learned space is not continuous
- How to sample from the latent space?
- Features could be highly correlated
- Need to define some constraints on the latent space...



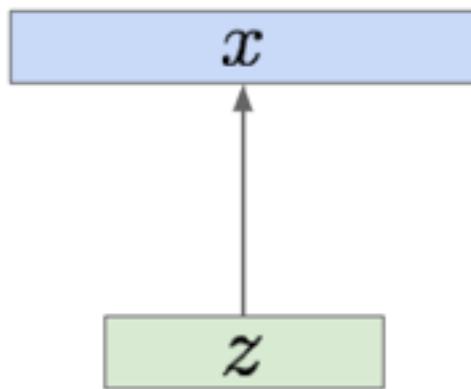
Reasonable constraints for $p(z)$?

Sample from
true conditional

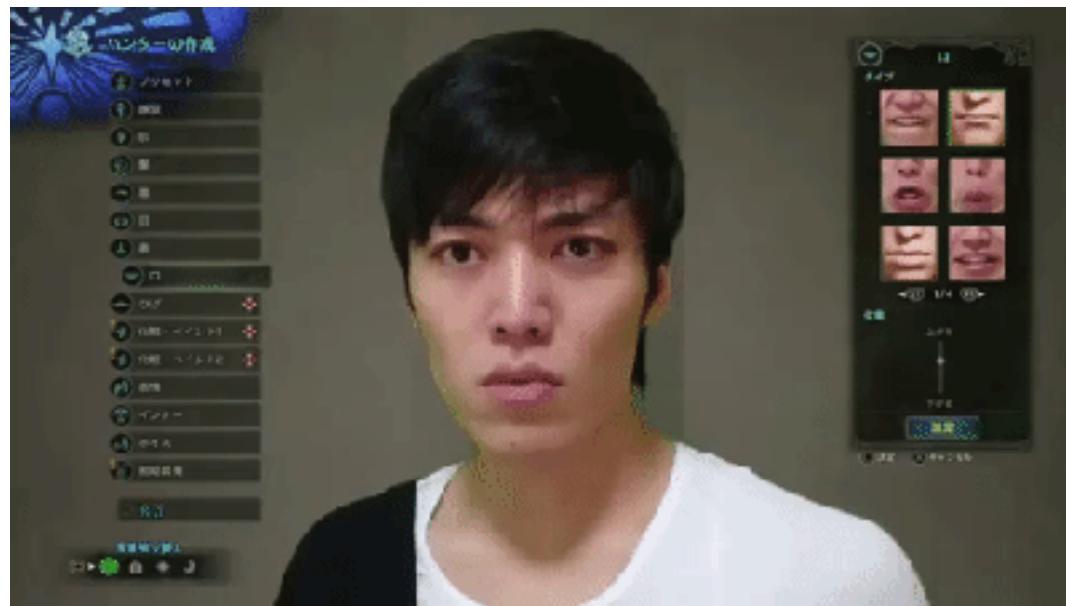
$$p_{\theta^*}(x \mid z^{(i)})$$

Sample from
true prior

$$p_{\theta^*}(z)$$



- Should be simple, easy to sample from: Normal
- Each component should be independent: Covariance
 - Encourages features that are semantic



Optimizing

$$p(z|x) = \frac{p(x|z)p(z)}{p(x)}$$

We need this inference in order to compute latent variable

$$p(x) = \int p(x|z)p(z)dz$$

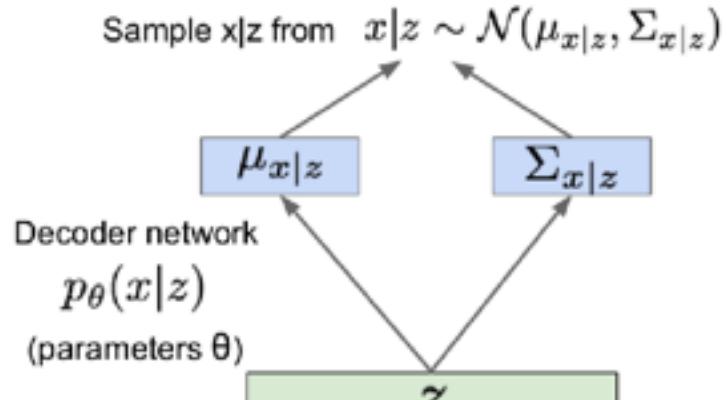
Denominator is of this form

- We can't compute! **Intractable computation** for all “ z ”
- So let's define this with **variational inference**:
 - Only needs to work for z **with observed** $x^{(i)}$
 - 1. **Encode** observed $x^{(i)}$ using network $q(z|x^{(i)})$,
 - 2. Use $q(z|x^{(i)})$ to sample z appropriately, then **decode** with another neural network
 - 3. Make $q(z|x^{(i)})$ as close as possible to $p(z|x)$

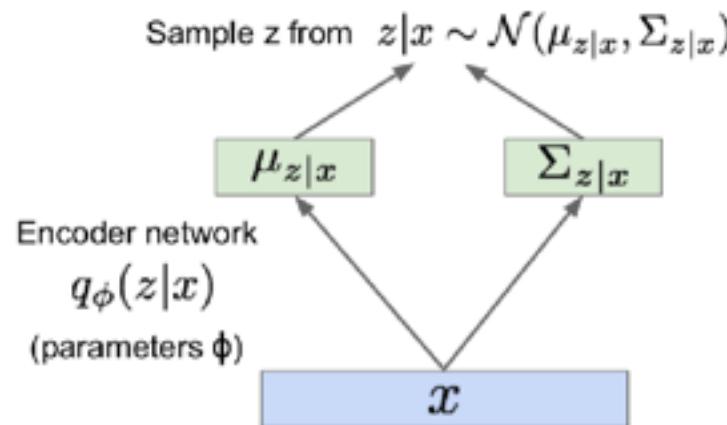


Need a new formulation

Step Two: Decode



Step One: Encode



Step Three: Make p and q Similar

$$D(q(z|x^{(i)}) \| p(z|x^{(i)})) = \mathbf{E}_{q(z|x)} \left[\log \left(\frac{q(z|x^{(i)})}{p(z|x^{(i)})} \right) \right]$$

we can manipulate this equation to yield the following:

$$\log p(x^{(i)}) = \mathbf{E}_{z \leftarrow q_\phi(z|x)} [\log p(x^{(i)})]$$

intuition MLE: maximize probability of observed $x^{(i)}$

Output of network, q , are the mean and covariance for sampling a variable z

$$\mu(x^{(i)}) \text{ and } \Sigma(x^{(i)})$$



Need a new formulation

$$= \mathbf{E}_{\mathbf{z} \leftarrow q(z|x)} [\log p(x^{(i)})]$$

$$= \mathbf{E}_q \left[\log \frac{p(x^{(i)}|z)p(z)}{p(z|x^{(i)})} \frac{q(z|x^{(i)})}{q(z|x^{(i)})} \right] \quad \text{Bayes rule + multiply by one}$$

$$= \mathbf{E}_q [\log p(x^{(i)}|z)] + \mathbf{E}_q \left[\log \frac{p(z)}{q(z|x^{(i)})} \right] + \mathbf{E}_q \left[\log \frac{q(z|x^{(i)})}{p(z|x^{(i)})} \right]$$
$$= \mathbf{E}_q [\log p(x^{(i)}|z)] - \mathbf{E}_q \left[\log \frac{q(z|x^{(i)})}{p(z)} \right] + \mathbf{E}_q \left[\log \frac{q(z|x^{(i)})}{p(z|x^{(i)})} \right]$$

$$= \mathbf{E}_q [\log p(x^{(i)}|z)] - D_{KL} [q(z|x^{(i)})||p(z)] + D_{KL} [q(z|x^{(i)})||p(z|x^{(i)})]$$

always negative

$$\log p(x^{(i)}) \geq \mathbf{E}_q [\log p(x^{(i)}|z)] - D_{KL} [q(z|x^{(i)})||p(z)] \quad \text{Maximize Lower Bound}$$

What have we really done here? Is this still MLE?



The Loss Function

Maximize through
Error of Reconstruction
This is just negative cross entropy

assume $p(z)$ is Normal,
zero mean and std=1
because we like Normal

$$\begin{aligned} D_{KL}[N(\mu(X), \Sigma(X)) \| N(0, 1)] &= \frac{1}{2} (\text{tr}(\Sigma(X)) + \mu(X)^T \mu(X) - k - \log \det(\Sigma(X))) \\ &\quad \text{outputs of q network} \\ D_{KL}[N(\mu(X), \Sigma(X)) \| N(0, 1)] &= \frac{1}{2} \left(\sum_k \Sigma(X) + \sum_k \mu^2(X) - \sum_k 1 - \log \prod_k \Sigma(X) \right) \\ &= \frac{1}{2} \left(\sum_k \Sigma(X) + \sum_k \mu^2(X) - \sum_k 1 - \sum_k \log \Sigma(X) \right) \\ &\geq \mathbf{E}_{q(z|x^{(i)})} [\log p(x^{(i)}|z)] - D_{KL} [q(z|x^{(i)}) \| p(z)] \\ &= \frac{1}{2} \sum_k (\Sigma(X) + \mu^2(X) - 1 - \log \Sigma(X)) \end{aligned}$$



The Covariance Output

$$\geq \mathbf{E}_{q(z|x^{(i)})} [\log p(x^{(i)} | z)] - D_{KL} [q(z | x^{(i)}) || p(z)]$$

Maximize through

Error of Reconstruction

This is just negative cross entropy

assume $p(z)$ is Normal,

zero mean and std=1

because we like Normal

$$= \frac{1}{2} \sum_k (\Sigma(X) + \mu^2(X) - 1 - \log \Sigma(X))$$

covariance is not numerically stable because of underflow

$$D_{KL}[N(\mu(X), \Sigma(X)) \| N(0, 1)] = \frac{1}{2} \sum_k \left(\exp(\Sigma(X)) + \mu^2(X) - 1 - \frac{\Sigma(X)}{\log \text{var}} \right)$$

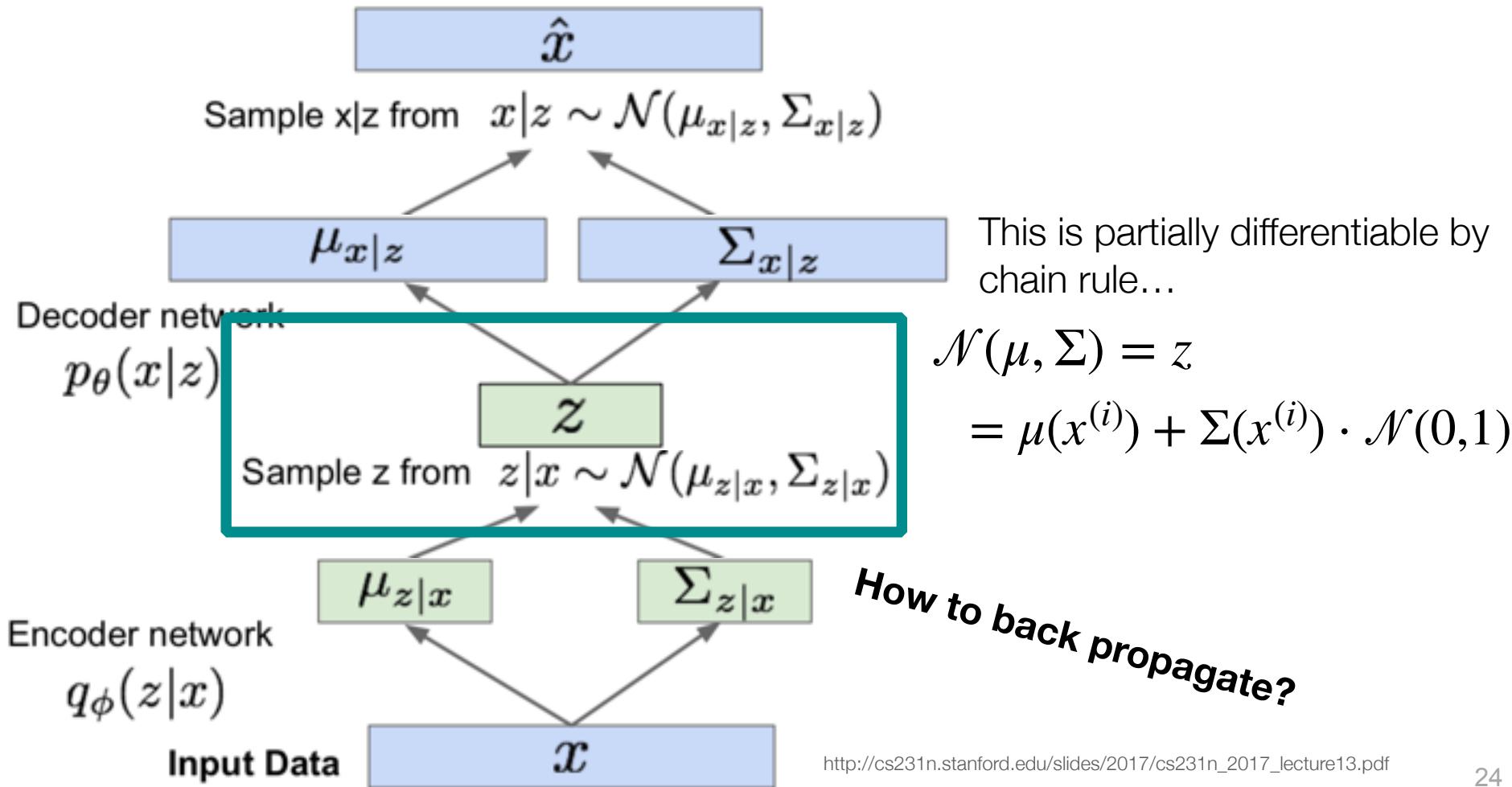
so we will have the neural network output log variance

Also, remember we assume **diagonal covariance**, so z 's are not correlated
This means covariance is only a vector of variances (the diagonal)



Back Propagating

$$\geq \mathbf{E}_{q(z|x^{(i)})} [\log p(x^{(i)}|z)] - D_{KL} [q(z|x^{(i)})||p(z)]$$



