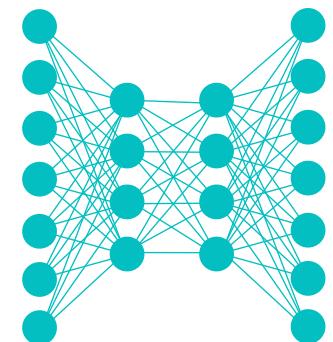


# Lecture Notes for **Neural Networks** **and Machine Learning**



## Adversarial Auto Encoders



# Logistics and Agenda

- Logistics
  - None
- Agenda
  - Student Paper Presentation
  - VAEs, Last time
  - AAE
  - Simple Generative Adversarial Networks
  - GANs Demo (next time)



# Paper Presentation

---

## **Masked Autoencoders As Spatiotemporal Learners**

---

**Christoph Feichtenhofer\***   **Haoqi Fan\***   **Yanhai Li**   **Kaiming He**  
Meta AI, FAIR

[https://github.com/facebookresearch/mae\\_st](https://github.com/facebookresearch/mae_st)



# Last Time: VAEs

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

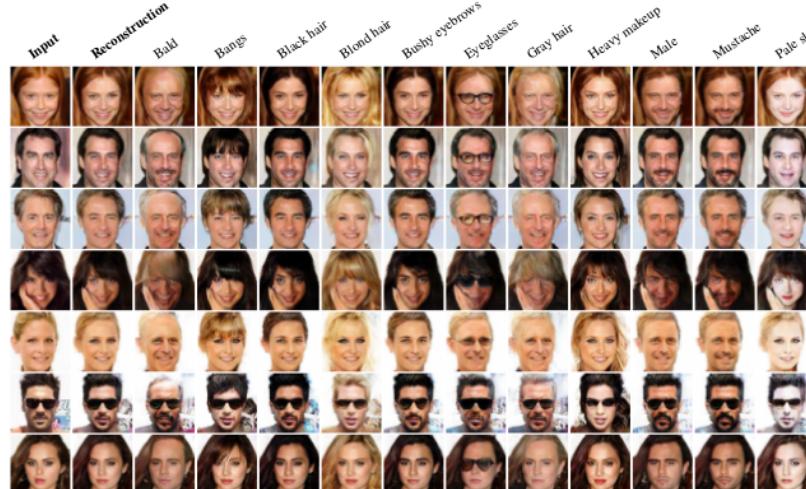
Sampling from variational auto encoder using MNIST



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

Geoff Hinton after writing the paper on backprop in 1986



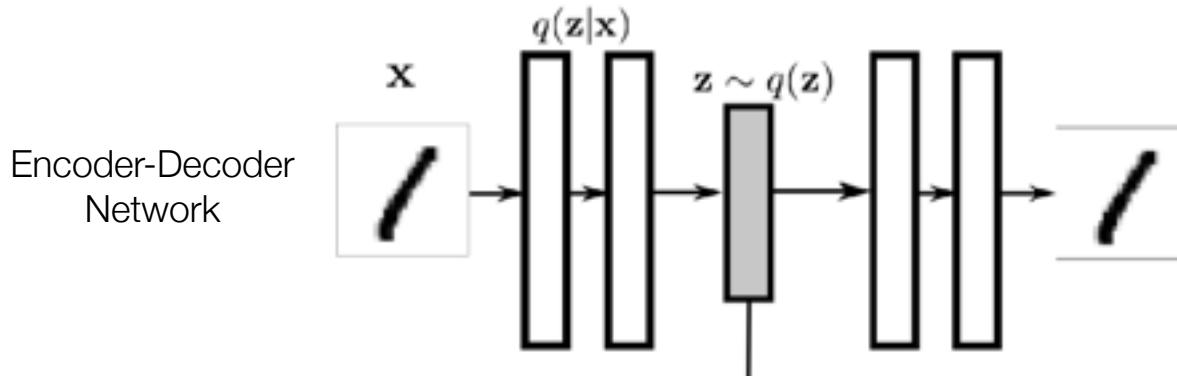
# Do we need something more than VAE?

- 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?



