

Lecture Notes for
Neural Networks
and Machine Learning



Generative Adversarial
Networks Overview



Logistics and Agenda

- Logistics
 - None
- Agenda
 - VAE Demo (last time)
 - AAE
 - Simple Generative Adversarial Networks
 - GANs Demo



Last Time: VAEs

```
# encode the input into a mean and variance parameter
x_mean, x_log_variance = encoder(input_img)
mu[x(i)] = E[μ(i)]
sigma[x(i)] = sqrt(E[σ(i)])
# Draw a latent point using a small random epsilon
z = x_mean + exp(x_log_variance) * epsilon
z = μ(x(i)) + exp(Σ(x(i))) · N(0,1)

# then decode z back to an image
reconstructed_img = decoder(z)
reconstructed_img = p(x(i)|z)

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

def vae_loss(psrf, x, z_decoded):
    x = K.flatten(x)
    z_decoded = K.flatten(z_decoded)
    xent_loss = keras.metrics.binary_crossentropy(x, z_decoded) - Eq(z|x(i))[log p(x(i)|z)]
    kl_loss = -1e-4 * K.mean(
        1 + x_log_var - K.square(x_mean) - K.exp(-1.0*x_log_var), axis=-1)
    return xent_loss + kl_loss

Note:
Hipped from maximization to minimization
```

$$-\sum_z 1 + \widehat{\Sigma(x^{(i)})} - \mu(x^{(i)})^2 - \exp(\widehat{\Sigma(x^{(i)})})$$

$$= -E_{q(z|x^{(i)})} [\log p(x^{(i)}|z)] - \sum_z 1 + \widehat{\Sigma(x^{(i)})} - \mu(x^{(i)})^2 - \exp(\widehat{\Sigma(x^{(i)})})$$



VAEs in Keras

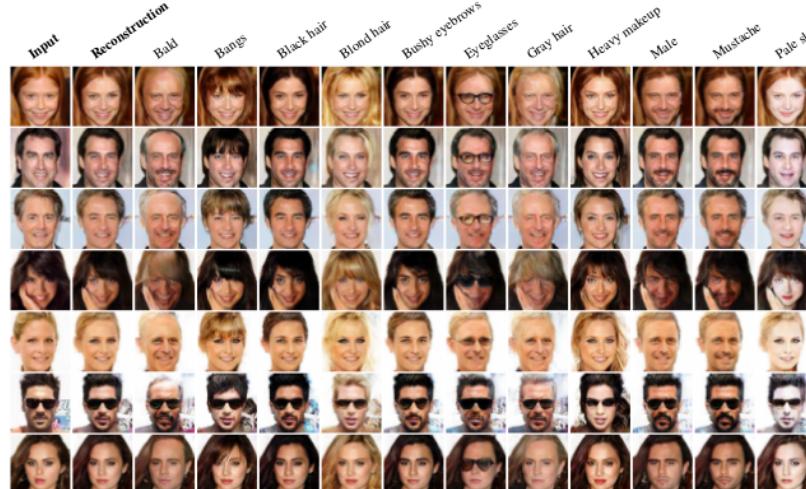
Sampling from variational auto encoder using MNIST



Demo by Francois Chollet

In Master Repo: [07a_VAEs_in_Keras.ipynb](#)

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



24



Adversarial Auto Encoding



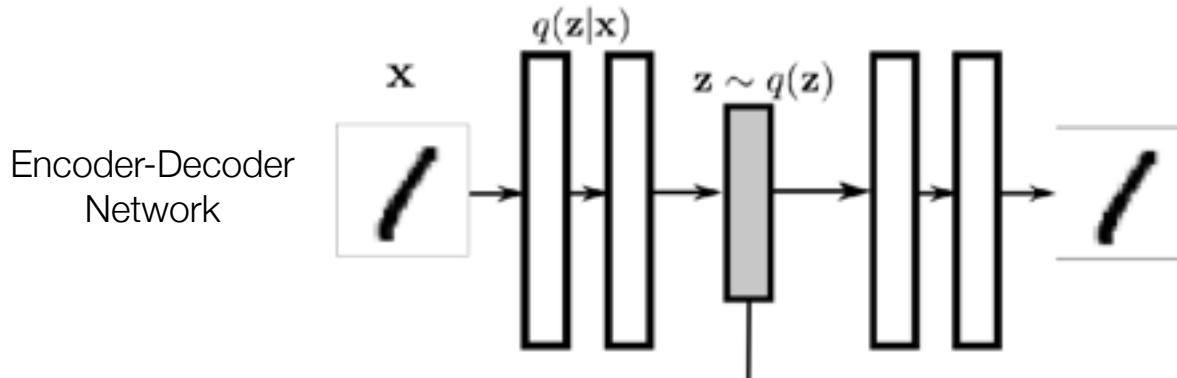
Do we need something more than VAE?

- Arguments for Yes:
 - ELBO is not global optimum! But... provides theory
 - Assumption of Normal distributions to $q(z)$ is limiting
 - Training tends to be slower (...so do GANs...)
 - Manifold of distributions do not cover the latent space completely (not guaranteed)
 - We can't incorporate distributions separately for different classes without reformulating loss function
- Arguments for No:
 - It seems hard, how can we research methods that aren't low hanging fruit? Plus the VAE math was like really hard for me to understand so this is not going to be very fun, guaranteed. Ah, fine lets look at it.



The Main Idea

- How can we enforce constraints on the latent space with a pair of networks?