The Loss Function Implementation

```
# Encode the input into a mean and variance parameter
z_mean, z_log_variance = encoder(input_img)

# Draw a latent point using a small random epsilon
z = z_mean + exp(z_log_variance) * epsilon

# Then decode z back to an image
reconstructed_img = decoder(z)

# Instantiate a model
model = Model(input_img, reconstructed_img)

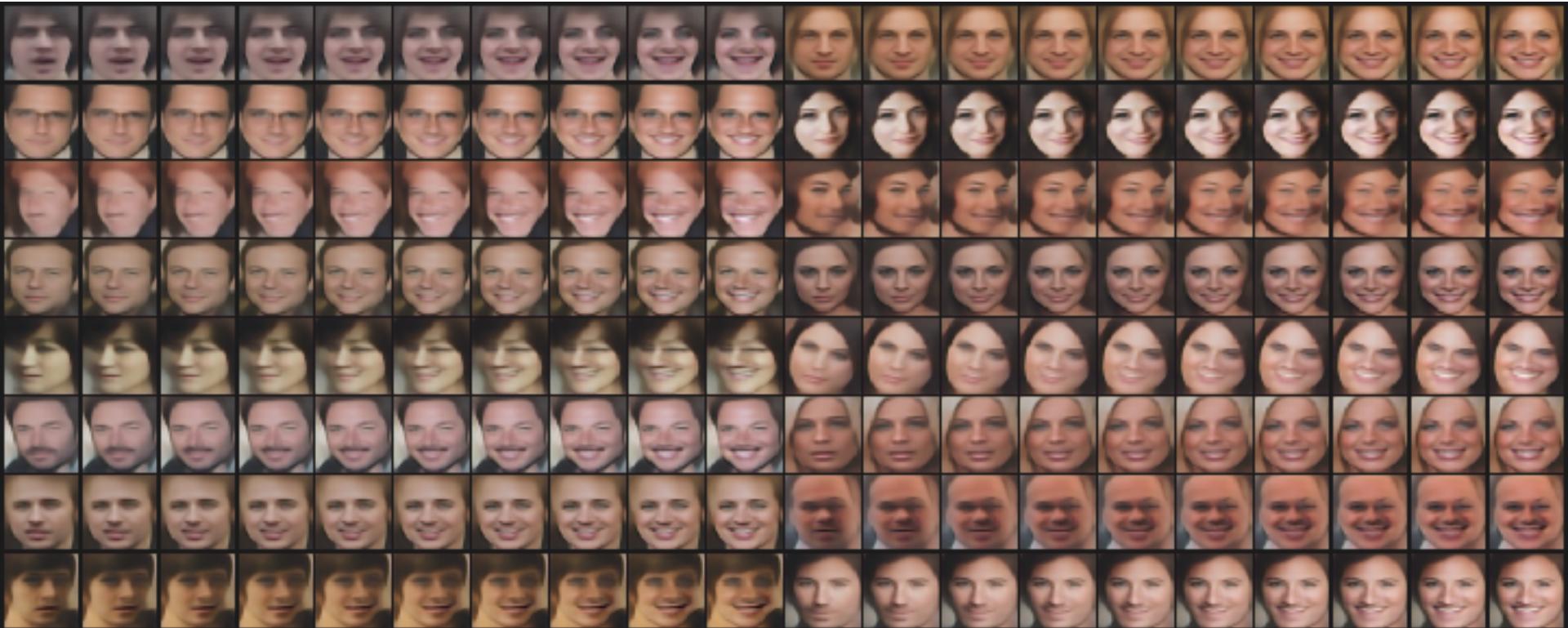
def vae_loss(self, x, z_decoded):
    x = K.flatten(x)
    z_decoded = K.flatten(z_decoded)
    xent_loss = keras.metrics.binary_crossentropy(x, z_decoded)
    kl_loss = -5e-4 * K.mean(
        1 + z_log_var - K.square(z_mean) - K.exp(z_log_var), axis=-1)
    return K.mean(xent_loss + kl_loss)
```

$$\begin{aligned}\mathcal{N}(\mu, \Sigma) &= z \\ &= \mu(x^{(i)}) + \Sigma(x^{(i)}) \cdot \mathcal{N}(0,1)\end{aligned}$$

$$\begin{aligned}&\mathbf{E}_{q(z|x^{(i)})} [\log p(x^{(i)}|z)] - D_{KL} [q(z|x^{(i)})||p(z)] \\&= \mathbf{E}_{q(z|x^{(i)})} [\log p(x^{(i)}|z)] - \sum_k \exp(\Sigma(x^{(i)})) + \mu(x^{(i)})^2 - 1 - \Sigma(x^{(i)})\end{aligned}$$



VAE Examples

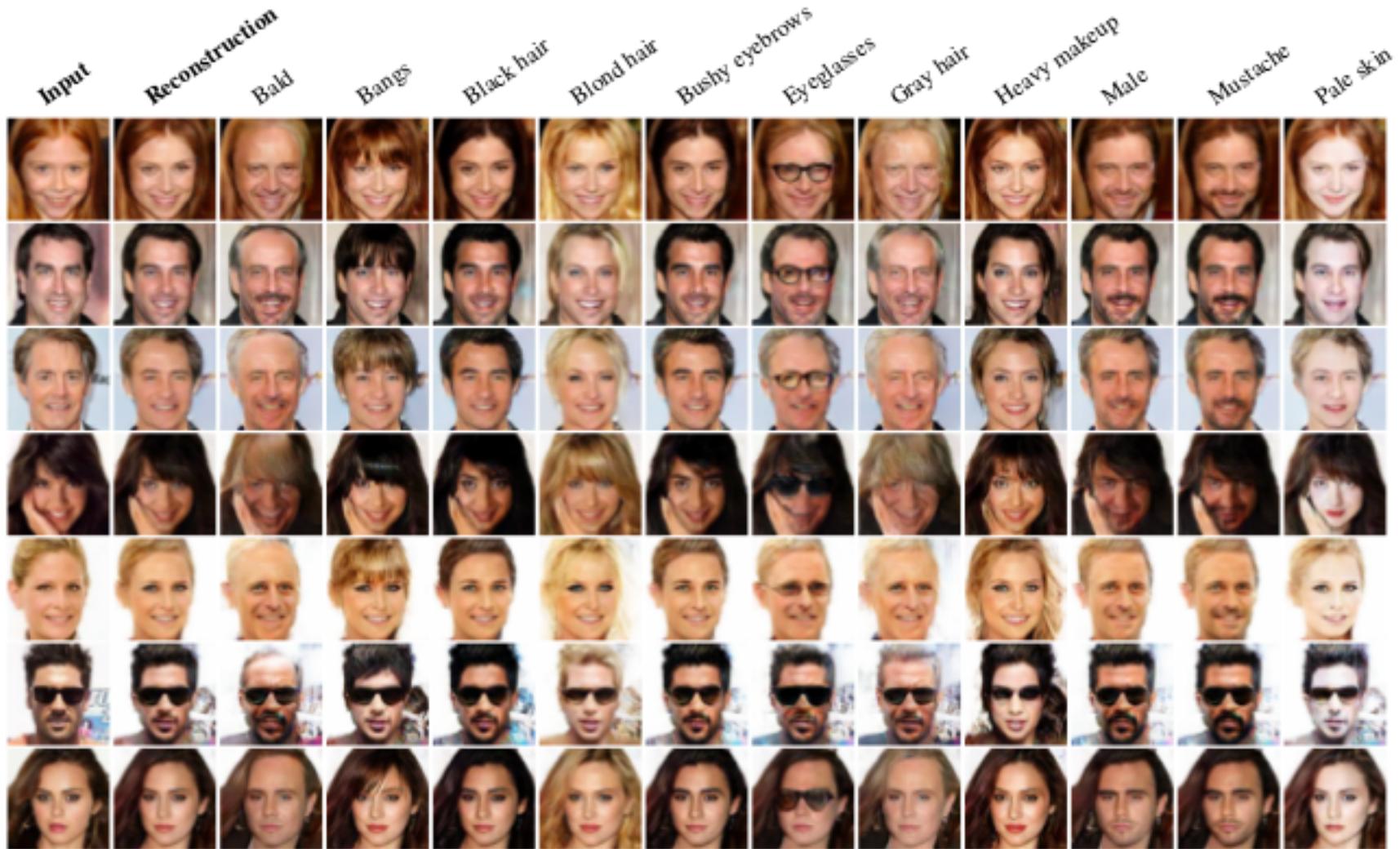


Encoding faces, then finding the “ z ” that relates to smiling.



VAE Examples

Different, automatically found z , latent variables



Lecture Notes for **Neural Networks** **and Machine Learning**

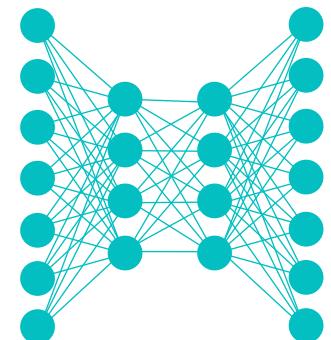
Generative Networks



Next Time:

GANs

Reading: Chollet CH8

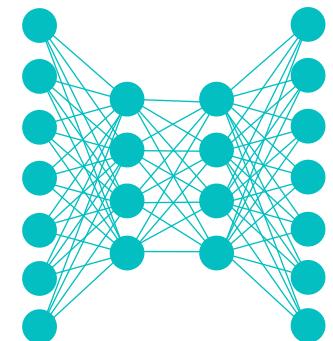




Lecture Notes for
Neural Networks
and Machine Learning



VAE Demo
Generative Adversarial
Networks



Logistics and Agenda

- Logistics
 - Online Learning Reminders
 - How is the format?
 - How is lab three going (discuss deadline)?
- Agenda
 - VAE Demo
 - Simple Generative Adversarial Networks



Back to VAEs

Official ACM @TheOfficialACM

Yoshua Bengio, Geoffrey Hinton and Yann LeCun, the fathers of #DeepLearning, receive the 2018 #ACMTuringAward for conceptual and engineering breakthroughs that have made deep neural networks a critical component of computing today. bit.ly/2HVJldV



Yoshua Bengio



Geoffrey Hinton



Yann LeCun



Geoffrey Hinton @geoffreyhinton · 23h

Replying to @JeffDean @ylecun and @TheOfficialACM

Thanks to my graduate students and postdocs whose work won a Turing award. Thanks to my visionary mentors Inman Harvey, David Rumelhart and Terry Sejnowski. And thanks to Jeff Dean for creating the brain team that turns basic research in neural nets into game-changing products.

40

313

2,208



Yann LeCun @ylecun · 20h

Congratulations Geoff, from one of your former postdocs ;-)

3

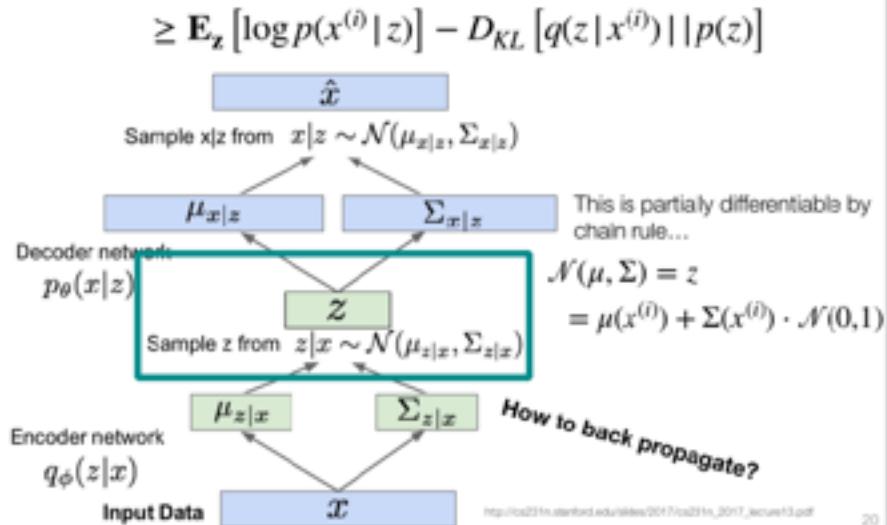
11

312



Last Time

- Variational auto encoding



```
# Encode the input into a mean and variance parameter
z_mean, z_log_variance = encoder(input_img)

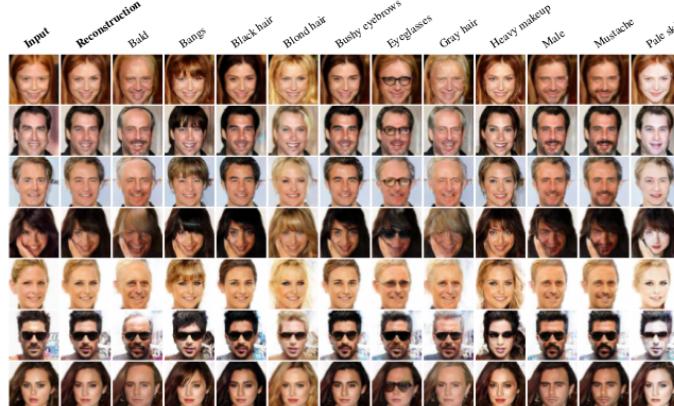
# Draw a latent point using a small random epsilon
z = z_mean + exp(z_log_variance) * epsilon

# Then decode z back to an image
reconstructed_img = decoder(z)

# Instantiate a model
model = Model(input_img, reconstructed_img)

def vae_loss(self, x, z_decoded):
    x = K.flatten(x)
    z_decoded = K.flatten(z_decoded)
    xent_loss = keras.metrics.binary_crossentropy(x, z_decoded)
    kl_loss = -5e-4 * K.mean(
        1 + z_log_var - K.square(z_mean) - K.exp(z_log_var), axis=-1)
    return K.mean(xent_loss + kl_loss)
```

$$\begin{aligned} & \mathbf{E}_z [\log p(x^{(i)} | z)] - D_{KL} [q(z | x^{(i)}) || p(z)] \\ &= \mathbf{E}_z [\log p(x^{(i)} | z)] - \sum \exp(\Sigma(x^{(i)})) + \mu(x^{(i)})^2 - 1 - \Sigma(x^{(i)}) \end{aligned}$$





VAEs in Keras

Implementation from Book on MNIST



Demo by Francois Chollet

Follow Along: <https://github.com/fchollet/deep-learning-with-python-notebooks/blob/master/8.4-generating-images-with-vae.ipynb>



The Simple GAN

Geoff Hinton after writing the paper on backprop in 1986



Latent Variable Approximation with GANS

- **Idea:** transform random input data into target of interest
 - Like an image!
- Rather than training an encoder with variational inference, we need a transformation protocol
 - Transformer will be a...
 - Neural Network (of course!)
- Transformer is a “complex” sampling method
- How to optimize and provide training examples?
 - Use two Neural Networks!
 - ... Neural Networks are like the chiropractic of the machine learning...