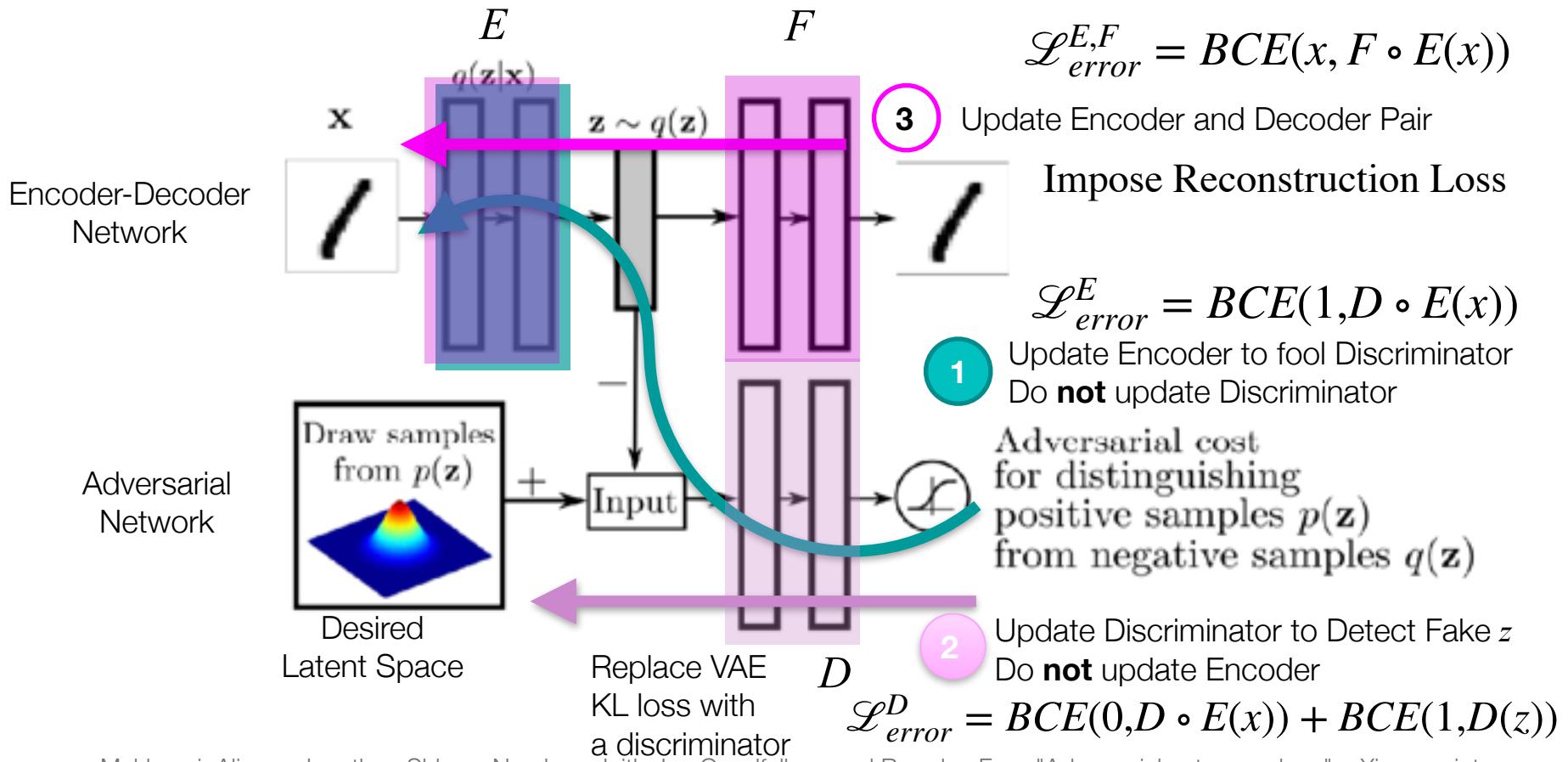
Makhzani, Alireza, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. "Adversarial autoencoders." arXiv preprint arXiv:1511.05644 (2015).

35



# The Main Idea

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

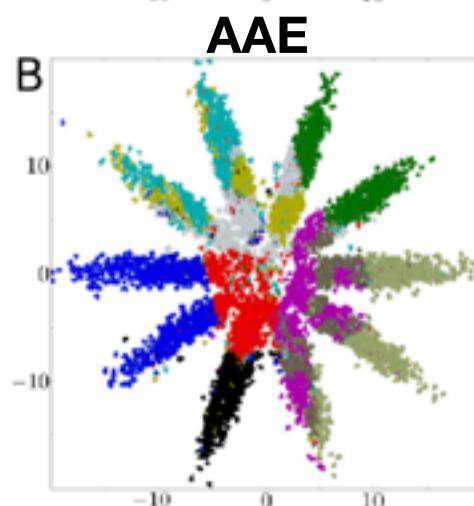
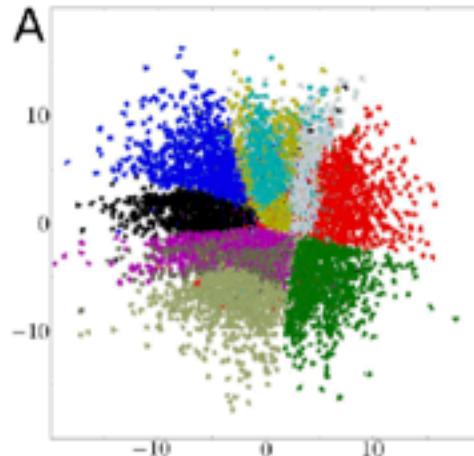


Makhzani, Alireza, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. "Adversarial autoencoders." arXiv preprint arXiv:1511.05644 (2015).