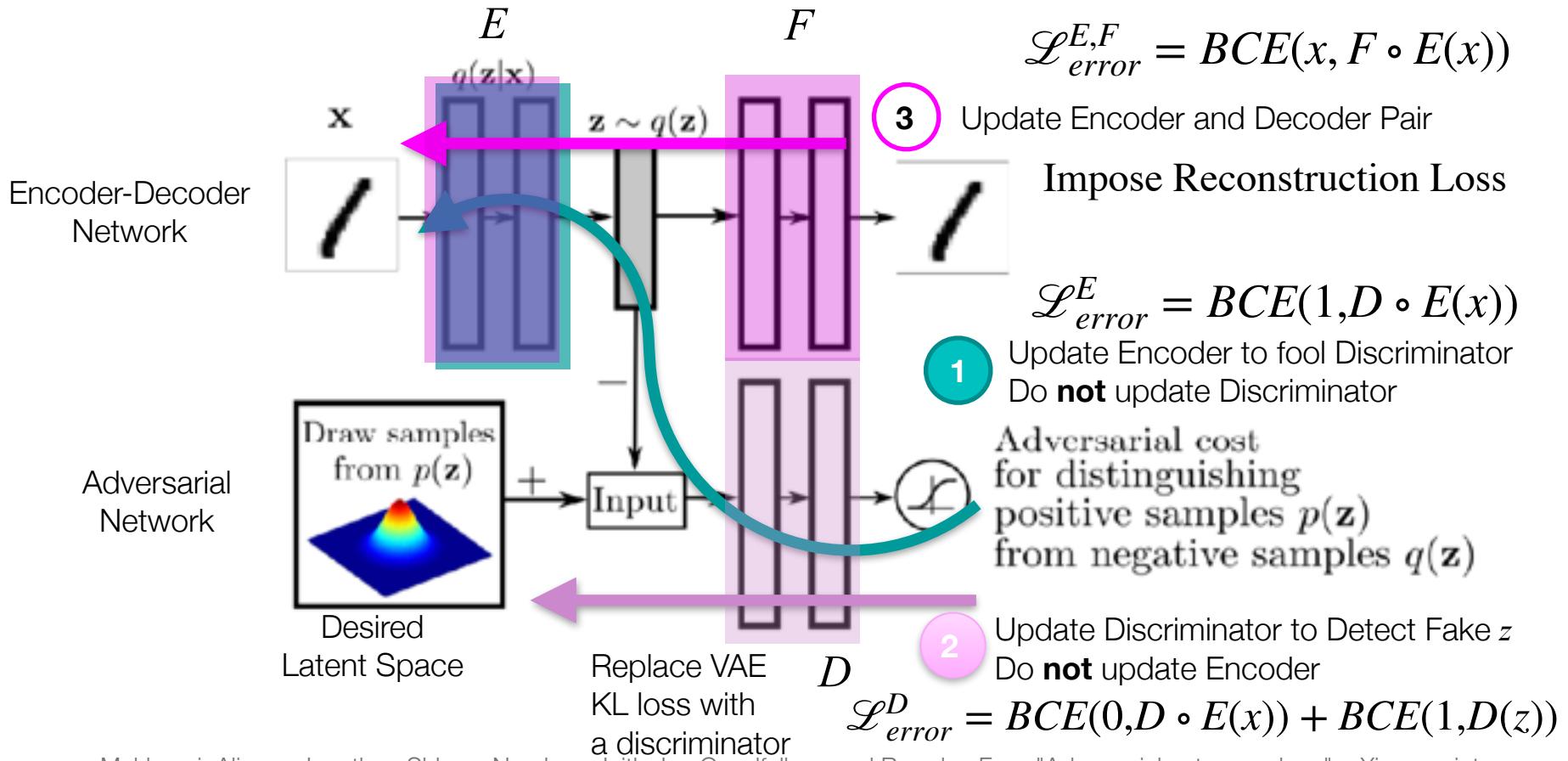
Makhzani, Alireza, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. "Adversarial autoencoders." arXiv preprint arXiv:1511.05644 (2015).

33



The Main Idea

- How can we enforce constraints on the latent space with a pair of networks?

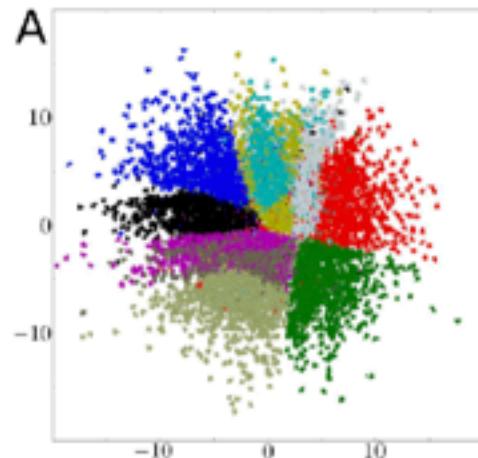


Makhzani, Alireza, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. "Adversarial autoencoders." arXiv preprint arXiv:1511.05644 (2015).