Generative Adversarial Network

- Generator: $x = g(z)$
- Discriminator: $\{0,1\} = d(x)$
- Mini-max, turn-based game:

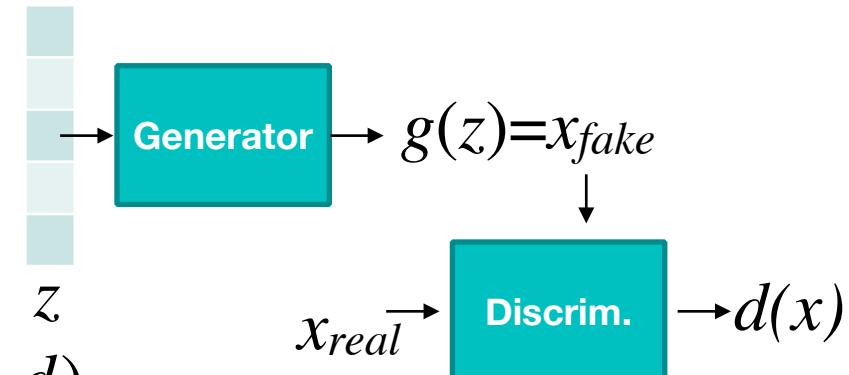
$$\hat{w} = \arg \min_g \max_d v(g, d)$$

- Nice differentiable choice for v (zero sum game):

$$v(g, d) = \mathbf{E}_{x \leftarrow p_{data}} [\log d(x_{real})] + \mathbf{E}_{x \leftarrow g(z)} [\log(1 - d(x_{fake}))]$$

only portion dependent on g
generator minimizes

everything dependent on d
discriminator maximizes



Taking turns

$$v(g, d) = \mathbf{E}_{x \leftarrow p_{data}} [\log d(x_{real})] + \mathbf{E}_{x \leftarrow g(z)} [\log(1 - d(x_{fake}))]$$

$$\hat{w} = \arg \min_g \max_d v(g, d)$$

- If only the problem was convex! This would be easy to solve...
- But $v(g,d)$ is not convex, not even close
 - Saddle points, saddle points everywhere
 - Like optimizing in an Ocean
- Taking turns on the gradient descent will probably never reach equilibrium
 - There is no convergence guarantee
 - But practically we like when discriminator == chance



Practical Objectives

- Discriminator objective, gradient ascent:

$$\max_d \mathbf{E}_{x \leftarrow p_{data}} [\log d(x_{real})] + \mathbf{E}_{x \leftarrow g(z)} [\log(1 - d(x_{fake}))]$$

- Generator objective, gradient descent:

$$\min_g \mathbf{E}_{x \leftarrow g(z)} [\log(1 - d(x))]$$

- But **practically** generator is really hard to train, amplified by a decreasing gradient
- Goodfellow tried to solve this mathematically through a number of different formulations
 - Nothing seemed to work, and then...



Practical Objectives

- Original Generator objective, gradient descent:

$$\min_g \mathbf{E}_{x \leftarrow g(z)} [\log(1 - d(x_{fake}))]$$

- Goodfellow *et al.* gave up on using rigorous mathematics (intractable) and relied on heuristic:
 - Instead of minimizing when discriminator is right
 - Why not maximize when its wrong?
 - ◆ not trying to win, just make the other person lose

$$\max_g \mathbf{E}_{x \leftarrow g(z)} [\log(d(x_{fake}))]$$

No longer a zero sum game?!



Final Loss Functions

- Discriminator:

$$\max_d \mathbf{E}_{x \leftarrow p_{data}} [\log d(x_{real})] + \mathbf{E}_{x \leftarrow g(z)} [\log(1 - d(x_{fake}))]$$

**Same as minimizing the binary cross entropy
of the model with: (1) labeled data and (2) generated data!**

- Generator:

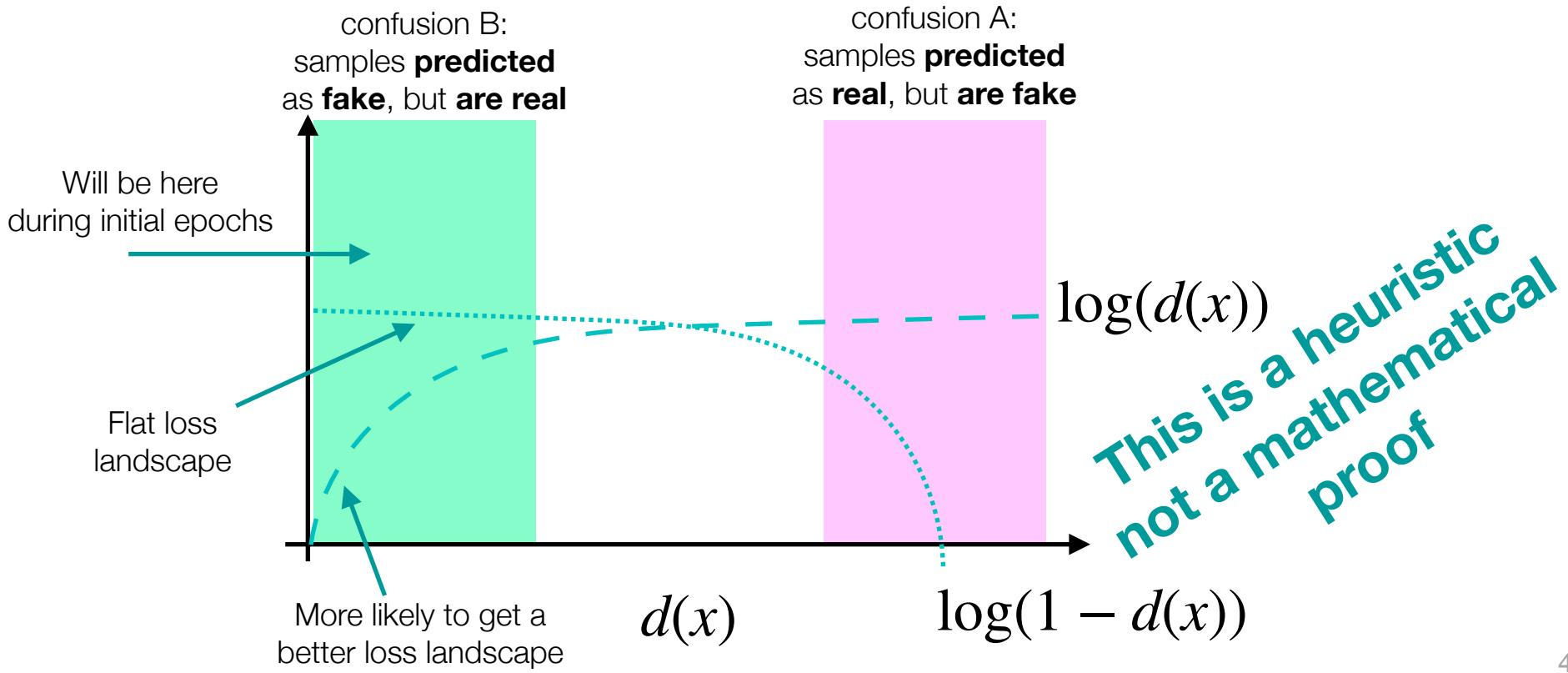
$$\max_g \mathbf{E}_{x \leftarrow g(z)} [\log(d(x_{fake}))]$$
 Freeze Discriminator Weights

**Same as minimizing the binary cross entropy
of “mislabeled” generated data!**



Why was minimizing incorrect?

- Training two networks is inherently unstable, need heuristic
- Let's assume generator is not any good for initial epochs!



How to Train your (dra) GAN

deeplearning.ai presents
Heroes of Deep Learning

Ian Goodfellow

Research Scientist at Google Brain

Apple



Every time Ian Goodfellow has a tutorial, It should be called:

“GAN-splaining with Ian”

—Geoffrey Hinton, Probably



Simple Training Approach

- Some caveats on why training will be difficult:
 - The latent space discovered by the Generator has almost no guarantees regarding continuity and structure (different than VAEs)
 - Convergence of GAN is not possible, only equilibrium
 - ◆ even saying discriminator is chance is not enough!
 - Each iteration through, the entire loss function changes because generator changes
 - ◆ Also we are sampling different points from generator...



Simple Training Approach

for number of training iterations **do**

for k steps **do**

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

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

end for

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

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

end for

Does this work?!

Not really! Why?

Goodfellow et al. “Generative Adversarial Networks” (NeurIPS 2014)



A bag of tricks from F. Chollet

The process of training GANs and tuning GAN implementations is notoriously difficult. There are a number of known tricks you should keep in mind. Like most things in deep learning, it's **more alchemy than science: these tricks are heuristics, not theory-backed guidelines**. They're supported by a level of intuitive understanding of the phenomenon at hand, and they're known to work well empirically, although not necessarily in every context.

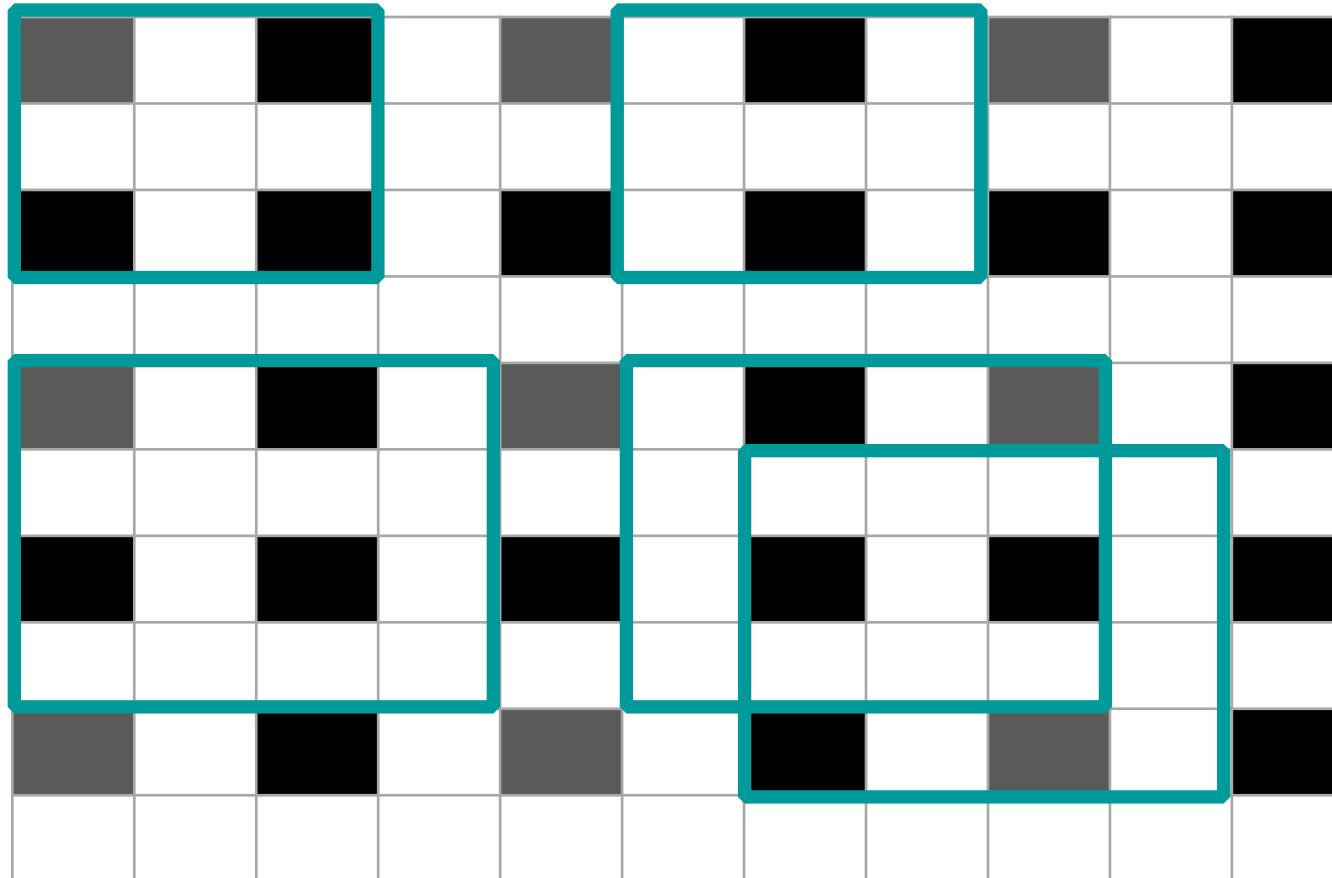
—Francois Chollet, DL with Python

- Use tanh for generator output, not a sigmoid
 - We typically normalize inputs to a NN to be from [-1, 1], generator needs to mirror this squashing
- Sample from Normal Distribution
 - Everyone else is doing it (even though no assumption of Gaussian in GAN formulation)
- Random is more robust in optimizer and labels
 - GANs get stuck a lot
- Sparse gradients are not your friend here
 - No max pooling, no ReLU 😞
- Make encoder and decoder symmetric
 - Unequal pixel coverage (checkerboards)



What do you mean checkerboard patterns?

- Kernel size should be a symmetric multiple of the stride



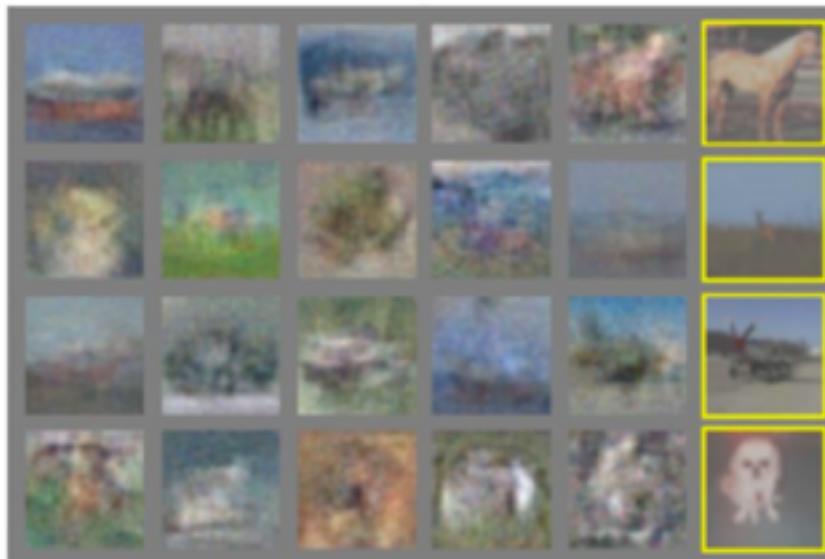
Bias needs to account for both when 2 pixels and four pixels are in the kernel space

Multiple of stride ensures that same number of active pixels are there.