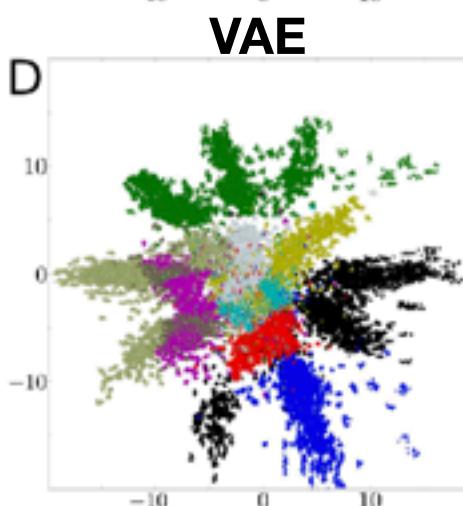
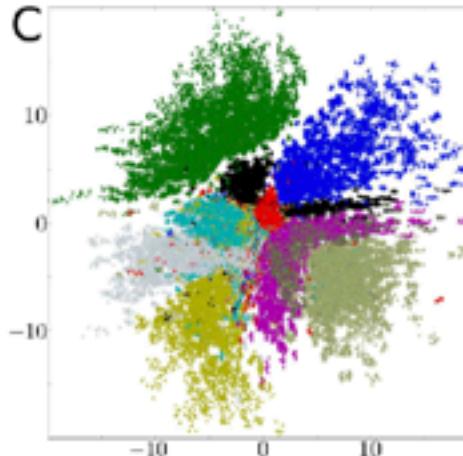


# Arbitrary Prior Distributions

Adversarial Autoencoder



Variational Autoencoder



Manifold of  
Adversarial Autoencoder

E

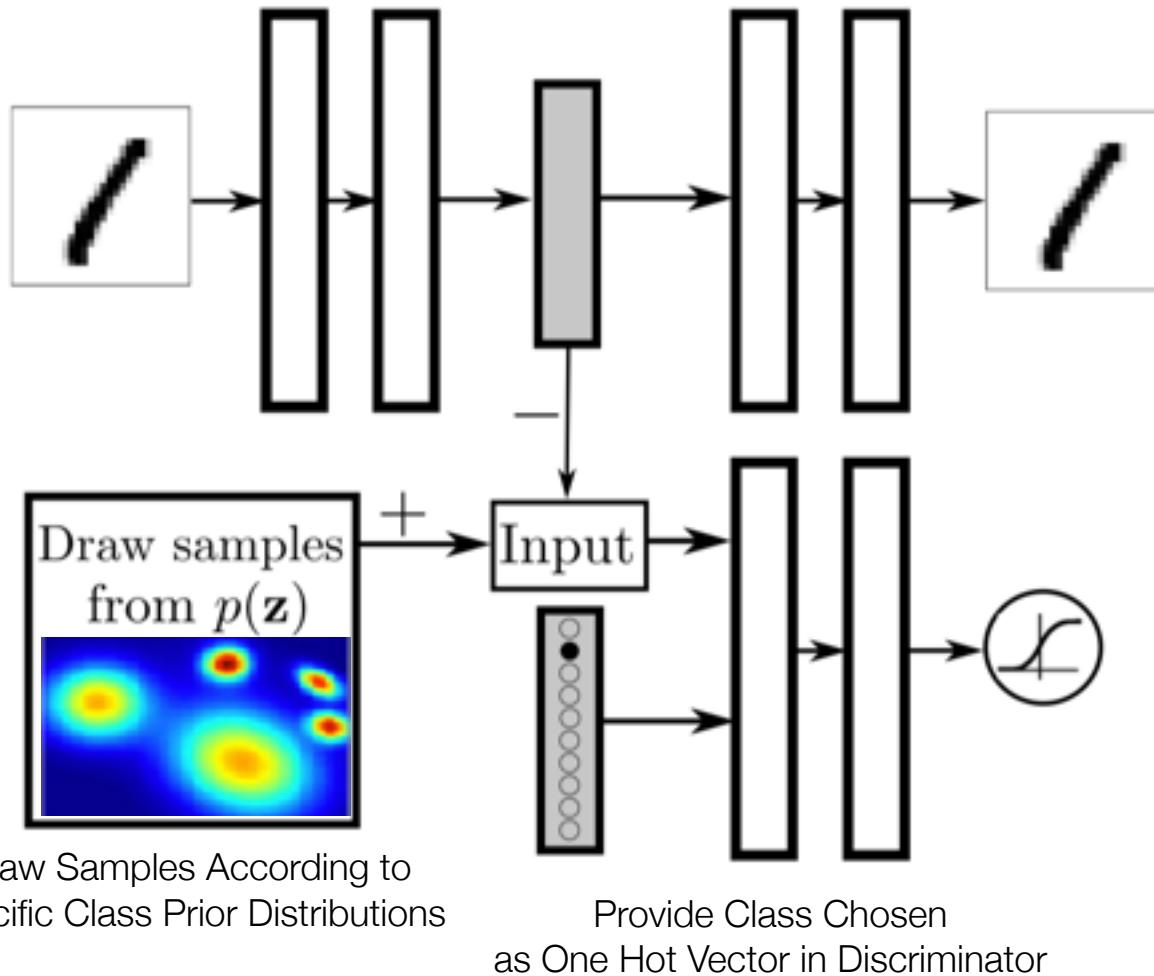
0	0	0	0	0	5	3	3	3	3	5	5	5	8	8	8	2
0	0	0	0	0	5	3	3	3	3	5	5	5	8	8	8	2
0	0	0	0	0	5	3	3	3	3	5	5	5	8	2	2	2
0	0	0	0	0	5	3	3	3	3	5	5	5	8	2	2	2
6	0	0	0	0	0	5	3	3	3	3	5	8	8	8	2	2
6	6	0	0	0	5	3	3	3	3	5	5	8	8	8	2	2
6	6	6	6	6	6	5	5	5	5	5	5	8	8	8	2	2
6	6	6	6	6	6	6	6	6	6	5	5	5	8	8	2	2
6	6	6	6	6	6	6	6	6	6	5	5	5	8	8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