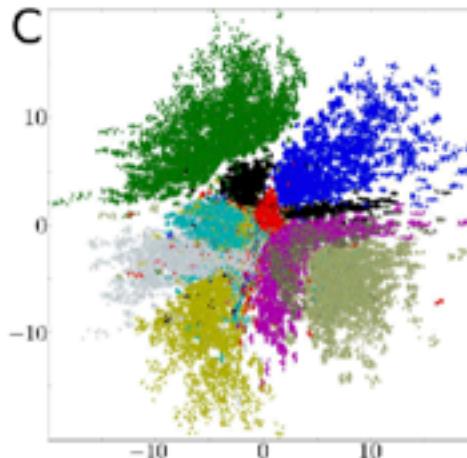


Arbitrary Prior Distributions

Adversarial Autoencoder



Variational Autoencoder

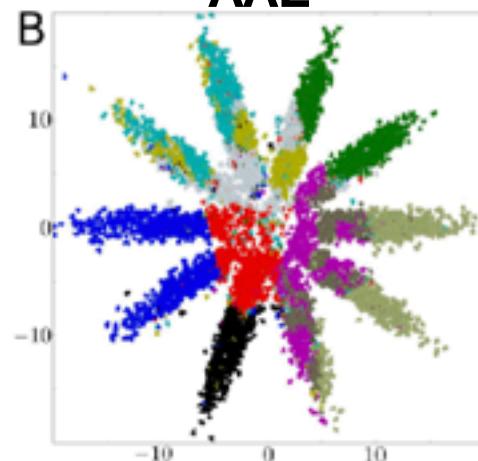


Manifold of
Adversarial Autoencoder

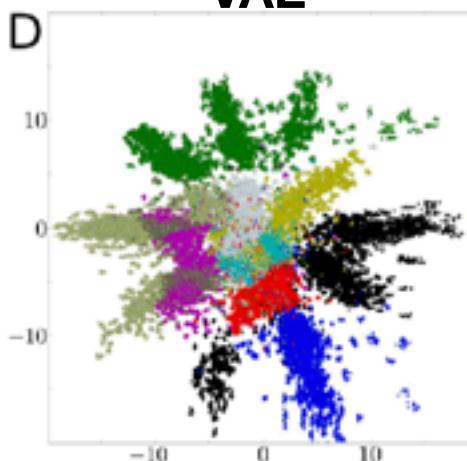
E

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

AAE



VAE

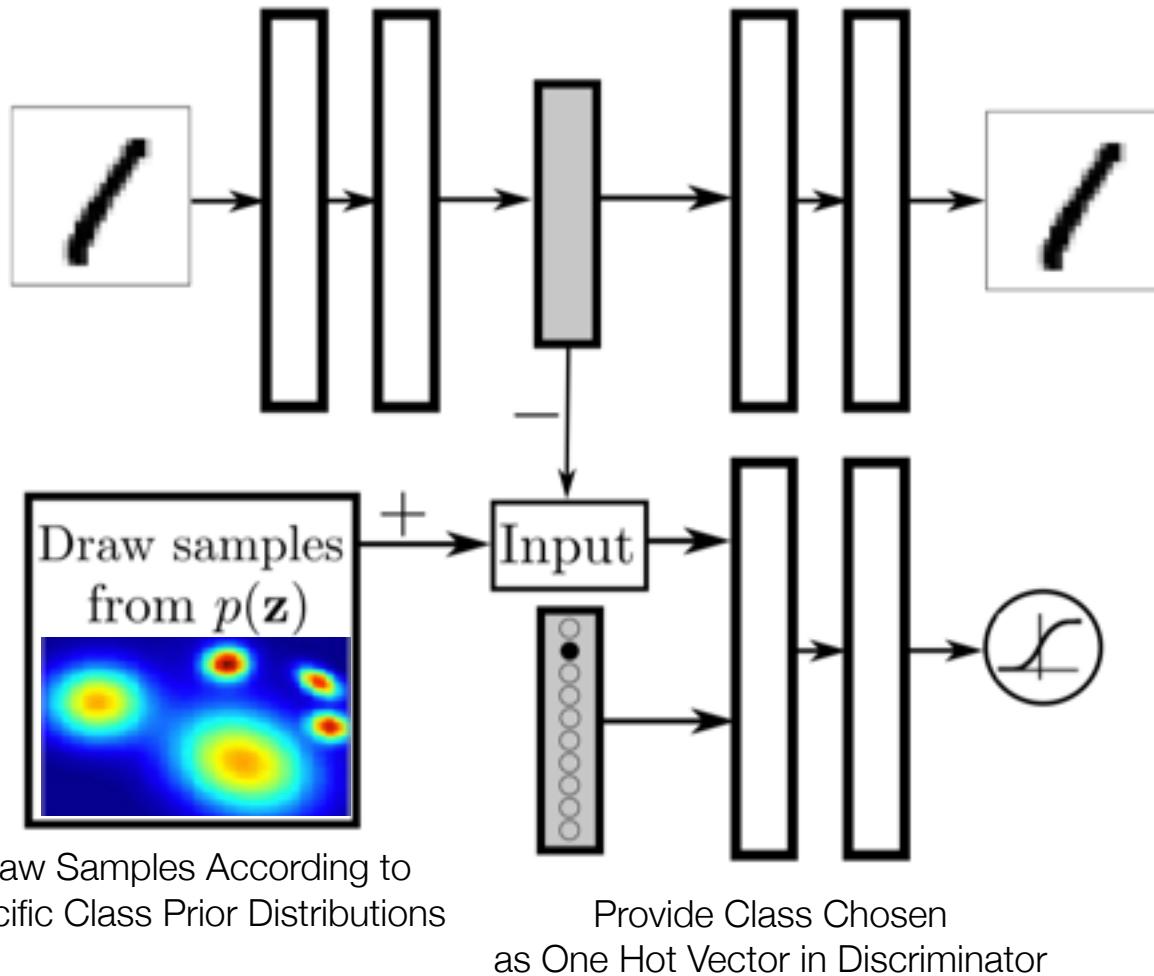


0	5
1	6
2	7
3	8
4	9

Makhzani, Alireza, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. "Adversarial autoencoders." arXiv preprint arXiv:1511.05644 (2015).



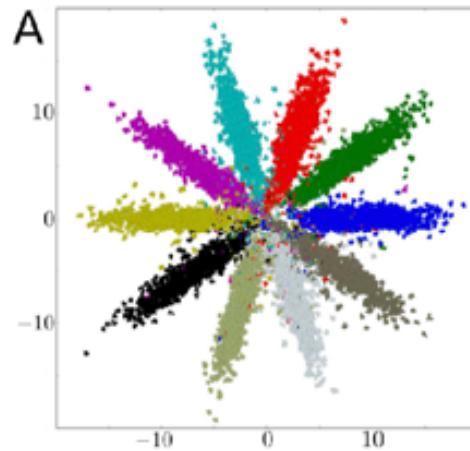
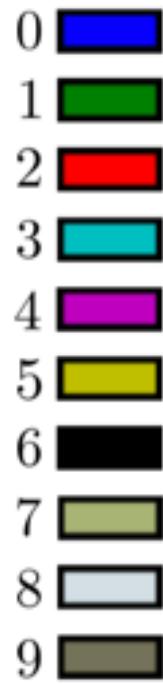
Sampling From Classes



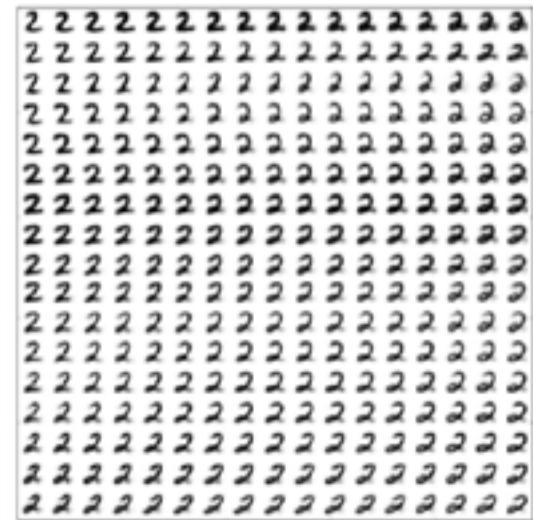
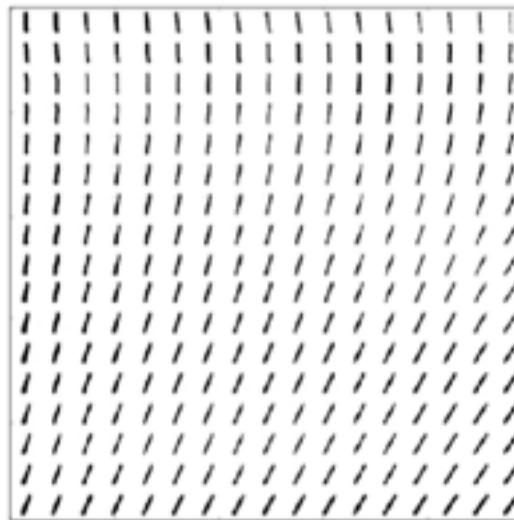
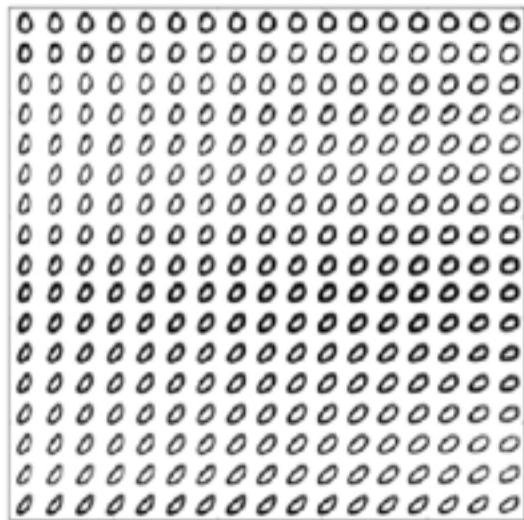
Makhzani, Alireza, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. "Adversarial autoencoders." arXiv preprint arXiv:1511.05644 (2015).