Some Results of Generation

Goodfellow et al. "Generative Adversarial Networks" (NeurIPS 2014)

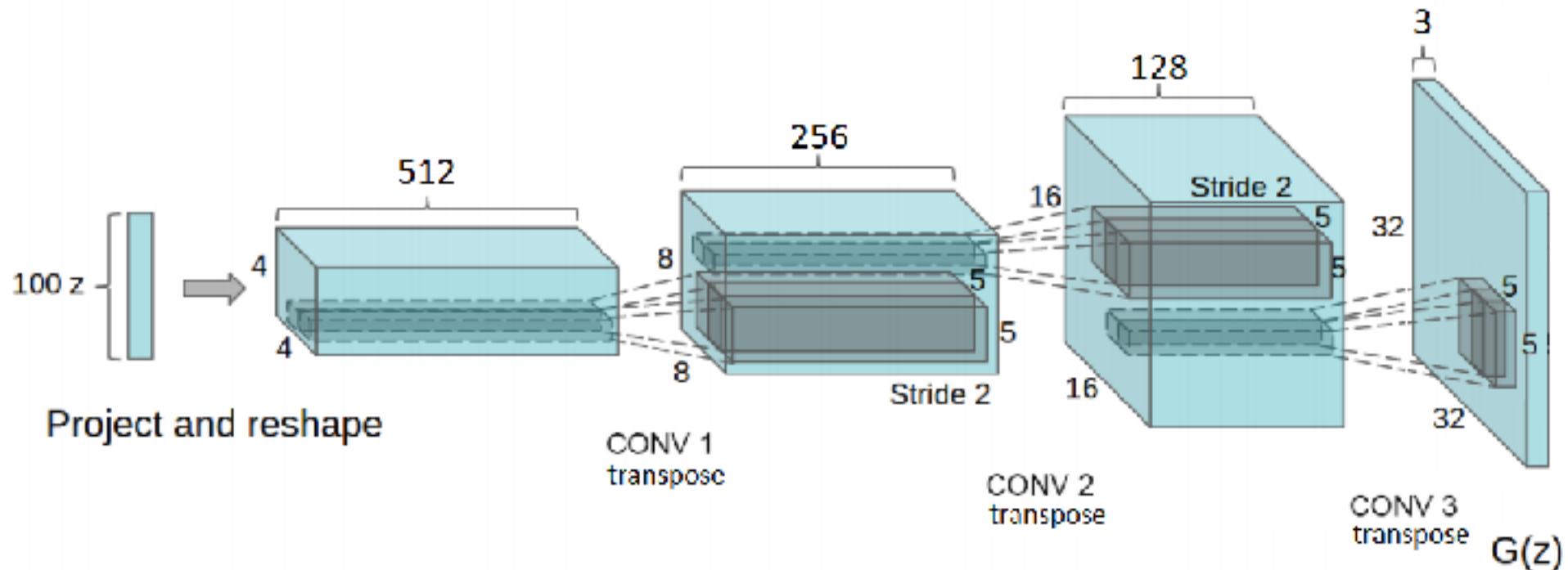


Highlight: Nearest Training Sample

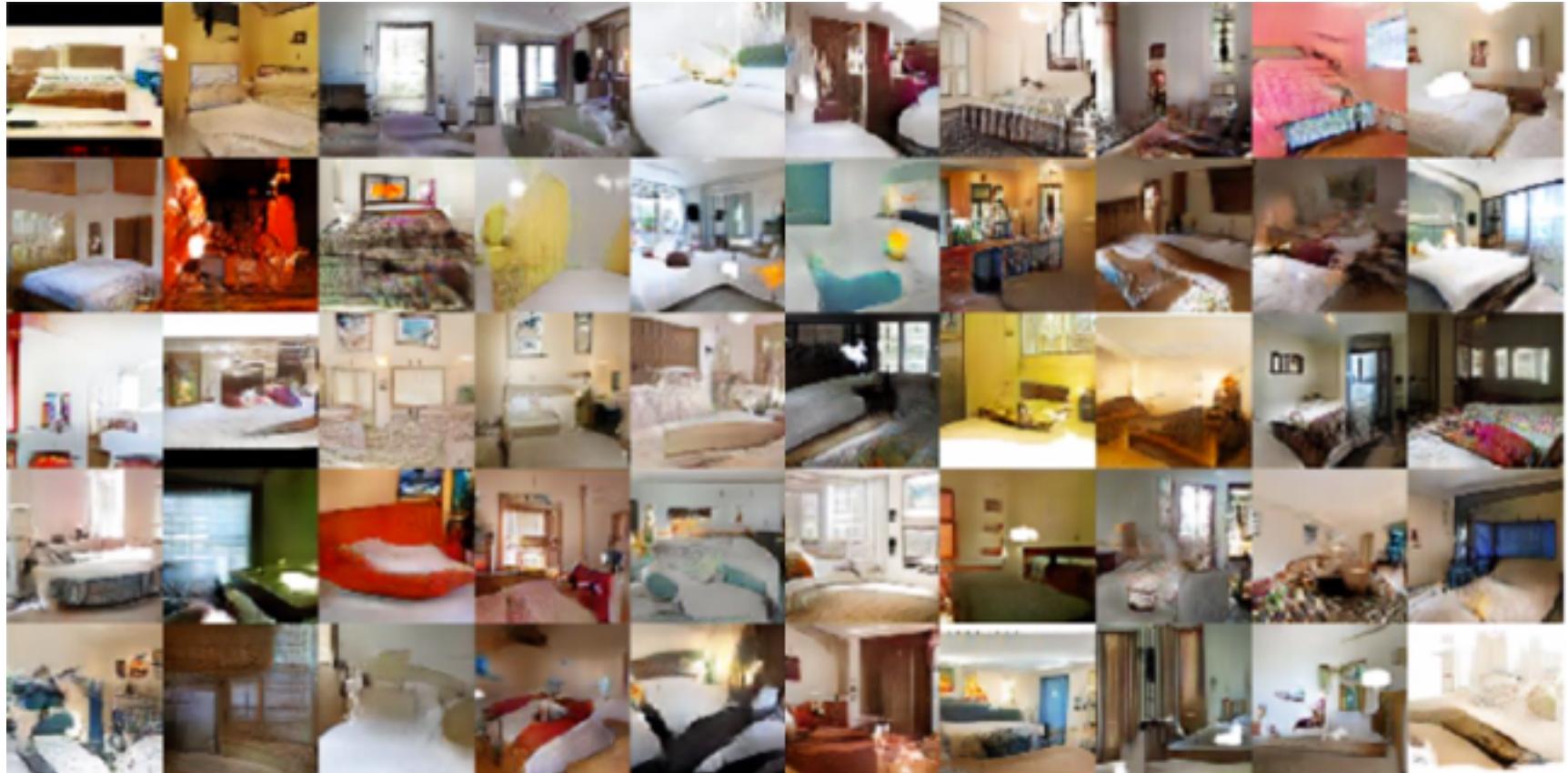


Deep Convolutional GANs

- Just need to reshape and use upsampling



Some Results of Generation



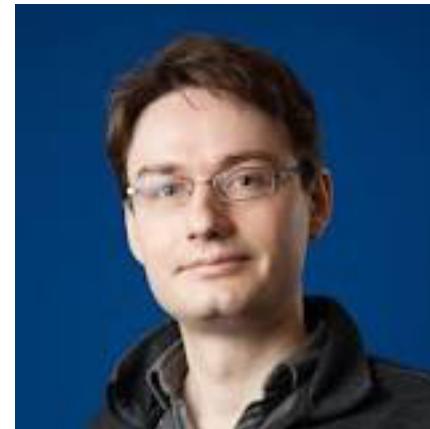
Radford, Alec, Luke Metz, and Soumith Chintala. "Unsupervised representation learning with deep convolutional generative adversarial networks." *arXiv preprint arXiv:1511.06434* (2015).





GANs in Keras

Implementation from Book on
“Frogs from CIFAR”



Demo by Francois Chollet

Follow Along: <https://github.com/fchollet/deep-learning-with-python-notebooks/blob/master/8.5-introduction-to-gans.ipynb>

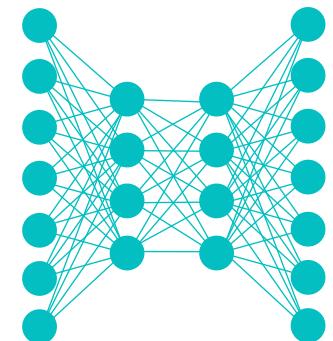


Lecture Notes for **Neural Networks** **and Machine Learning**

GANs

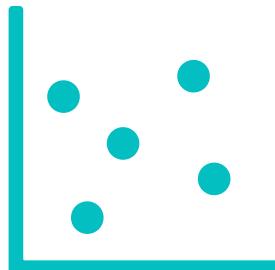


Next Time:
Practical GANs
Reading: Chollet CH8

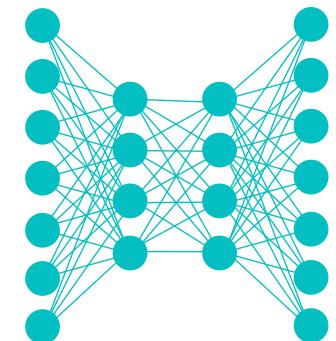




Lecture Notes for
Neural Networks
and Machine Learning



Practical
GANs



Logistics and Agenda

- Logistics
 - Student Presentations Today and Next Lecture
- Agenda
 - Review from Last Time, GAN Demo
 - Training GANs in different loss landscapes
 - **Student Presentation:** Representational Learning with Deep Convolutional GANs, Radford
 - Wasserstein GAN



Last Time

- Simple GANS

- Discriminator:

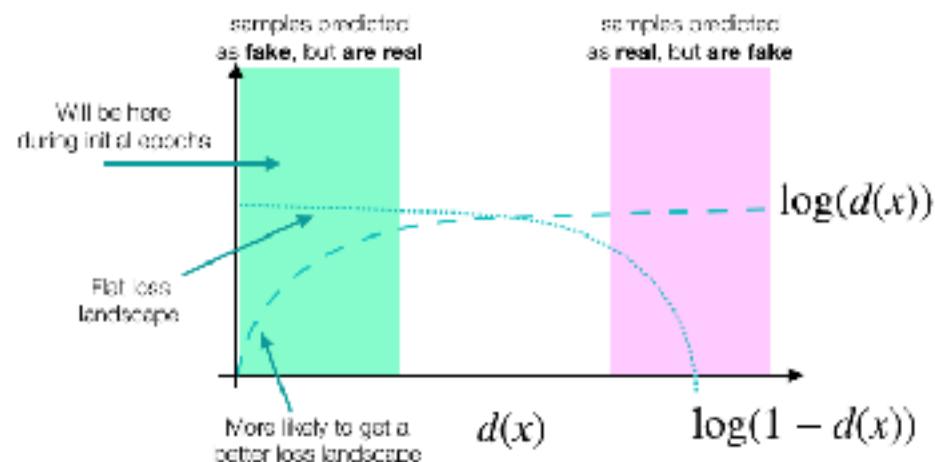
$$\max_d \mathbb{E}_{x \sim p_{data}} [\log d(x)] + \mathbb{E}_{x \sim g(z)} [\log(1 - d(x))]$$

Same as minimizing the binary cross entropy
of the model with: (1) labeled data and (2) generated data!

- Generator:

$$\max_g \mathbb{E}_{x \sim g(z)} [\log(d(x))] \quad \text{Freeze Discriminator Weights}$$

Same as minimizing the binary cross entropy
of "mislabeled" generated data!





GANs in Keras

Implementation from Book on
“Frogs from CIFAR”



Demo by Francois Chollet

Game:
Can you catch any errors!?

Follow Along: <https://github.com/fchollet/deep-learning-with-python-notebooks/blob/master/8.5-introduction-to-gans.ipynb>



GANs Loss Landscape

Abstract

Meta-meta-learning for Neural Architecture Search through arXiv Descent

Antreas Antoniou
MetaMind
aa@mm.ai

Nick Pawlowski
Google π^2
nick@x.z

Jack Turner
slow.ai
jack@slow.ai

James Owers
Facebook AI Research Team
jim@fart.org

Joseph Mellor
Institute of Yellow Jumpers
joe@anditwasall.yellow

Elliot J. Crowley
ClosedAI
elliot@closed.ai

Recent work in meta-learning has set the deep learning community alight. From minute gains on few-shot learning tasks, to discovering architectures that are slightly better than chance, to solving intelligence itself¹, meta-learning is proving a popular solution to every conceivable problem ever conceivable conceived ever.

In this paper we venture deeper into the computational insanity that is meta-learning, and potentially risk exiting the simulation of reality itself, by attempting to meta-learn at a third learning level. We showcase the resulting approach—which we call *meta-meta-learning*—for neural architecture search. Crucially, instead of *meta-learning* a neural architecture differentiably as in DARTS (Liu et al., 2018) we *meta-meta-learn* an architecture by searching through arXiv. This *arXiv descent* is GPU-free and only requires a handful of graduate students. Further, we introduce a regulariser, called *college-drapour*, which works by randomly removing a single graduate student from our system. As a consequence, procrastination levels decrease significantly, due to the increased workload and sense of responsibility each student attains.

The code for our experiments is publicly available at [REDACTED]
Edit: we have decided not to release our code as we are concerned that it may be used for malicious purposes.