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).



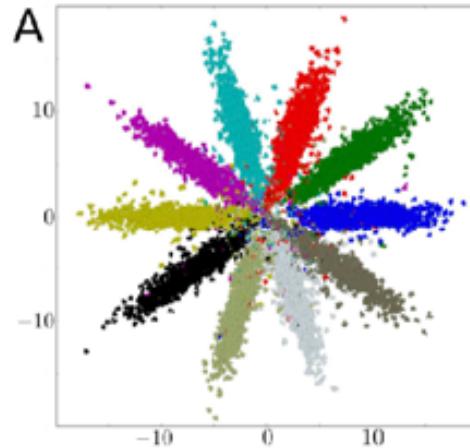
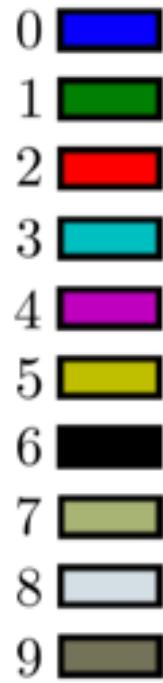
# Sampling From Classes



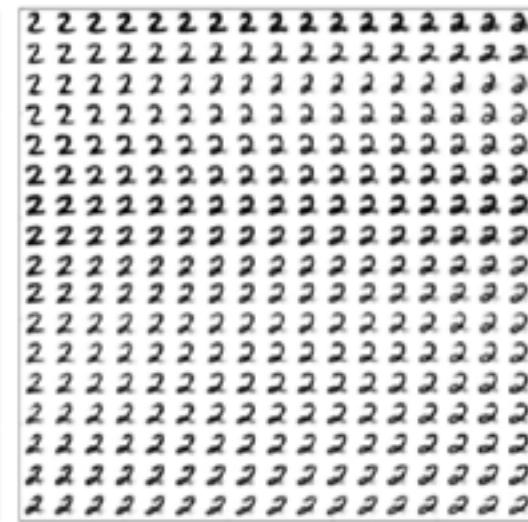
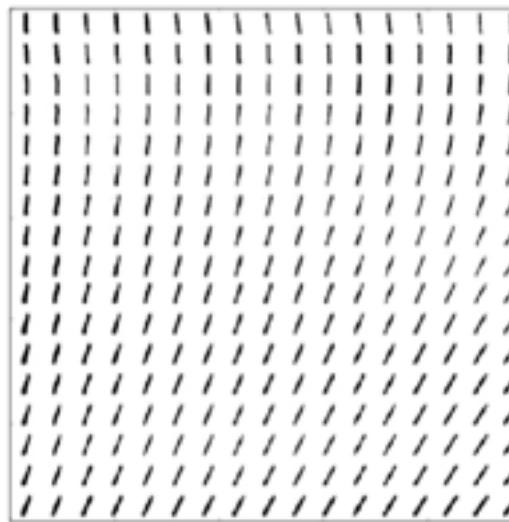
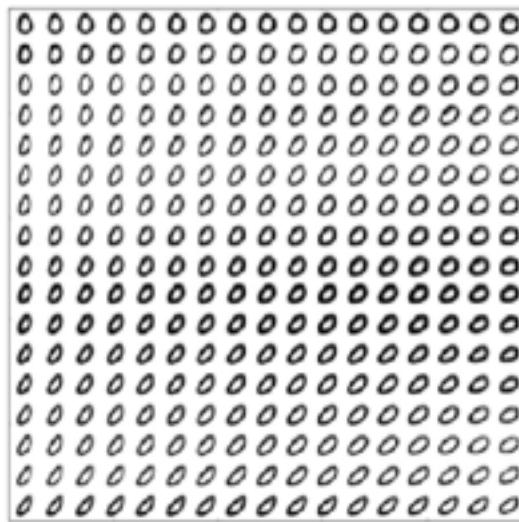
Makhzani, Alireza, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. "Adversarial autoencoders." arXiv preprint arXiv:1511.05644 (2015).