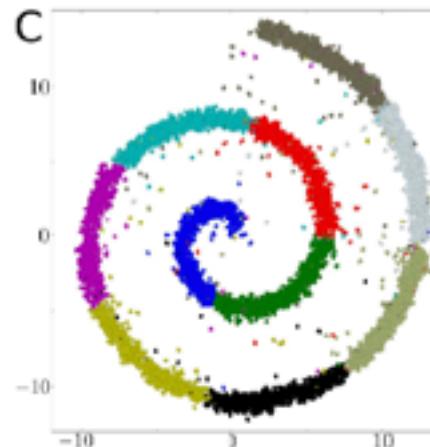
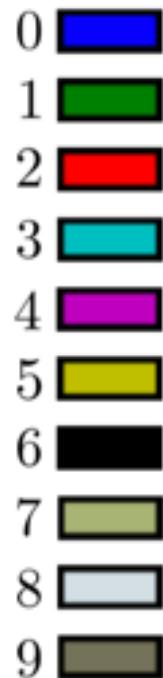
Conditional Class Latent Spaces



Sample Along Main Axis of the Gaussian Component for Each Digit



Conditional Class Latent Spaces



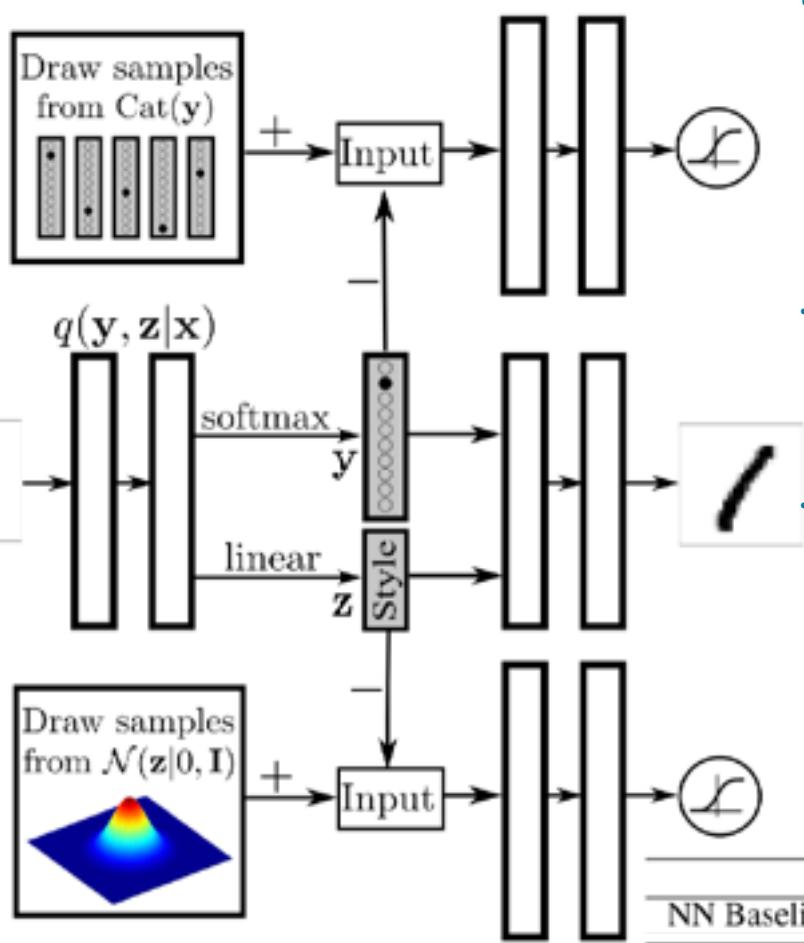
Sample Along Swiss Roll Axis



Makhzani, Alireza, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. "Adversarial autoencoders." *arXiv preprint arXiv:1511.05644* (2015): 38