The GAN Loss Landscape

- Problems:
 - GANs might not converge
 - GAN output might be the same sample (mode collapse)
 - Slow training: gradient vanishes
- A collection of Heuristics and Observations that is more Alchemy than Science
 - Feature Matching loss
 - Mini-batch discrimination
 - Historical averaging
 - Virtual Batch normalization
 - Experience Replay
 - Manipulating the Loss Function

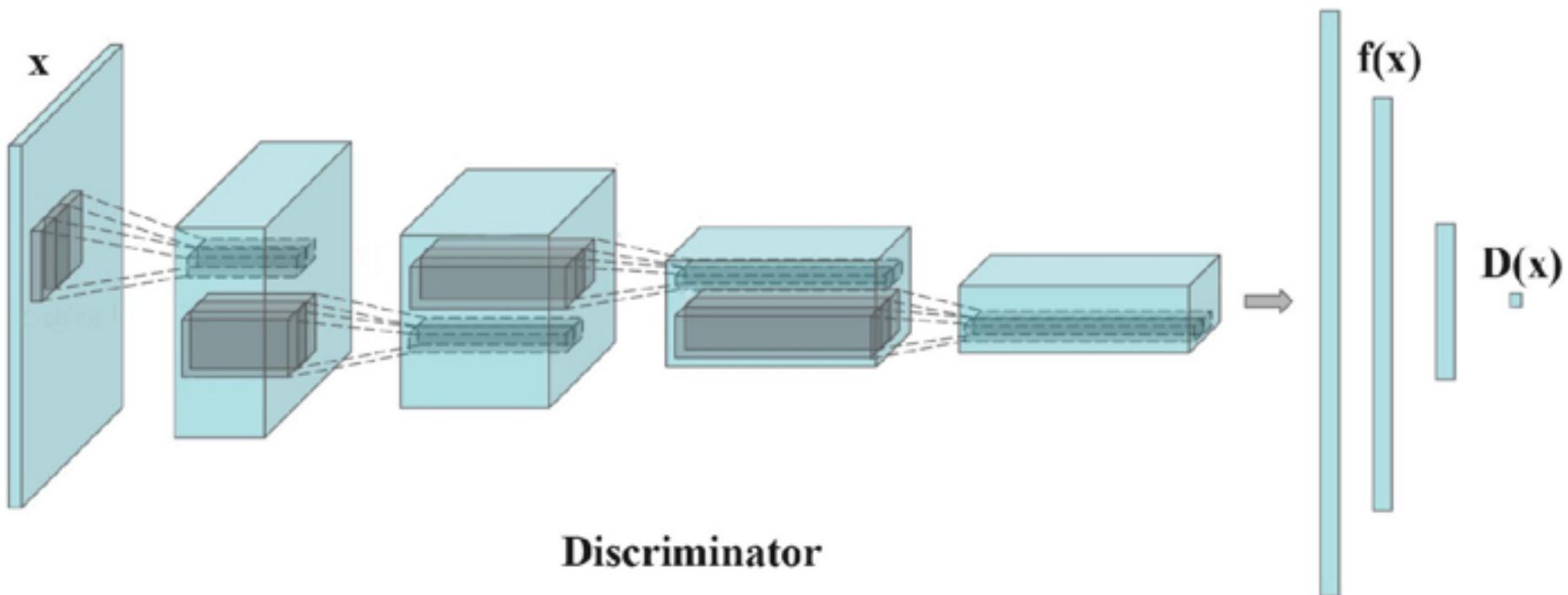


There are no GuAraNtees

- From Salimans, Goodfellow, et al., :
 - “*In this work, we introduce several techniques intended to encourage convergence of the GANs game. These techniques are motivated by a heuristic understanding of the non-convergence problem. They lead to improved semi-supervised learning performance and improved sample generation. We hope that some of them may form the basis for future work, providing formal guarantees of convergence.*”

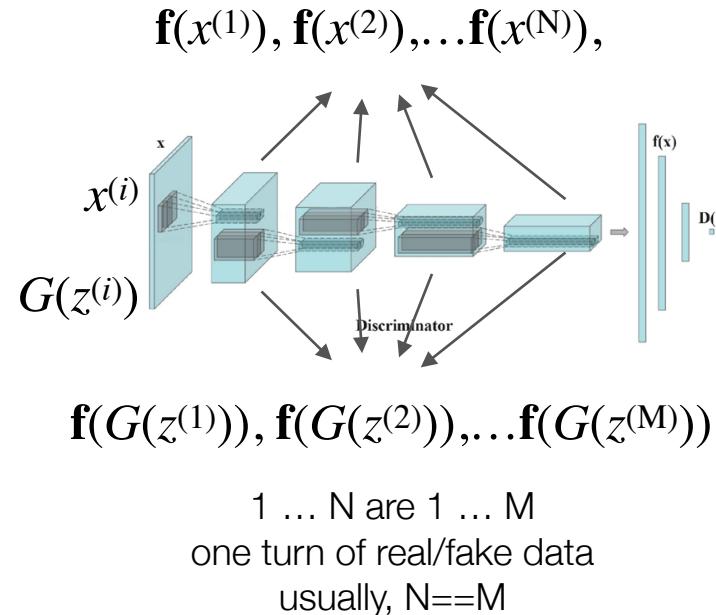


Before we go too far



Feature Matching

- Loss function for generator is really difficult and unstable
 - hard to optimize based on feedback only from discriminator output!
- Since generator is trying to statistically match features inside the discriminator (to fool it)
 - try to make generator match statistics of discriminator activations



$$\|\mathbf{E}_{x \leftarrow data}[\mathbf{f}(x)] - \mathbf{E}_{x \leftarrow q(z)}[\mathbf{f}(G(z))]\|_2^2$$

mean discriminator activations
from **real** data

mean discriminator activations
from **fake** data

```
real_mean = K.reduce_mean(real_features, axis = 0)
fake_mean = K.reduce_mean(fake_features, axis = 0)
feat_match = K.reduce_mean(tf.square(real_mean-fake_mean))
```

Advances in neural information

62

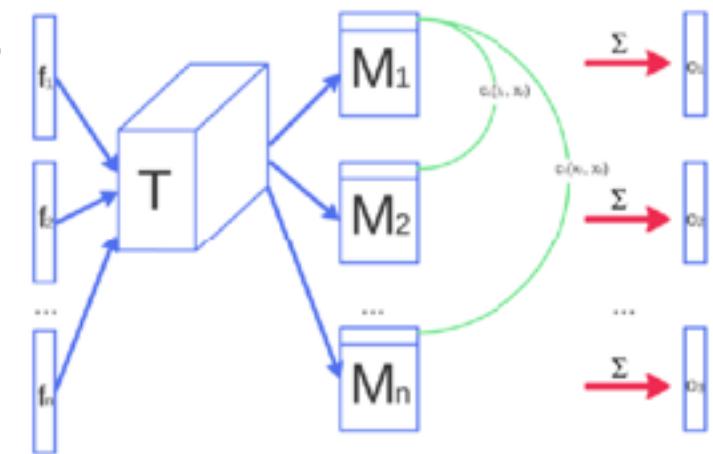


Mini-Batch Discrimination

- Mode collapse of generator happens when parameters emit the same point $G(z^{(i)}) \sim G(z^{(j)})$
- Solution: Incentivize dissimilarity of generated samples by letting discriminator see batches (detect mode collapse early, incentivizes generator to prevent it)
- Use discriminator activation $\mathbf{f}(G(z))$

$$o(x_i) = \exp \left(\sum_{j \in \text{Batch}} -\|\mathbf{f}(x_i) \cdot \mathbf{T} - \mathbf{f}(x_j) \cdot \mathbf{T}\| \right)$$

i^{th} sample
in batch j^{th} sample
in batch

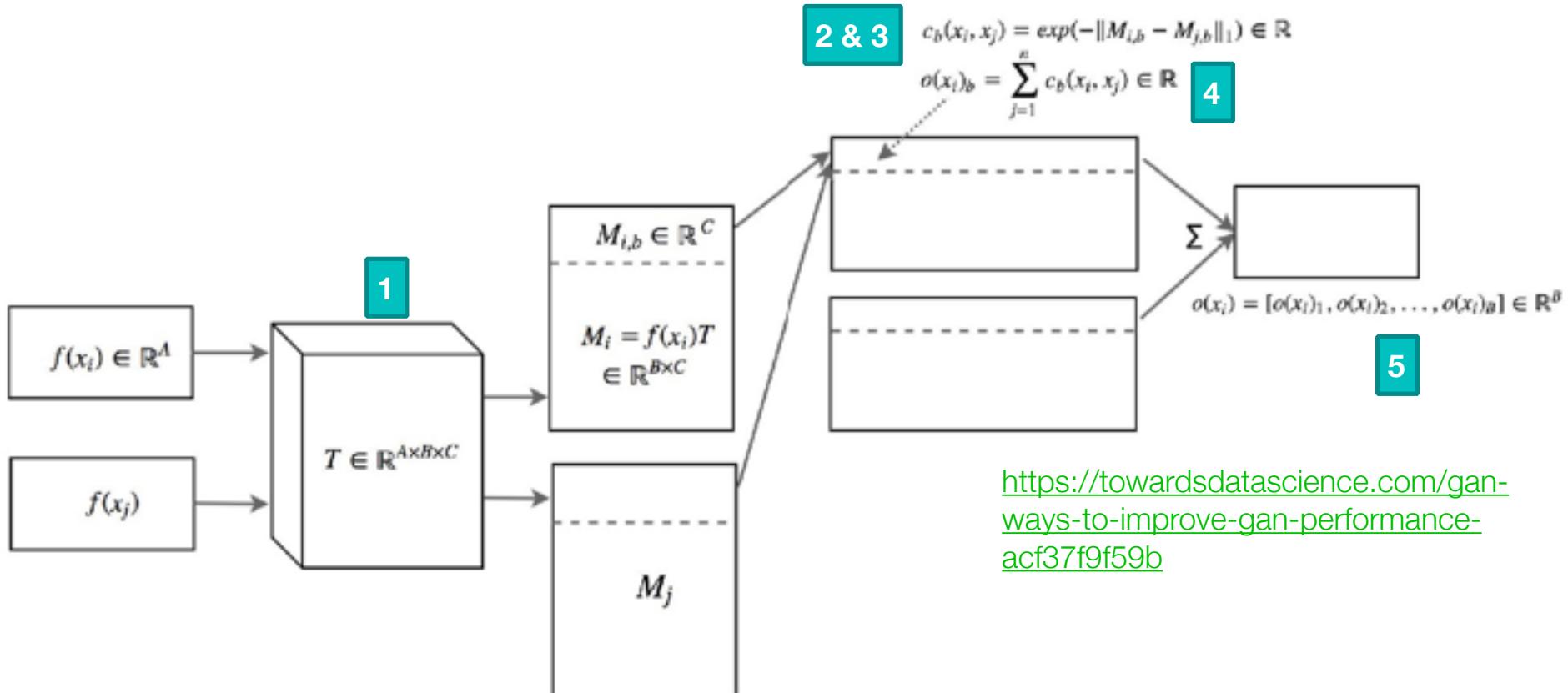


- Concatenate and Feed $o(x_{1:n})$ as a feature into next discriminator layer

Salimans, Tim, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. "Improved techniques for training gans." In *Advances in neural information processing systems*, pp. 2234-2242. 2016.



Mini-Batch Discrimination



```

M = K.reshape(K.dot(x, self.T), (-1, self.num_kernels, self.kernel_dim))
1
L1 = K.abs(K.expand_dims(M, 3) - K.expand_dims(K.permute_dimensions(M, [1, 2, 0]), 0))
2
c = K.exp(-K.sum(L1, axis=2))
3
o = K.sum(c, axis=2)
4
return K.concatenate([x, o], 1)
5

```

This broadcast is tricky, but powerful!

Salimans, Tim, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. "Improved techniques for training gans." In *Advances in neural information processing systems*, pp. 2234-2242. 2016.



Historical Averaging

- Because we cannot converge, parameters should constantly be exploring different solutions and exploiting current found solutions
- Solution: Add explore/exploit tradeoff directly in loss function

$$\lambda_t \left\| w[T] - \frac{1}{n} \sum_{i=T-n}^{T-1} w[t] \right\|^2$$



Virtual Batch Normalization

- Batch normalization is really swell
- But for GANs, it makes current x 's dependent on other generated x 's (and it was already hard to train...)
- Why not use a reference batch for statistics? And just use statistics of the reference batch activations for normalizing?
 - Every time we update W 's in the network, we need to recompute the statistics for the reference batch
 - That's it!

$$\mu^{ref} = \frac{1}{M} \sum_i a_i^{ref} \quad \sigma^{ref2} = \frac{1}{M} \sum_i (a_i^{ref} - \mu^{ref})^2 \quad z_i = \gamma \cdot \frac{(a_i - \mu^{ref})}{\sqrt{\sigma^{ref2} + \epsilon}} + \beta$$

Salimans, Tim, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. "Improved techniques for training gans." In *Advances in neural information processing systems*, pp. 2234-2242. 2016.



Experience Replay

- The discriminator can overpower the generator quickly and it over-learns to a particular epoch
- Solution: feed it more than just the current epoch
- Save previous epochs and randomly grab from them when updating the generator!
- Grab a coffee and chill...