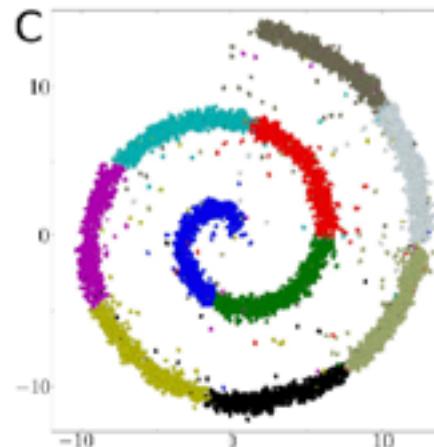
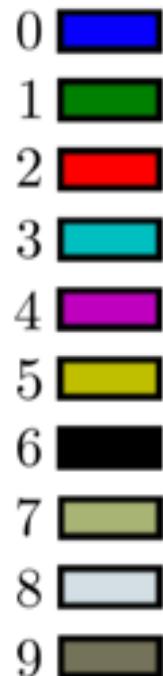
# Conditional Class Latent Spaces



Sample Along Main Axis of the Gaussian Component for Each Digit



# Conditional Class Latent Spaces



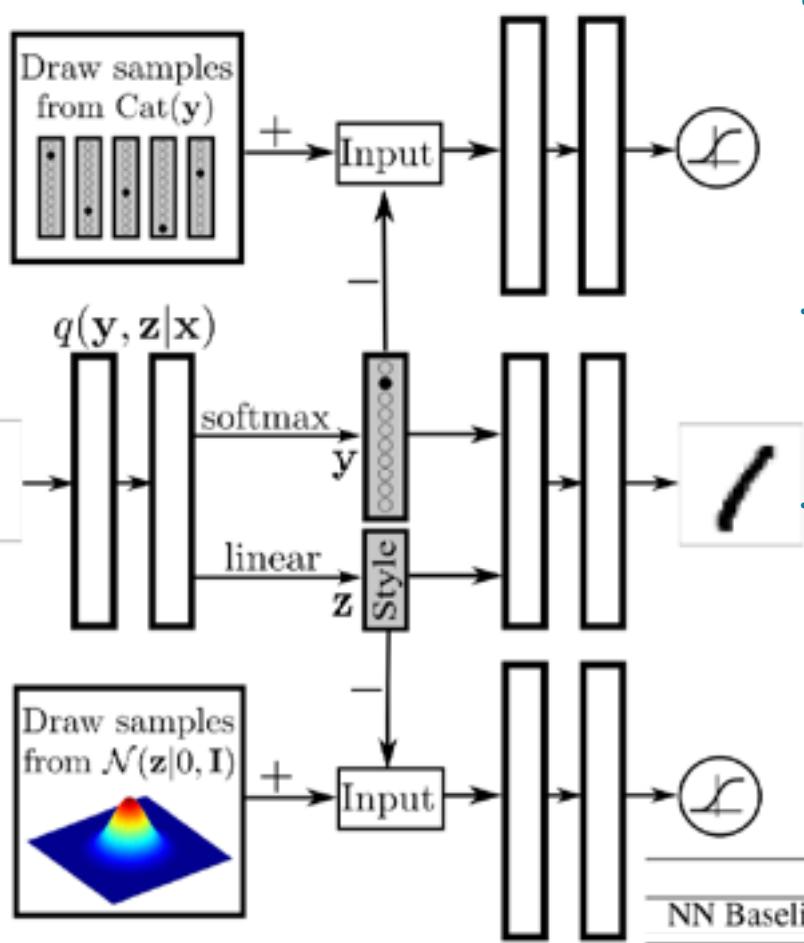
Sample Along Swiss Roll Axis



Makhzani, Alireza, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. "Adversarial autoencoders." *arXiv preprint arXiv:1511.05644* (2015).40