Semi-Supervised Classification



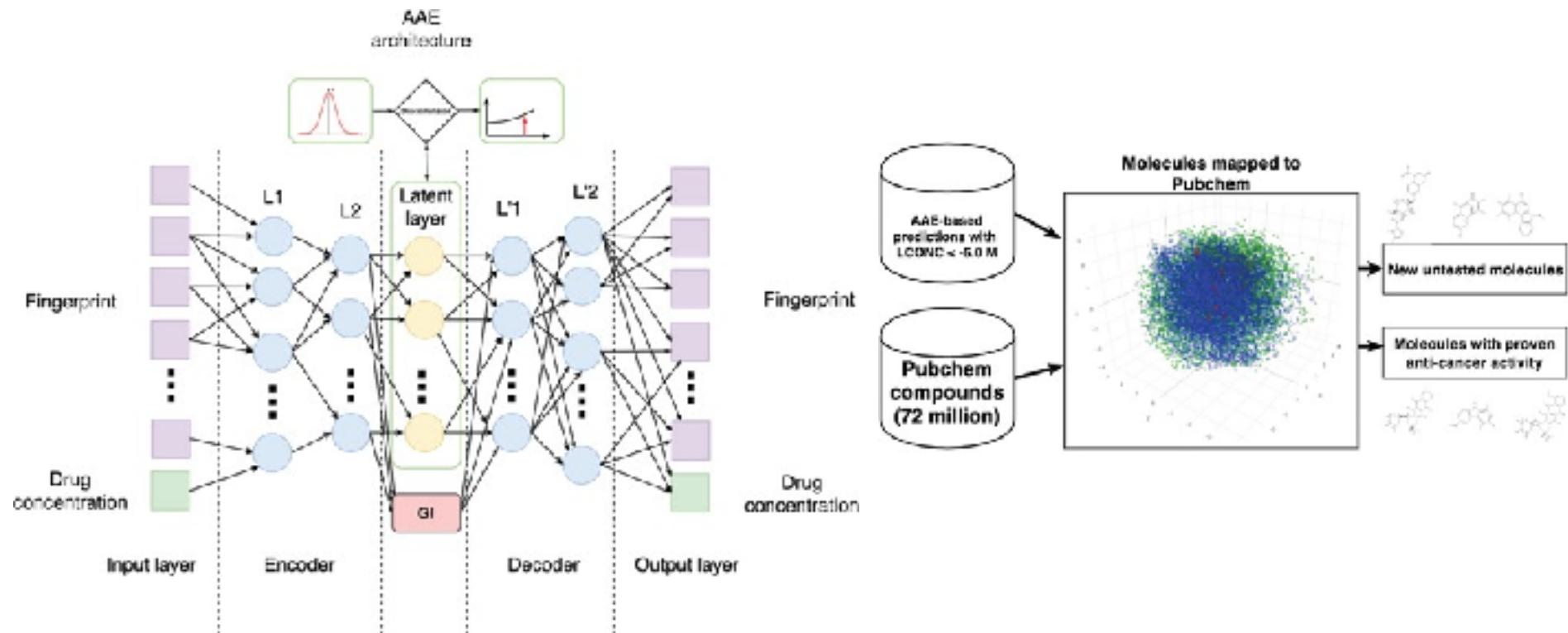
- Top Discriminator
 - Draw samples from one hot encoded labels, $\text{Cat}(y)$
 - Ensures created vector is categorical, $\text{softmax}(z_{\text{prev}}) \rightarrow \text{"one hot"}$
- Bottom Discriminator
 - Same as previous
 - Constrains latent representation
- Supervised Training (disentangled?)
 - Update AAE networks with a few batches
 - Use small labeled mini-batches to update $q(y|x)$ with binary cross entropy
 - Repeat

	MNIST (100)	MNIST (1000)	MNIST (All)
NN Baseline	25.80	8.73	1.25
Adversarial Autoencoders	1.90 (± 0.10)	1.60 (± 0.08)	0.85 (± 0.02)



Other Uses For AAE

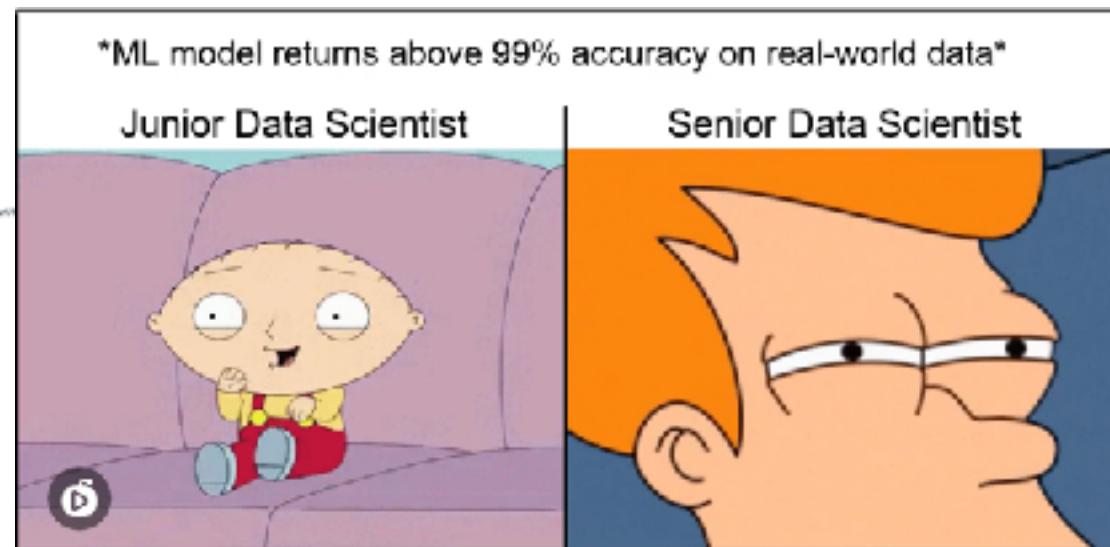
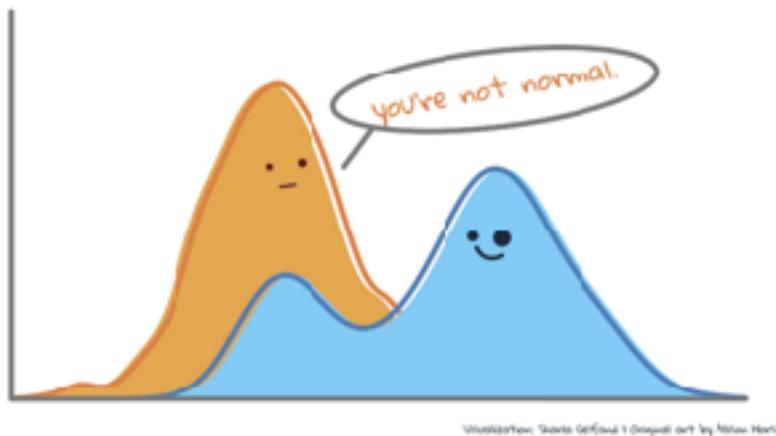
- Molecular Fingerprinting



Kadurin, Artur, Alexander Aliper, Andrey Kazennov, Polina Mamoshina, Quentin Vanhaelen, Kuzma Khrabrov, and Alex Zhavoronkov. "The cornucopia of meaningful leads: Applying deep adversarial autoencoders for new molecule development in oncology." *Oncotarget* 8, no. 7 (2017): 10883.