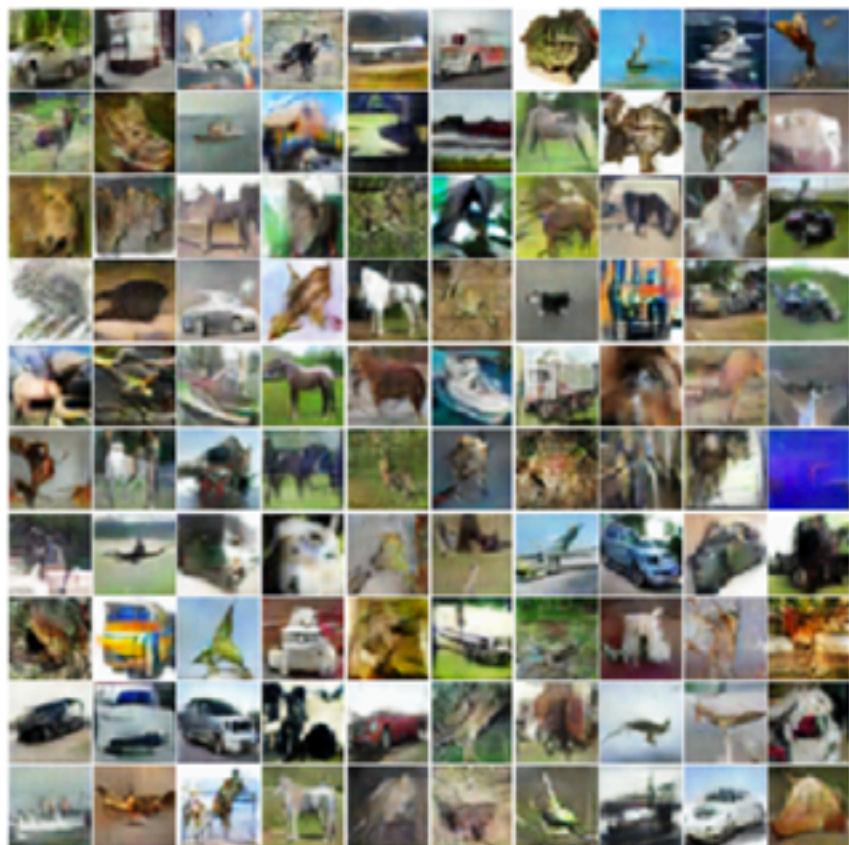
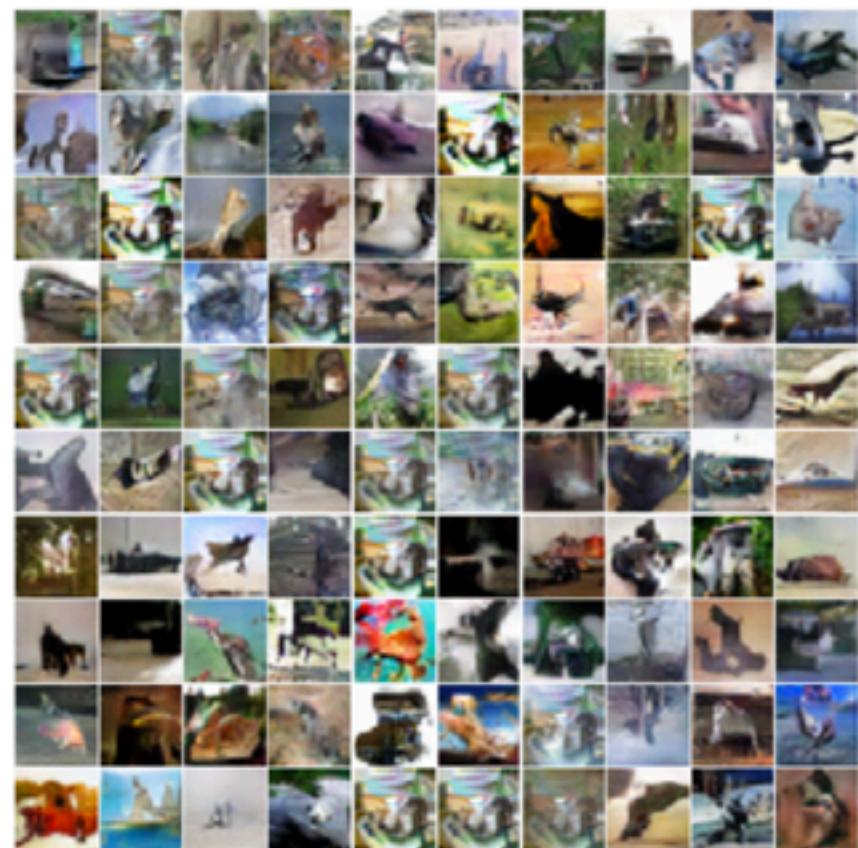


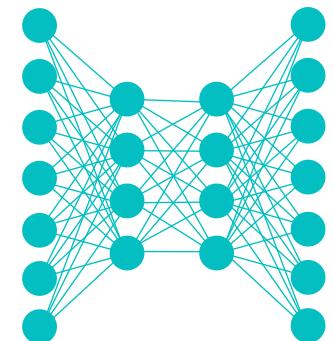
Figure 4: Samples generated during semi-supervised training on CIFAR-10 with feature matching (Section 3.1, *left*) and minibatch discrimination (Section 3.2, *right*).



Lecture Notes for
Neural Networks
and Machine Learning



Wasserstein
GANs



Logistics and Agenda

- Logistics
 - Student Presentations Today
- Agenda
 - **Student Presentation:** Representational Learning with Deep Convolutional GANs, Radford
 - **Student Presentation:** Towards Principled Methods for Training GANs, Arjovsky and Bottou 2017
 - Wasserstein GAN
 - Demo



Paper Presentation

UNSUPERVISED REPRESENTATION LEARNING WITH DEEP CONVOLUTIONAL GENERATIVE ADVERSARIAL NETWORKS

Alec Radford & Luke Metz
indico Research
Boston, MA
[\(alec,luke}@indico.io](mailto:(alec,luke}@indico.io)

Soumith Chintala
Facebook AI Research
New York, NY
soumith@fb.com

ABSTRACT

In recent years, supervised learning with convolutional networks (CNNs) has seen huge adoption in computer vision applications. Comparatively, unsupervised learning with CNNs has received less attention. In this work we hope to help bridge the gap between the success of CNNs for supervised learning and unsupervised learning. We introduce a class of CNNs called deep convolutional generative adversarial networks (DCGANs), that have certain architectural constraints, and demonstrate that they are a strong candidate for unsupervised learning. Training on various image datasets, we show convincing evidence that our deep convolutional adversarial pair learns a hierarchy of representations from object parts to scenes in both the generator and discriminator. Additionally, we use the learned features for novel tasks - demonstrating their applicability as general image representations.



Change the loss function entirely

Name	Value Function
GAN	$L_D^{GAN} = E[\log(D(x))] + E[\log(1 - D(G(x)))]$ $L_G^{GAN} = E[\log(D(G(x)))]$
LSGAN	$L_D^{LSGAN} = E[(D(x) - 1)^2] + E[D(G(x))^2]$ $L_G^{LSGAN} = E[(D(G(x)) - 1)^2]$
WGAN	$L_D^{WGAN} = E[D(x)] - E[D(G(z))]$ $L_G^{WGAN} = E[D(G(z))]$ $W_D = \text{clip_by_value}(W_D, 0, 1)$
WGAN_GP	$L_D^{WGAN_GP} = L_D^{WGAN} + \lambda E[VD(\alpha x + (1 - \alpha)G(x)) - 1]$ $L_G^{WGAN_GP} = L_G^{WGAN}$
DRAGAN	$L_D^{DRAGAN} = L_D^{GAN} + \lambda E[(VD(\alpha x - (1 - \alpha)x_p) - 1)^2]$ $L_G^{DRAGAN} = E[\log(D(G(x))) + \log(1 - D(G(x)))]$
CGAN	$L_G^{CGAN} = E[\log(D(G(x), c))]$
InfoGAN	$L_{D,Q}^{InfoGAN} = L_D^{GAN} - \lambda L_1(c, c')$ $L_G^{InfoGAN} = L_G^{GAN} - \lambda L_1(c, c')$
ACGAN	$L_{D,Q}^{ACGAN} = L_{D,Q}^{DRAGAN} + E[P(class = c D(x))] + E[P(class = c G(z))]$ $L_G^{ACGAN} = L_G^{GAN} + E[P(class = c G(z))]$
EBCGAN	$L_D^{EBCGAN} = D_{AE}(x) + \max(0, m - D_{AE}(G(x)))$ $L_G^{EBCGAN} = D_{AE}(G(x)) + \lambda \cdot PT$
BEGAN	$L_D^{BEGAN} = D_{AE}(x) - k_t D_{AE}(G(x))$ $L_G^{BEGAN} = D_{AE}(G(x))$ $k_{t+1} = k_t + \lambda(y D_{AE}(x) - D_{AE}(G(x)))$

The best loss function to use:

There has to be a more principled way of doing this!!!!

Academic blasphemy:
Maybe it doesn't really matter



Paper Presentation

TOWARDS PRINCIPLED METHODS FOR TRAINING GENERATIVE ADVERSARIAL NETWORKS

Martin Arjovsky

Courant Institute of Mathematical Sciences
martinarjovsky@gmail.com

Léon Bottou

Facebook AI Research
leonb@fb.com

ABSTRACT

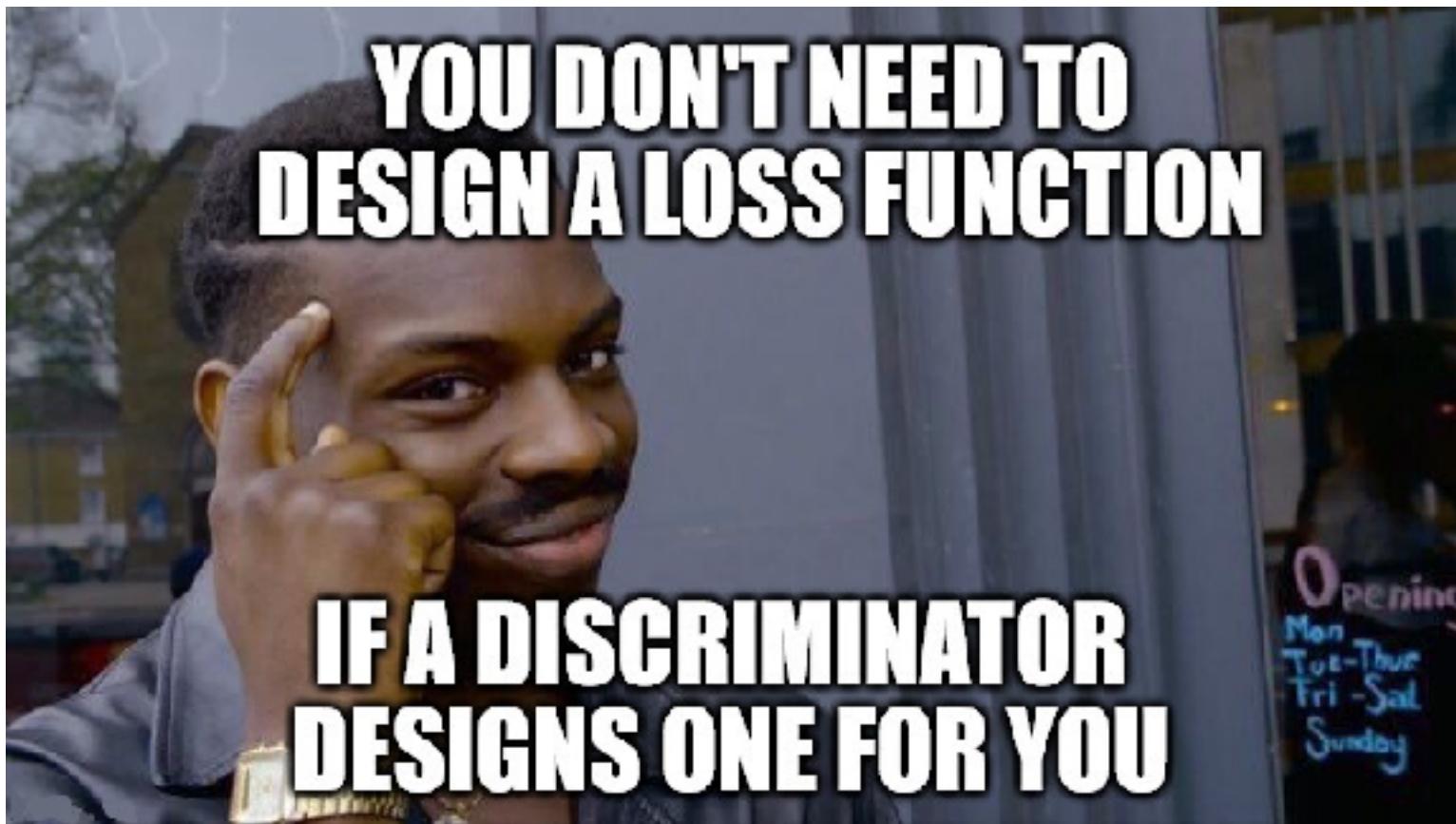
The goal of this paper is not to introduce a single algorithm or method, but to make theoretical steps towards fully understanding the training dynamics of generative adversarial networks. In order to substantiate our theoretical analysis, we perform targeted experiments to verify our assumptions, illustrate our claims, and quantify the phenomena. This paper is divided into three sections. The first section introduces the problem at hand. The second section is dedicated to studying and proving rigorously the problems including instability and saturation that arise when training generative adversarial networks. The third section examines a practical and theoretically grounded direction towards solving these problems, while introducing new tools to study them.



This is still a good paper, its just
not going to give you what you want:
answers



Wasserstein GANs

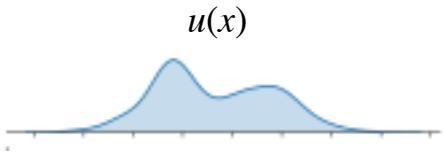


Caveat

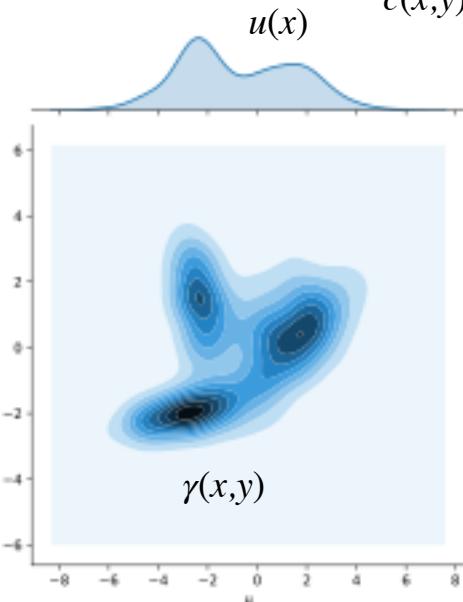
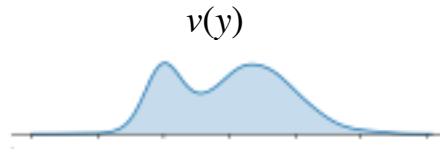
- Wasserstein Loss is actually very useful
 - But the mathematics pinning it down are based on approximations (like always)
 - So its **not really true** to say its actually working for the motivated reasons
 - But **it might be**
- So until we have a better theory, **lets assume it works** because it truly is a great distance measure



Wasserstein Distance (Earth Mover)



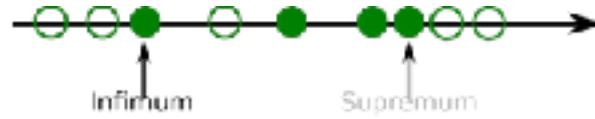
how to move dirt
from u to v ?