# Semi-Supervised Classification



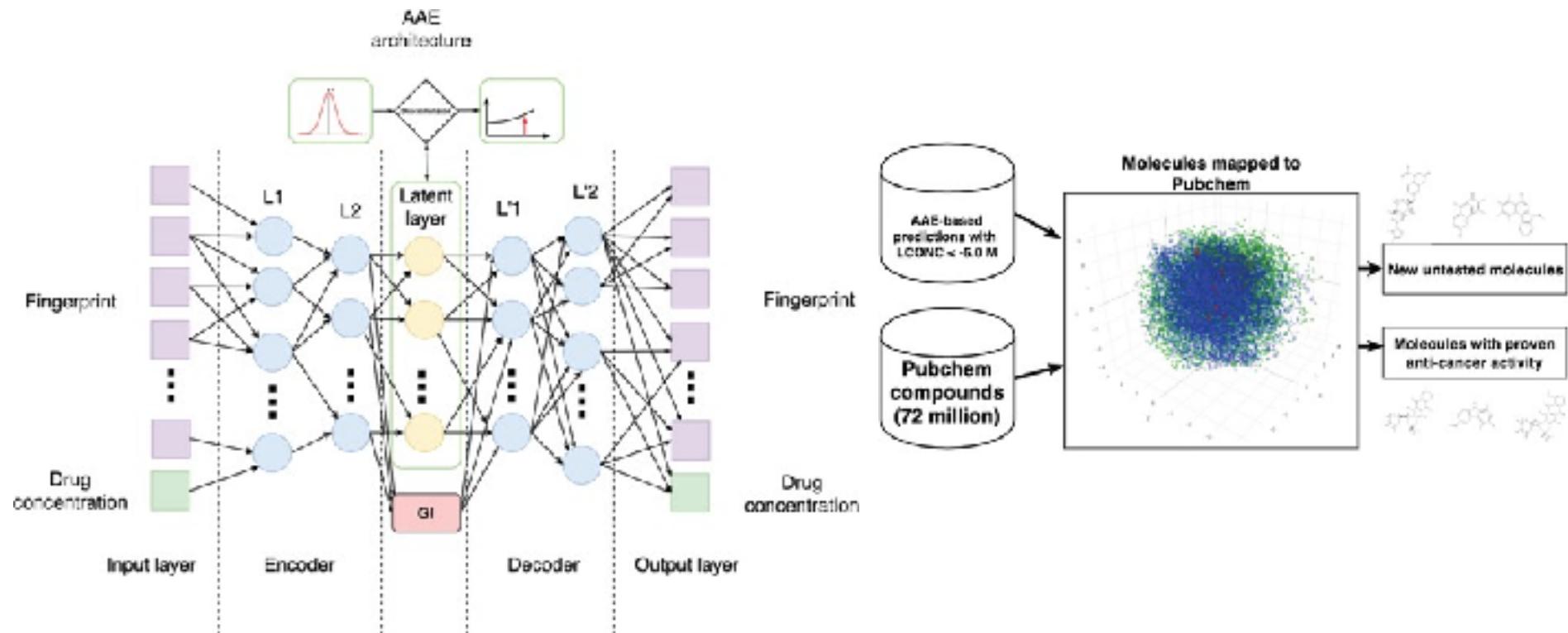
- Top Discriminator
  - Draw samples from one hot encoded labels,  $\text{Cat}(y)$
  - Ensures created vector is categorical,  $\text{softmax}(z_{\text{prev}}) \rightarrow \text{"one hot"}$
- Bottom Discriminator
  - Same as previous
  - Constrains latent representation
- Supervised Training (disentangled?)
  - 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
<b>Adversarial Autoencoders</b>	1.90 ( $\pm 0.10$ )	1.60 ( $\pm 0.08$ )	0.85 ( $\pm 0.02$ )