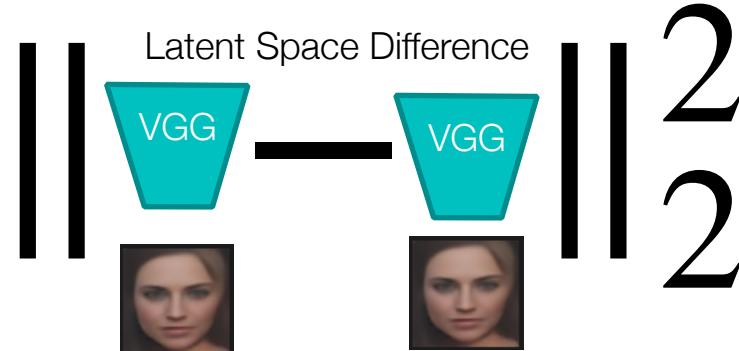
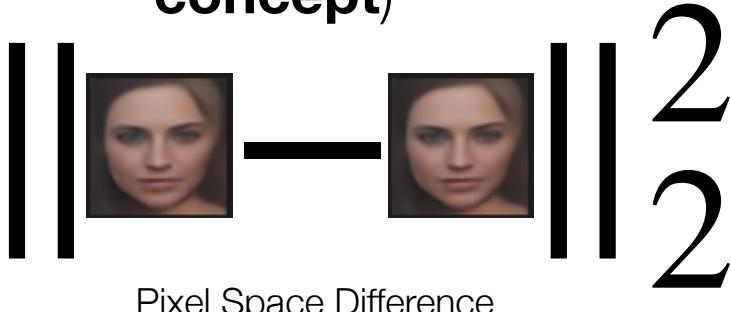


Adversarial Latent Auto-Encoders

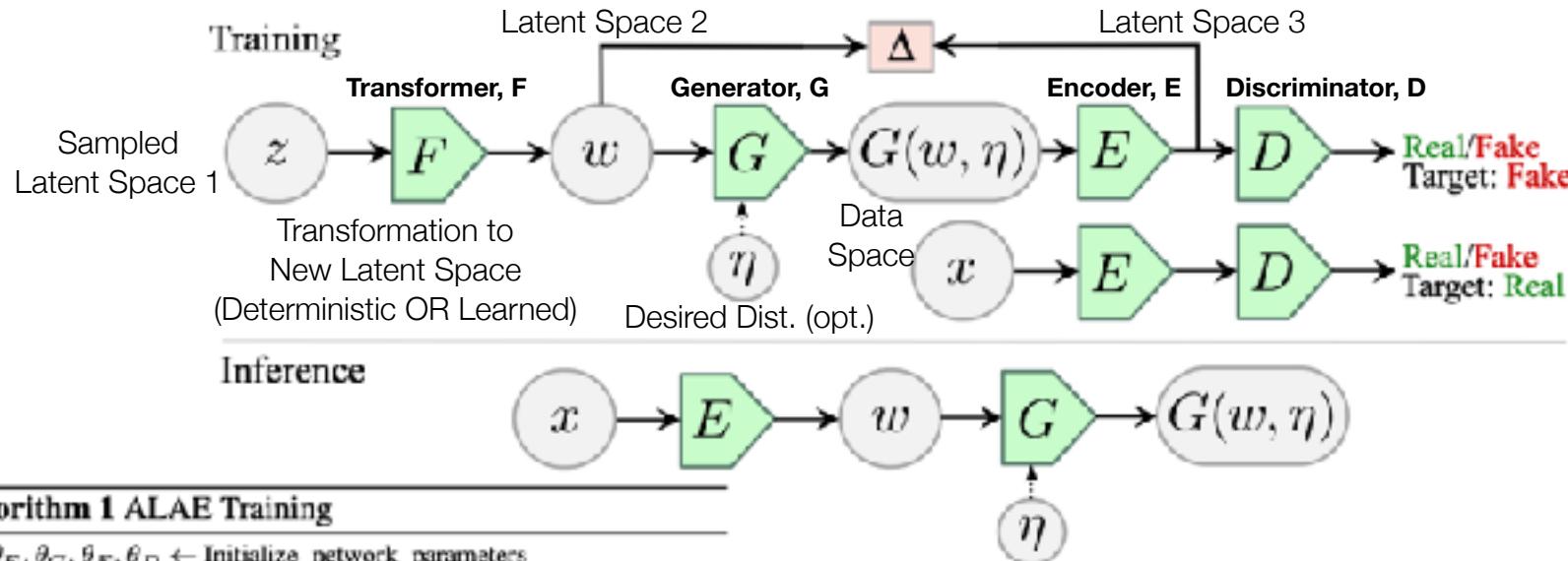


Drawbacks of AAEs

- AAEs capture representational properties of an encoder-decoder pair, with latent space following a specified structure
- Drawbacks:
 - Not entirely generative because we need examples
 - Not guaranteed to create **disentangled** representations (VAE *is better for this...*)
 - Requires that latent space distribution be selected apriori
 - Reconstruction loss calculated in data space, which may not be informative for optimizing latent space (**important concept**)



Adversarial Latent Auto-Encoders, ALAE



Algorithm 1 ALAE Training

```

1:  $\theta_F, \theta_G, \theta_E, \theta_D \leftarrow$  Initialize network parameters
2: while not converged do
3:   Step I. Update  $E$ , and  $D$ 
4:    $x \leftarrow$  Random mini-batch from dataset
5:    $z \leftarrow$  Samples from prior  $\mathcal{N}(0, I)$ 
6:    $L_{adv}^{E,D} \leftarrow$  softplus( $D \circ E \circ G \circ F(z)$ ) + softplus( $-D \circ E(x)$ ) +
 $\frac{\lambda}{2} E_{PD}(x) [\|\nabla D \circ E(x)\|^2]$ 
7:    $\theta_E, \theta_D \leftarrow$  ADAM( $\nabla_{\theta_D, \theta_E} L_{adv}^{E,D}, \theta_D, \theta_E, \alpha, \beta_1, \beta_2$ )
8:   Step II. Update  $F$ , and  $G$ 
9:    $z \leftarrow$  Samples from prior  $\mathcal{N}(0, I)$ 
10:   $L_{adv}^{F,G} \leftarrow$  softplus( $-D \circ E \circ G \circ F(z)$ )
11:   $\theta_F, \theta_G \leftarrow$  ADAM( $\nabla_{\theta_F, \theta_G} L_{adv}^{F,G}, \theta_F, \theta_G, \alpha, \beta_1, \beta_2$ )
12:  Step III. Update  $E$ , and  $G$ 
13:   $z \leftarrow$  Samples from prior  $\mathcal{N}(0, I)$ 
14:   $L_{error}^{E,G} \leftarrow \|F(z) - E \circ G \circ F(z)\|_2^2$ 
15:   $\theta_E, \theta_G \leftarrow$  ADAM( $\nabla_{\theta_E, \theta_G} L_{error}^{E,G}, \theta_E, \theta_G, \alpha, \beta_1, \beta_2$ )
16: end while

```

- Train $G(F)$ to fool discriminator
- Minimize difference F and $E(G(F))$
- Reconstruction loss: x and $E(D(x))$, *never actually calculated, just inferred*
- Therefore, $w \sim E(G(w))$ and $x \sim G(E(x))$

Pidhorskyi, Stanislav, Donald A. Adjeroh, and Gianfranco Doretto. "Adversarial latent autoencoders." In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 14104-14113. 2020.