$\gamma(x,y)$ one plan to move u to v
 $c(x,y)$ cost of moving u to v

$$\iint c(x, y) \cdot \gamma(x, y) \, dx \, dy \quad \text{Total cost of plan } \gamma$$

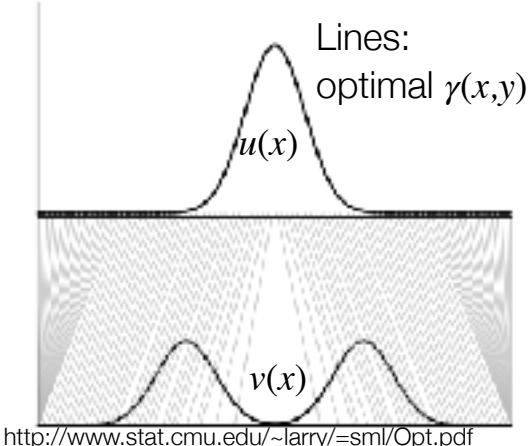
$$\inf_{\gamma \in \Pi} \iint c(x, y) \cdot \gamma(x, y) \, dx \, dy$$



Optimal plan
has greatest lower bound
(minimum cost in set Π)

$$\inf \mathbf{E}_{x,y \sim \gamma} [|x - y|] = \inf_{\gamma \in \Pi} \iint |x - y| \cdot \gamma(x, y) \, dx \, dy$$

If cost is difference
to move from x to y

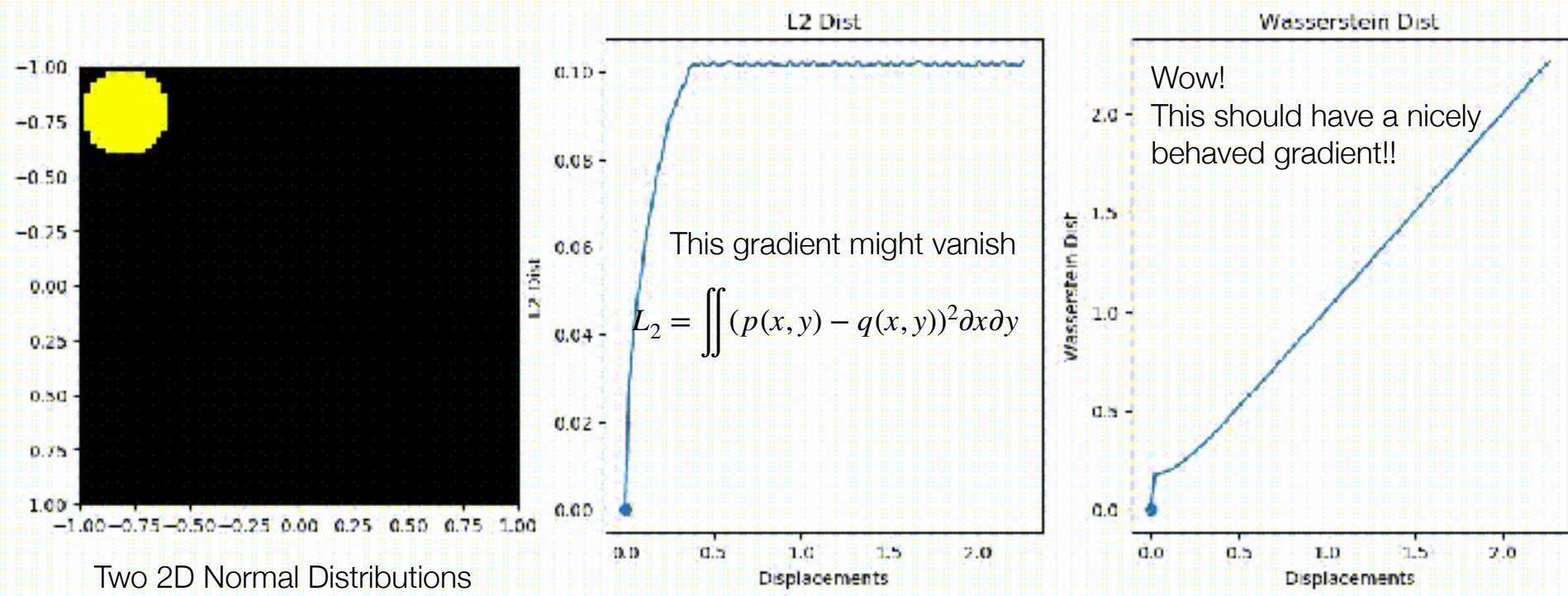


<http://www.stat.cmu.edu/~larry/sml/Opt.pdf>

image: wikipedia



Wasserstein Distance



<https://gist.github.com/ctralie/66352ae6ab06c009f02c705385a446f3>



Wasserstein

- *Wasserstein GAN adds few tricks to allow the Discriminator to approximate Wasserstein (aka Earth Mover's) distance between real and model distributions. Wasserstein distance roughly tells “how much work is needed to be done for one distribution to be adjusted to match another” and is remarkable in a way that it is defined even for non-overlapping distributions.*
- That should help training even when the generator and discriminator are not matching
- ...but Wasserstein Distance is not tractable to calculate for generative distributions so its time to start approximating!



Wasserstein Distance

- How much do I need to change one distribution to get it to match another? Moving dirt...
- Could use KL divergence, but we already know that has many problems associated with vanishing gradients
- Also, we will only have **samples** from the distributions (not the actual equations)
- **Wasserstein Distance** for continuous probabilities:

$$W(p_{data}, p_g) = \inf_{\gamma \in \prod(p_{data}, p_g)} \mathbf{E}_{x,y \leftarrow \gamma} [\|x - y\|]$$

- inf is greatest lower bound (min of a set)
- γ is completely and utterly **intractable** to **compute**



Wasserstein Duality

$$\begin{aligned} W(p_{data}, p_g) &= \inf_{\gamma \in \prod(p_{data}, p_g)} \mathbf{E}_{x,y \leftarrow \gamma} [\|x - y\|] \\ &= \inf_{\gamma} \mathbf{E}_{x,y \leftarrow \gamma} \left[\|x - y\| + \underbrace{\sup_f \mathbf{E}_{x \leftarrow p_{data}} [f(x)] - \mathbf{E}_{x \leftarrow p_g} [f(x)] - (f(x) - f(y))}_{\text{enforce constraint } =0, \text{ if } \gamma \in \Pi, \text{ else } +\infty} \right] \\ &= \inf_{\gamma} \sup_f \mathbf{E}_{x,y \leftarrow \gamma} \left[\|x - y\| + \mathbf{E}_{x \leftarrow p_{data}} [f(x)] - \mathbf{E}_{x \leftarrow p_g} [f(x)] - (f(x) - f(y)) \right] \\ &= \sup_f \mathbf{E}_{x \leftarrow p_{data}} [f(x)] - \mathbf{E}_{x \leftarrow p_g} [f(x)] - \inf_{\gamma} \underbrace{\mathbf{E}_{x,y \leftarrow \gamma} [\|x - y\| + (f(x) - f(y))]}_{\text{new constraint } =0, \text{ if } |f| \leq 1, \text{ else } -\infty} \\ &= \sup_{|f| \leq 1} \mathbf{E}_{x \leftarrow p_{data}} [f(x)] - \mathbf{E}_{x \leftarrow p_g} [f(x)] \end{aligned}$$

read this: <https://vincentherrmann.github.io/blog/wasserstein/>

80



Wasserstein Duality

- 1-Lipschitz constraint, critic: $|f(x_1) - f(x_2)| \leq |x_1 - x_2|$
- Wasserstein Duality Formula (approx for samples):

$$\mathcal{L}_{wass} = \sup_{|f| \leq 1} \mathbf{E}[f(x_{real})] - \mathbf{E}[f(x_{fake})] = \frac{1}{m} \sum_{i=1}^m f(x_{real}^{(i)}) - \frac{1}{n} \sum_{j=1}^n f(g(z^{(j)}))$$

- where \sup is least upper bound (maximization within set)
- f will be the critic, learned as a neural network, $m=n$

$$\frac{1}{m} \sum_{i=1}^m f(x_{real}^{(i)}) - f(g(z^{(i)}))$$

Maximize f (called critic),
freeze generator weights

$$\frac{1}{m} \sum_{i=1}^m f(g(z^{(i)}))$$

Maximize g ,
freeze weights of critic

Sometimes this for critic

$$\text{minimize: } \frac{1}{m} \sum_{i=1}^m f(x_{real}^{(i)}) \cdot f(g(z^{(i)}))$$

when *label* is -1 and 1

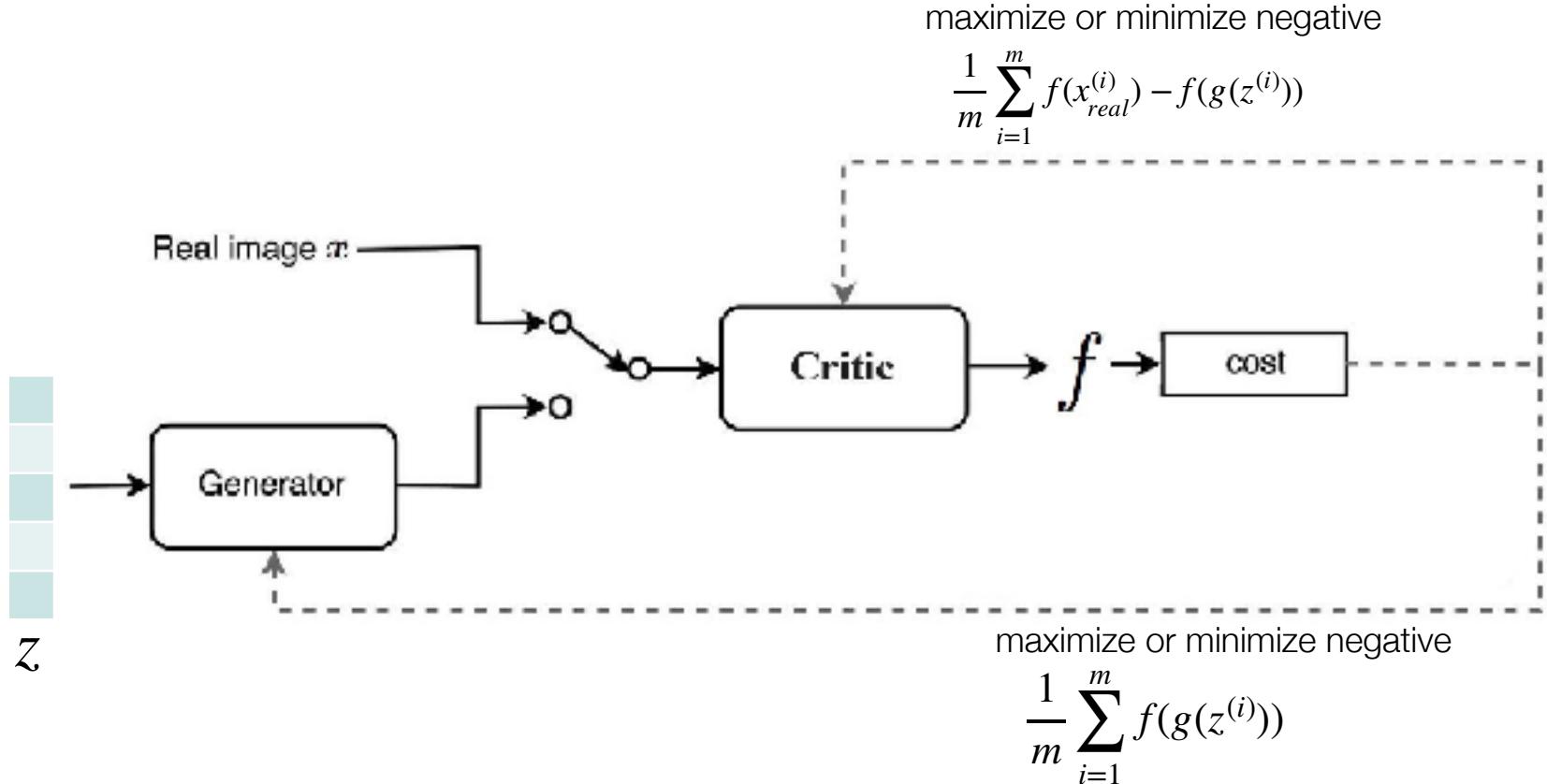
real	fake	$f(x)-g(x)$	$f(x)*g(x)$
-1	-1	0	1
-1	1	-2	-1
1	-1	2	-1
1	1	0	1

max min

https://medium.com/@jonathan_hui/gan-wasserstein-gan-wgan-gp-6a1a2aa1b490



Wasserstein Architecture



Practical Wasserstein Loss

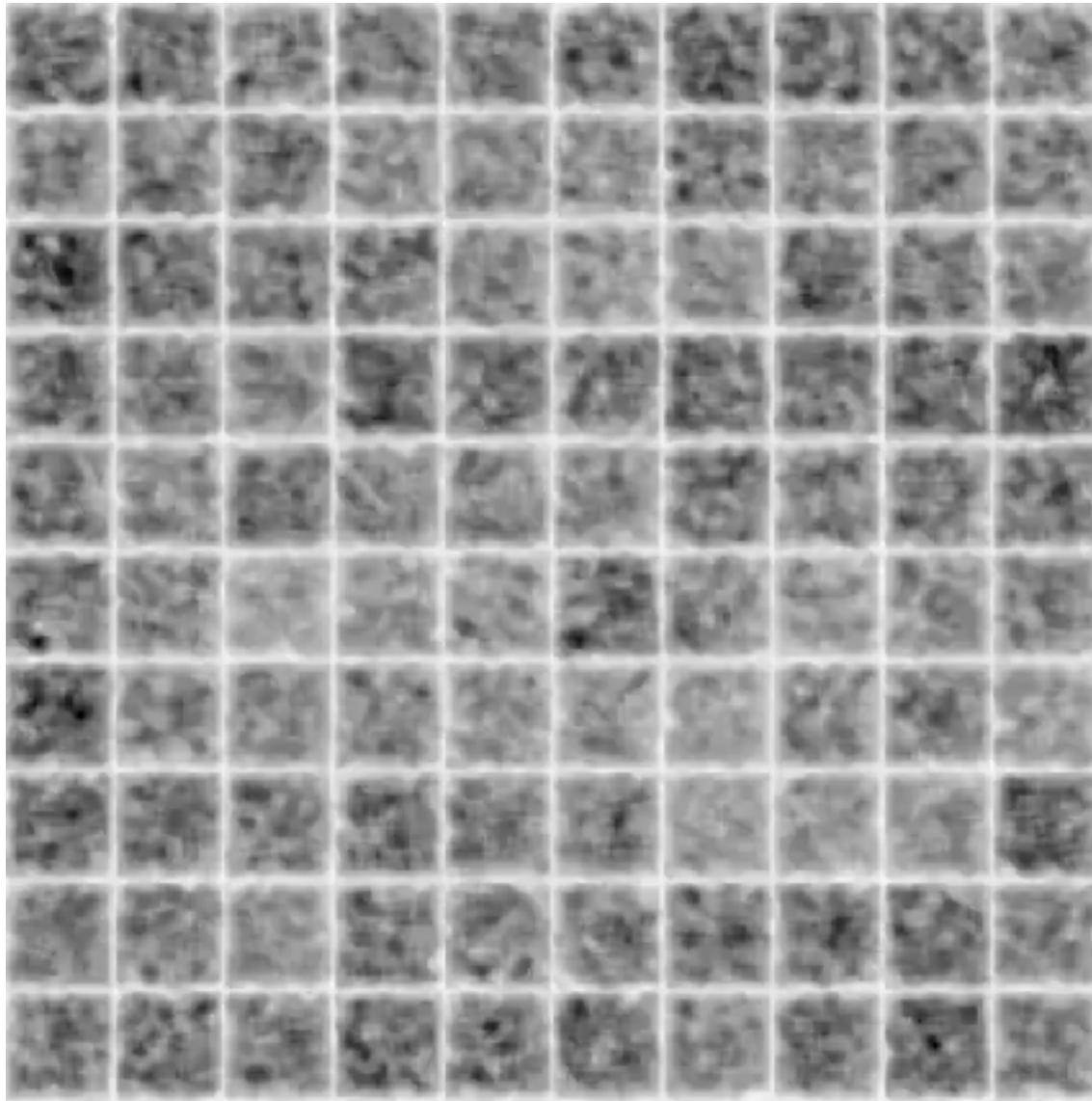
- Discriminator (now critic) becomes \mathbf{f} , and its output can be Lipschitz if all weights are small (squashes inputs)
 - so we clip them to $(-c \text{ to } c)$ (where $c \sim 0.01$)
 - critic output should be linear

```
# define the standalone critic model
def define_critic(in_shape=(28,28,1)):
    # weight initialization
    init = RandomNormal(stddev=0.02)
    # weight constraint
    const = ClipConstraint(0.01)
    # define model
    model = Sequential()
    # downsample to 14x14
    model.add(Conv2D(64, (4,4), strides=(2,2), padding='same'))
    model.add(BatchNormalization())
    model.add(LeakyReLU(alpha=0.2))
    # downsample to 7x7
    model.add(Conv2D(64, (4,4), strides=(2,2), padding='same'))
    model.add(BatchNormalization())
    model.add(LeakyReLU(alpha=0.2))
    # scoring, linear activation
    model.add(Flatten())
    model.add(Dense(1))
```