# Other Uses For AAE

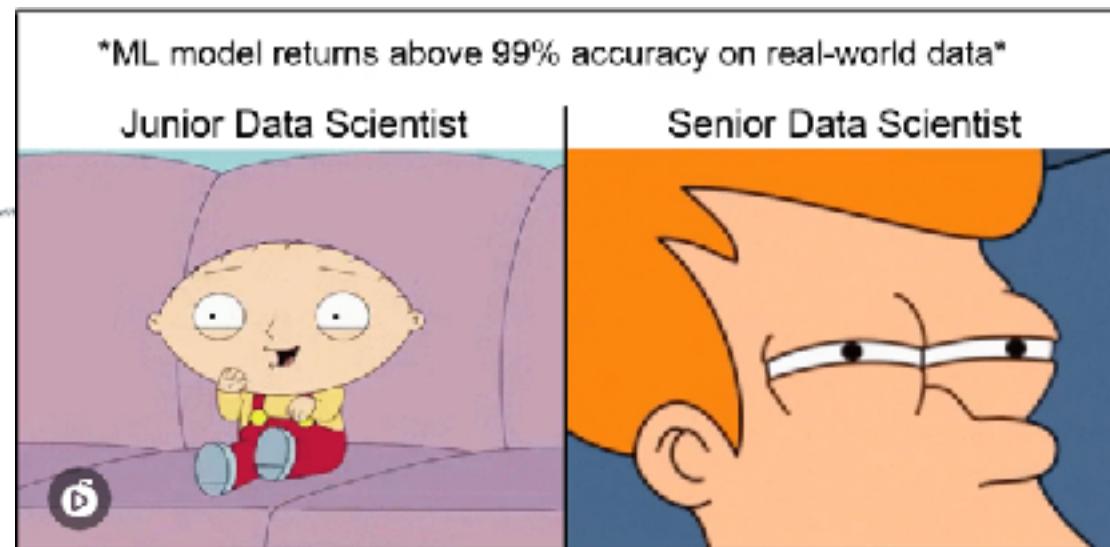
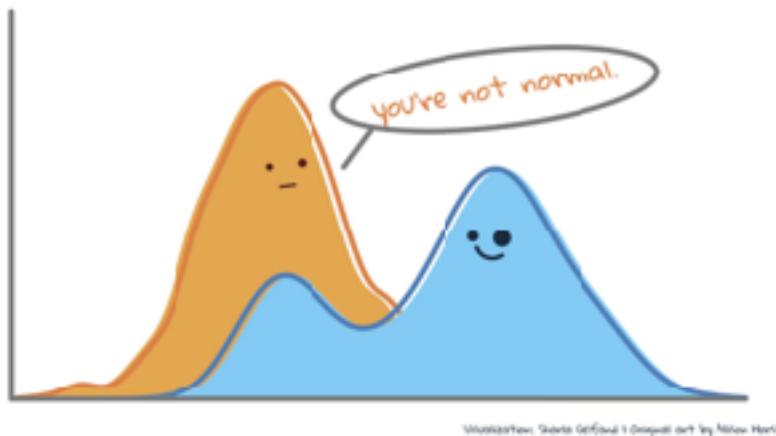
- Molecular Fingerprinting



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.

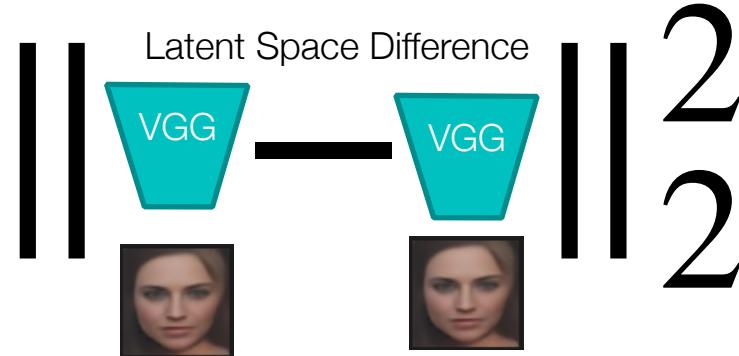
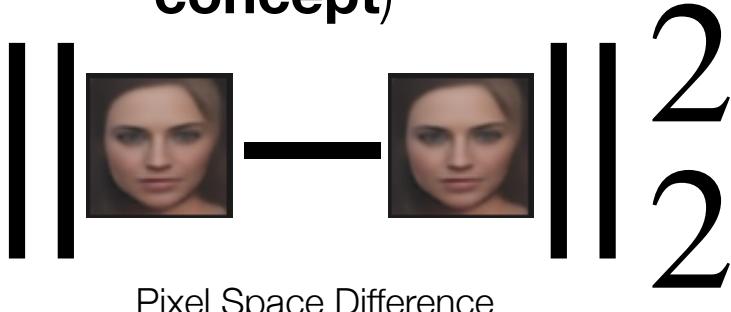