Adversarial Latent Auto-Encoders, ALAE

$D \rightarrow 1$ = Real data, $D \rightarrow 0$ = Sampled Data (Fake)

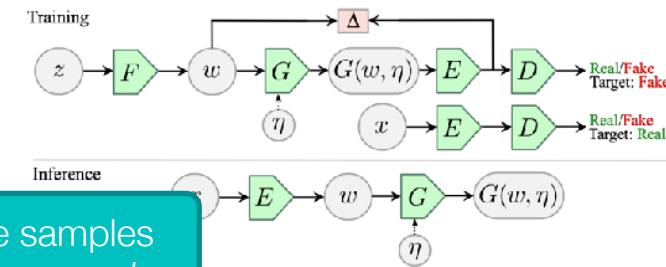
Algorithm 1 ALAE Training

```

1:  $\theta_F, \theta_G, \theta_E, \theta_D \leftarrow$  Initialize network parameters
2: while not converged do
3:
4:   Related to Nash Equilibrium,  
   but shifted...
5:
6:    $L_{adv}^{E,D} \leftarrow \text{softplus}(D \circ E \circ G \circ F(z)) + \text{softplus}(-D \circ E(x)) +$   

 $\frac{\gamma}{2} \mathbb{E}_{P_D(x)} [\|\nabla D \circ E(x)\|^2]$ 
7:    $\theta_E, \theta_D \leftarrow \text{ADAM}(\nabla_{\theta_D, \theta_E} L_{adv}^{E,D}, \theta_D, \theta_E, \alpha, \beta_1, \beta_2)$ 
8:   Step II. Update  $F$ , and  $G$ 
9:    $z \leftarrow$  Samples from prior  $\mathcal{N}(0, I)$ 
10:   $L_{adv}^{F,G} \leftarrow \text{softplus}(D \circ E \circ G \circ F(z)) + \text{softplus}(-D \circ E(x))$ 
11:   $\theta_F, \theta_G \leftarrow \text{ADAM}(\nabla_{\theta_F, \theta_G} L_{adv}^{F,G}, \theta_F, \theta_G, \alpha, \beta_1, \beta_2)$ 
12:  Based On Wasserstein Distance,  
... which is really helpful...
13:   $z \leftarrow$  Samples from prior  $\mathcal{N}(0, I)$ 
14:   $L_{error}^{E,G} \leftarrow \|F(z) - E \circ G \circ F(z)\|_2^2$ 
15:   $\theta_E, \theta_G \leftarrow \text{ADAM}(\nabla_{\theta_E, \theta_G} L_{error}^{E,G}, \theta_E, \theta_G, \alpha, \beta_1, \beta_2)$ 
16: end while

```



E,D: Detect fake samples
minimize D for fake samples

E,D: Detect real samples
minimize $-D$ for real samples

E,D: Gradient Penalty
Keep Gradient Magnitude Small
Keep in mind for later

F,G: Try to fool discriminator,
minimize $-D$ for fake samples

E,G: Keep latent spaces similar

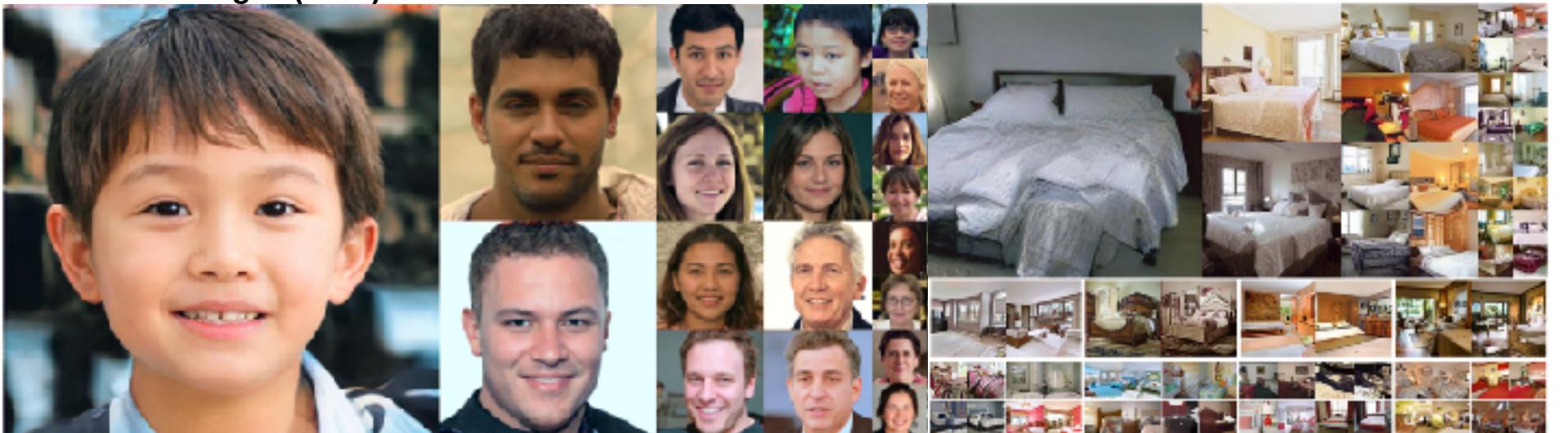


Do we really need to learn this?

Reconstructions



Generated Images (Fake)



Pidhorskyi, Stanislav, Donald A. Adjeroh, and Gianfranco Doretto. "Adversarial latent autoencoders." In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 14104-14113. 2020.



Yep, We need to study GANs

- The ALAE used some notation and processes that we need to study in order to understand:
 - why these are advantageous
 - how to do them properly
 - the tradeoffs and computational cost
- Therefore, the remaining topics:
 - Nash Equilibrium (Vanilla GAN)
 - GAN Training Tricks
 - Least Squares GAN
 - Wasserstein GAN
 - Other Tips and Tricks leading to BigGAN, StyleGAN
 - And Others!



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, we need only a transformation algorithm (generator)
 - Transformer will be a...
 - Neural Network (of course!)
- Transformer is equivalent to “complex” sampling method
- How to optimize and provide training examples?
 - Use two Neural Networks!
 - ... Deep Learning researchers are like the chiropractors of the machine learning community...