<https://machinelearningmastery.com/how-to-code-a-wasserstein-generative-adversarial-network-wgan-from-scratch/>



WGAN Example from Keras (MNIST)



This is a conditional GAN, which we have not yet discussed.

Intuitively it adds some control over the class being generated.



- In their empirical findings, no mode collapse problems
- And training longer resulted in good quality images (motivates that distributions overlap)
- Still some of the most competitive GAN results



WGAN with DCGAN generator



GAN with DCGAN generator



Without batch normalization & constant number of filters at each layer



Using a MLP as the generator

So we should always use Wasserstein!

- **Actually not really**
- Its really, really slow
- Clipping:

Weight clipping is a clearly terrible way to enforce a Lipschitz constraint. If the clipping parameter is large, then it can take a long time for any weights to reach their limit, thereby making it harder to train the critic till optimality. If the clipping is small, this can easily lead to vanishing gradients when the number of layers is big, or batch normalization is not used (such as in RNNs).

- Others have made improvements by incorporating gradient constraints
- New critic: WGAN with gradient penalty
- But... its still really slow compared to others



WGAN-GP (Gradient Penalty)

$$\frac{1}{m} \sum_{i=1}^m f(g(z^{(i)}))$$

Maximize g ,
freeze weights of critic
No change

- Theorem: a function is 1-Lipschitz if and only if its gradient has a norm less than or equal to one.

Maximize f
freeze generator weights,
incentivize critic gradient norm to be one

Choose large lambda
like lambda = 10

$$\frac{1}{m} \sum_{i=1}^m f(x_{real}^{(i)}) - f(g(z^{(i)})) + \lambda \frac{1}{W} \sum_{i=1}^W (\|\nabla f(\hat{x}^{(i)})\|_2 - 1)^2$$

where for ϵ in $U[0,1]$ $\hat{x}^{(i)} = \epsilon \cdot x_{real}^{(i)} + (1 - \epsilon) \cdot g(z^{(i)})$

randomly mix together real and fake images



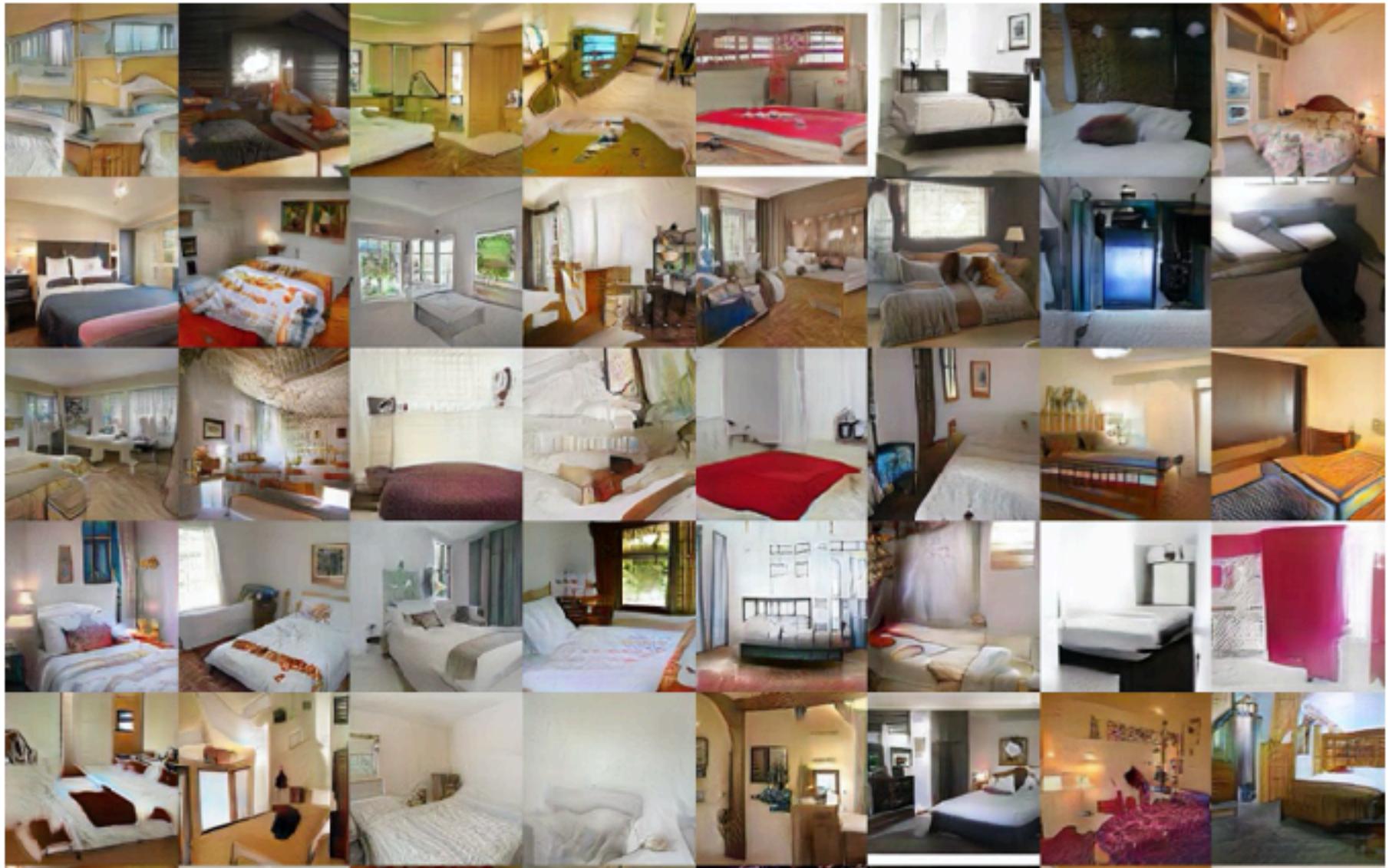
WGAN-GP Heuristics

- Choose large penalty coefficient
- Don't use batch normalization (in critic)
- Enforce gradient norm to be close to one, rather than less than one

Two-sided penalty We encourage the norm of the gradient to go towards 1 (two-sided penalty) instead of just staying below 1 (one-sided penalty). Empirically this seems not to constrain the critic too much, likely because the optimal WGAN critic anyway has gradients with norm 1 almost everywhere under \mathbb{P}_r and \mathbb{P}_g and in large portions of the region in between (see subsection 2.3). In our early observations we found this to perform slightly better, but we don't investigate this fully. We describe experiments on the one-sided penalty in the appendix.



WGAN-GP Results



89



WGAN-GP Comparative Results

WGAN (clipping)

Baseline: G: DCGAN Crit: DCGAN



WGAN-GP (ours)

G: DCGAN no BN Crit: DCGAN



G: MLP Crit: DCGAN



WGAN-GP Comparative Results

WGAN (clipping)

G/Crit: Tanh Non-linearities



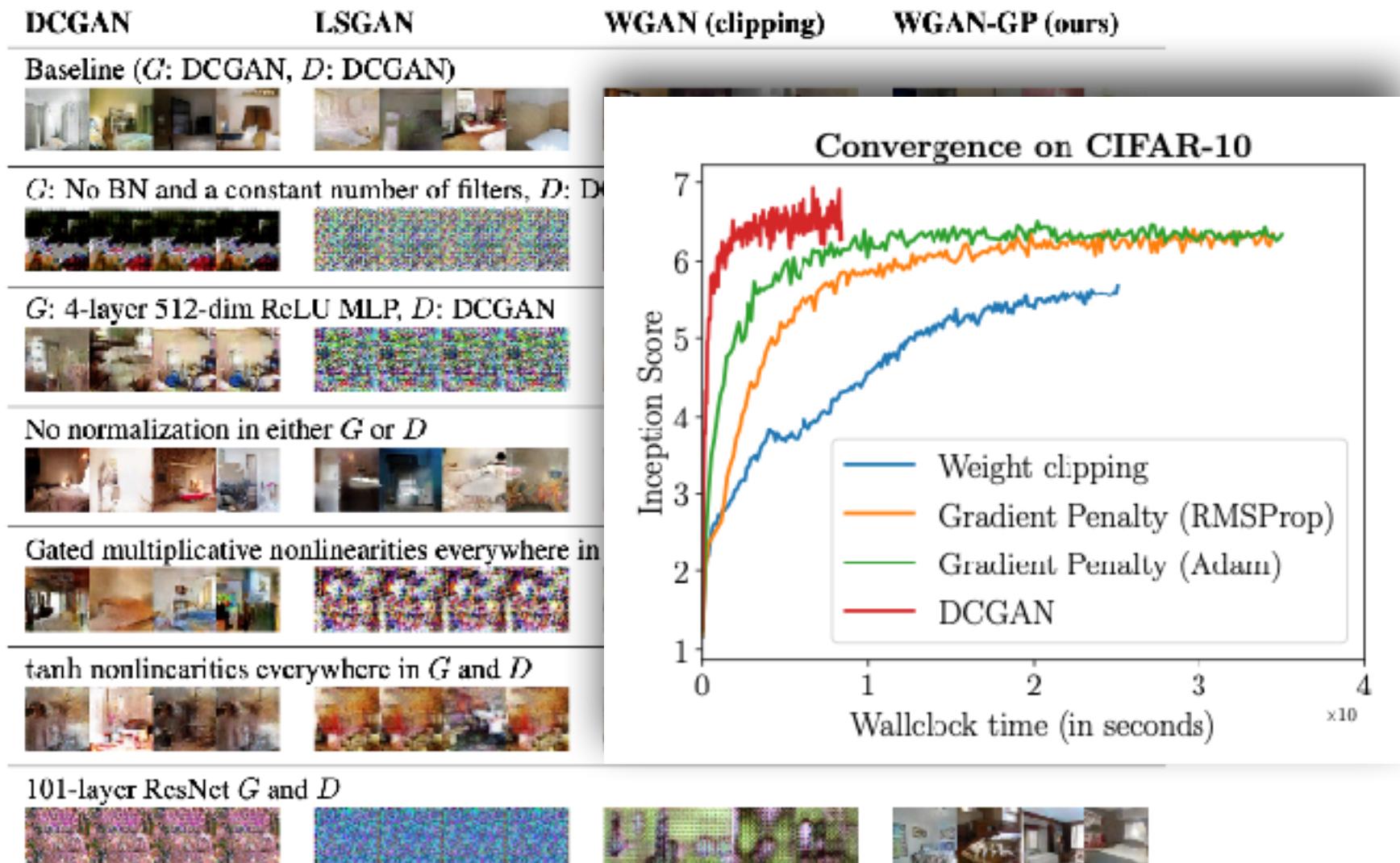
WGAN-GP (ours)



G/Crit: 101 Layer ResNet



WGAN-GP Comparative Results





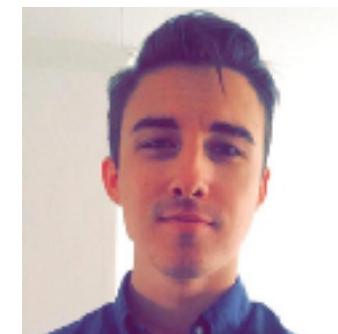
WGAN-GPs in Keras/PyTorch

Implementation from Github on
“MNIST”

Keras Follow Along: https://github.com/keras-team/keras-contrib/blob/master/examples/improved_wgan.py DCGAN

Demo by
Keras Contrib Community

PyTorch Follow Along: https://github.com/eriklindernoren/PyTorch-GAN/blob/master/implementation/wgan_gp/wgan_gp.py MLP-GAN



Demo by Erik Linder-Norén

93

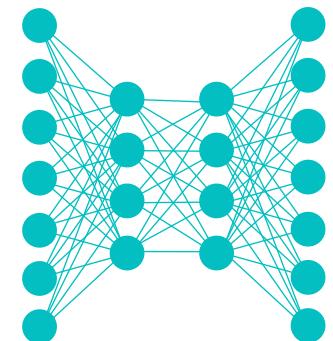


Lecture Notes for **Neural Networks** **and Machine Learning**

Wasserstein GANs



Next Time:
Beyond Generation
Reading: Chollet CH8





Backup slides



Title Between Topics



Example Slide





Title

Subtitle

Follow Along: Notebook Name

99