# Adversarial Latent Auto-Encoders

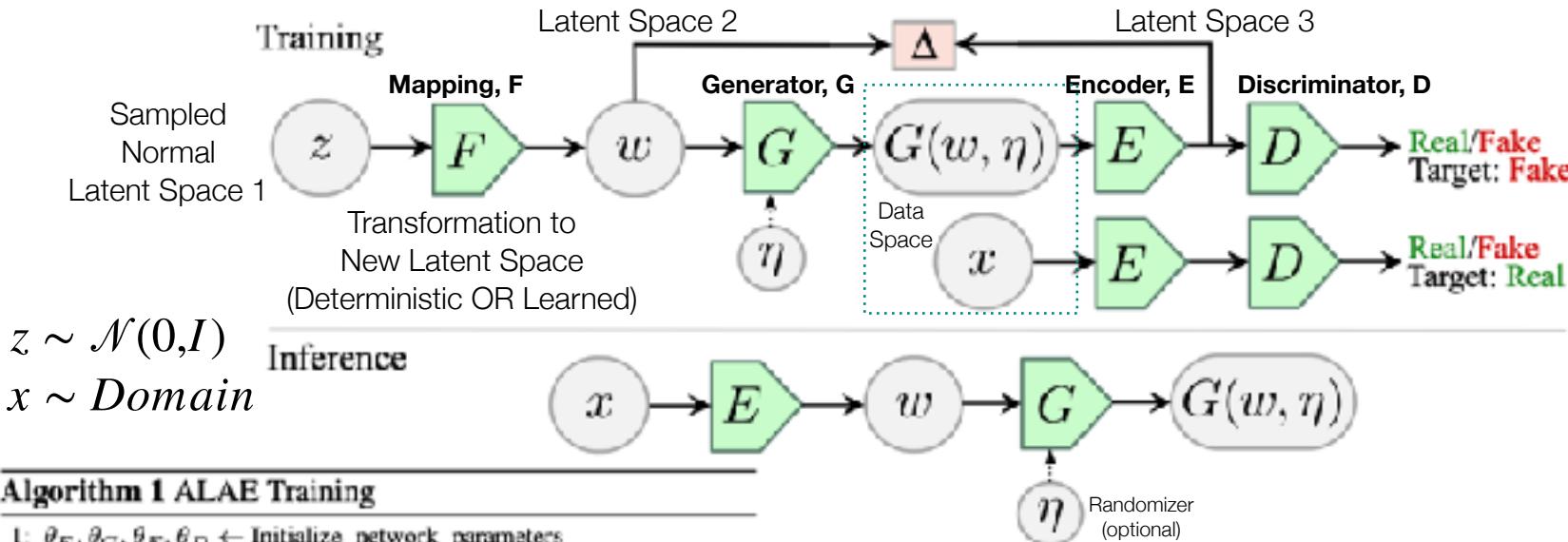


# Drawbacks of AAEs

- 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( $-D \circ E(x)$ ) +
 $\frac{\beta}{2} E_{PD}(x) [\|\nabla D \circ E(x)\|^2]$ 
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 Mapping/Generator to fool discriminator
- Train Encoder/Discriminator to find fake
- Minimize latent space encoder & mapper
- Reconstruction loss:  $x$  and  $G(E(x))$ , *never actually calculated, just inferred*
- Therefore, encoder and mapper are adversaries

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

$$\sigma(x) = \text{softplus} = \log(1 + \exp(x)) \quad \text{smoothed ReLU } [0, \infty)$$

**E,D:** Detect fake samples  
 $\min -D$  ( $\max D$ ) for fake samples

**E,D:** Detect real samples  
 $\min -D$  ( $\max D$ ) for real samples

$$6 : \mathcal{L}_{disc}^{E,D} = \sigma(D \circ E \circ G \circ F(z)) + \sigma(-D \circ E(x)) \\ + \frac{\gamma}{2} \cdot \mathbf{E}[\|\nabla D \circ E(x)\|]$$

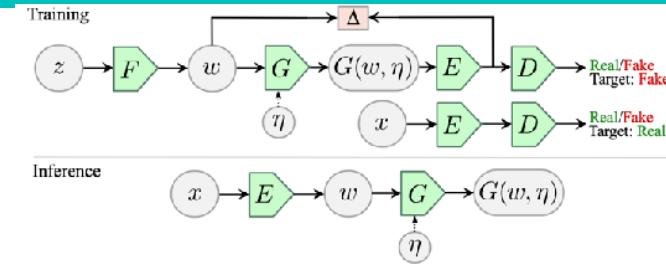
**F,G:** Try to fool discriminator,  
 $\min -D$  ( $\max D$ ) for fake samples

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

$$10 : \mathcal{L}_{gen}^{F,G} = \sigma(-D \circ E \circ G \circ F(z))$$

$$14 : \mathcal{L}_{latent}^{E,G} = \|F(z) - E \circ G \circ F(z)\|^2$$

**E,G:** Keep latent spaces similar



## 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( $-D \circ E(x)$ ) +
 $\frac{\gamma}{2} \mathbf{E}_{p_D(x)} [\|\nabla D \circ E(x)\|^2]$ 
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

```

**6: 10: Based On Wasserstein Distance, ... which is really helpful...**



# 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 (condensed)
  - PyTorch (skipping)
  - Least Squares GAN (skipping)
  - Wasserstein GAN
  - BigGAN, StyleGAN
  - Stable Diffusion



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



ALAEs

**Next Time:**

GANs

**Reading:** Chollet CH8