Generative Adversarial Network

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

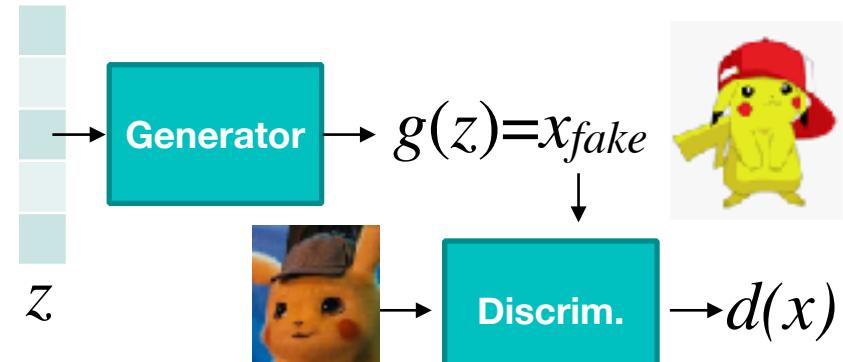
$$\hat{w} = \arg \min_g \max_d v(g, d)$$

- 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



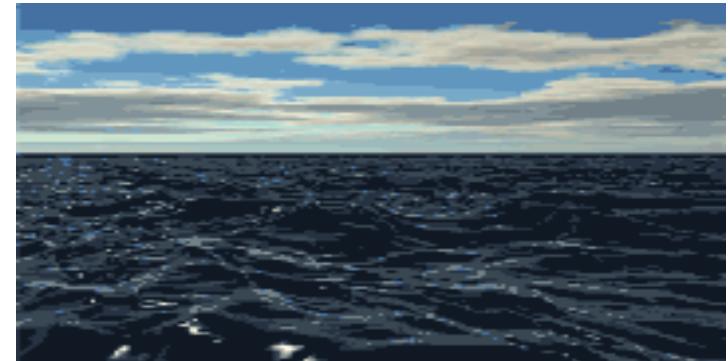
Taking turns: Nash Equilibrium

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

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

- If only the problem was convex! This would be easy to solve...

- But $v(g, d)$ is not convex
 - Saddle points everywhere
 - Like optimizing in an Ocean



- Taking turns on the gradient descent may never converge (Nash equilibrium is not convergence)
 - **So we have 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
- Ian 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 it is wrong?
 - ◆ not trying to win, just make the other person lose

$$\max_g \mathbf{E}_{x \leftarrow g(z)} [\log(d(x_{fake}))] \quad \text{No longer a zero sum game?!"}$$

**Now we are not guaranteed convergence or equilibrium when training the GAN...
But it trains...**



Final Loss Functions

- Discriminator:

Freeze Generator Weights

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

Same as minimizing the binary cross entropy

of the model with: (1) real data=1 and (2) generated data=0!

- Generator:

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

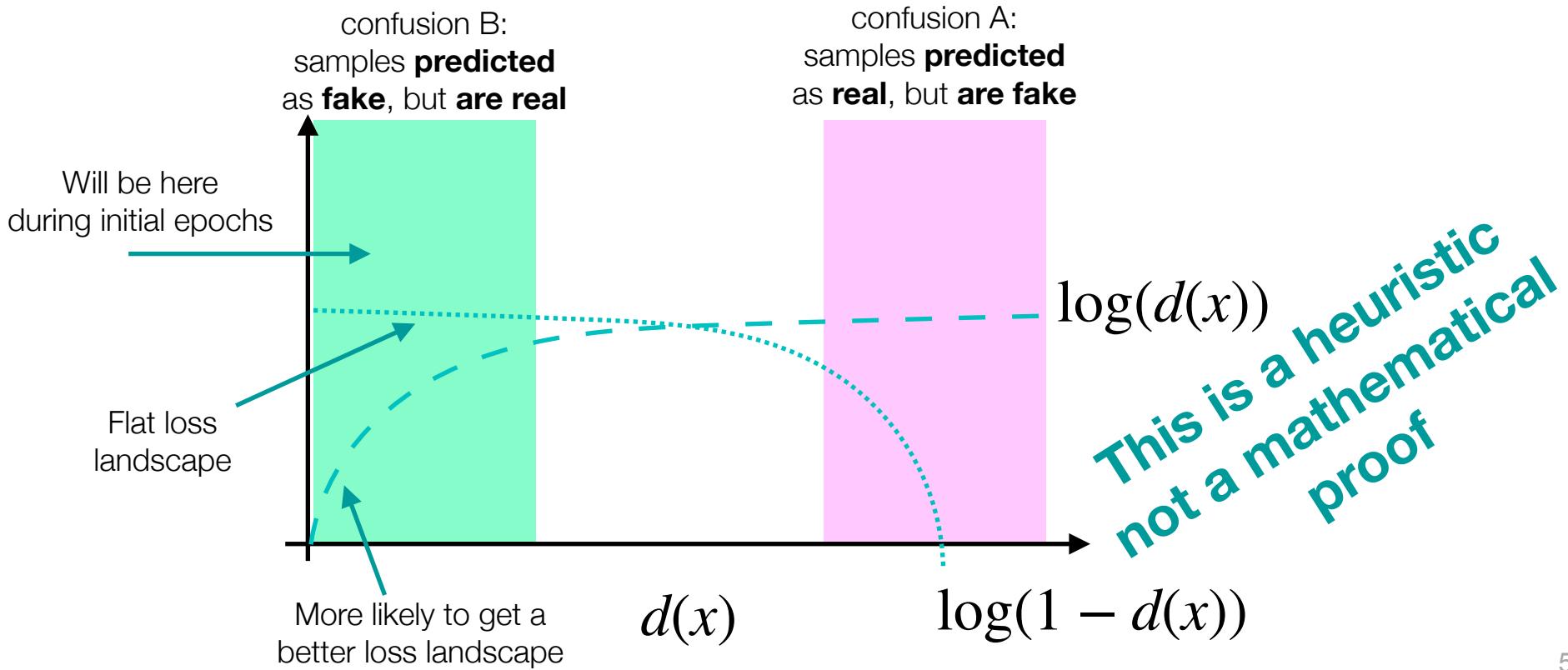
Freeze Discriminator Weights

**Same as minimizing the binary cross entropy
of “mislabeled” generated data=1!**



Intuition for why this really helps...

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



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

GANs



Next Time:
Practical GANs
Reading: Chollet CH8

