

Auto-Encoders

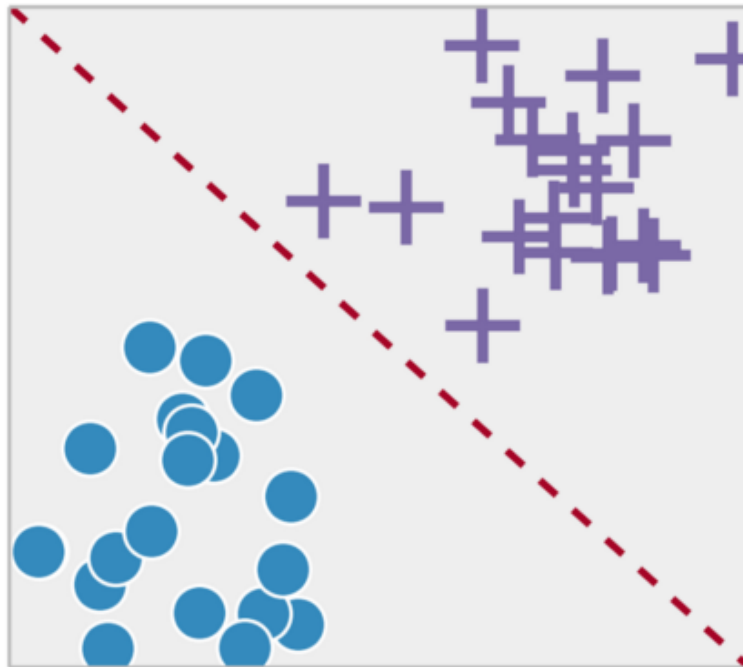
主讲：龙良曲

Outline

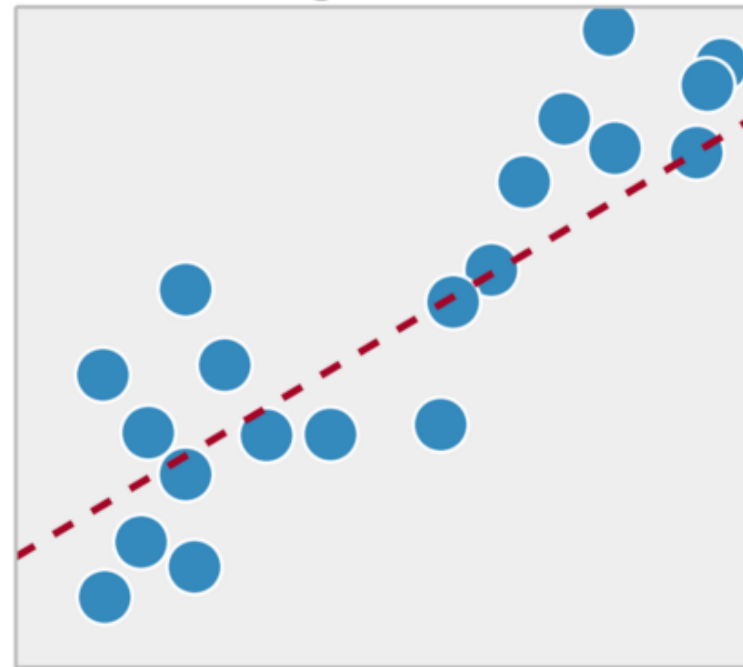


Supervised Learning

Classification



Regression



Massive Unlabeled data



Unsupervised Learning

■ "Pure" Reinforcement Learning (cherry)

- ▶ The machine predicts a scalar reward given once in a while.
- ▶ **A few bits for some samples**

■ Supervised Learning (icing)

- ▶ The machine predicts a category or a few numbers for each input
- ▶ Predicting human-supplied data
- ▶ **10→10,000 bits per sample**

■ Unsupervised/Predictive Learning (cake)

- ▶ The machine predicts any part of its input for any observed part.
- ▶ Predicts future frames in videos
- ▶ **Millions of bits per sample**

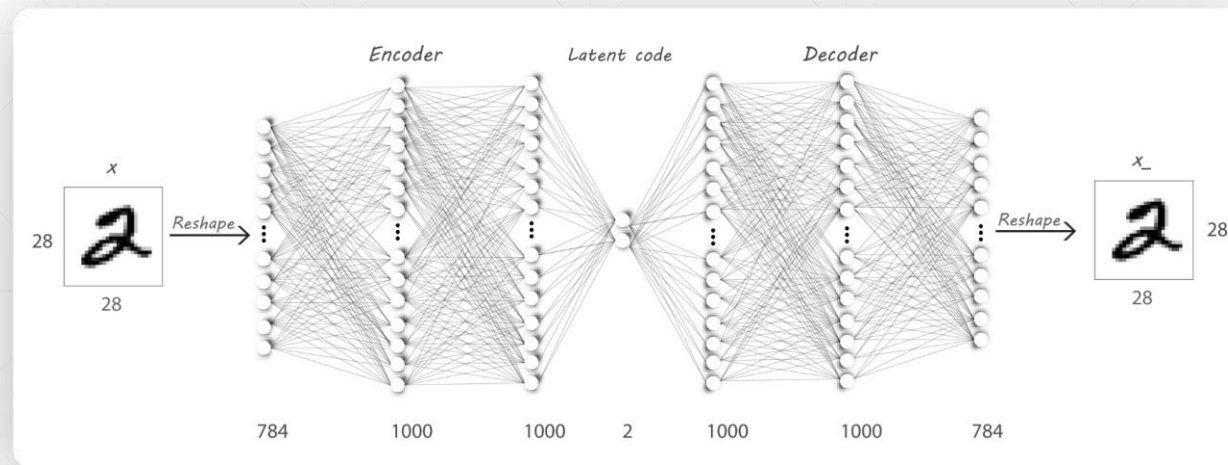
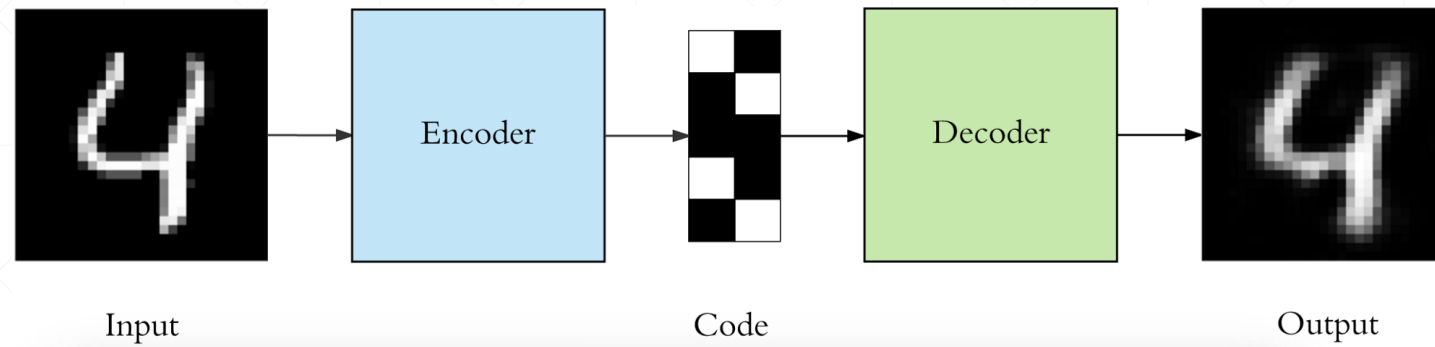


■ (Yes, I know, this picture is slightly offensive to RL folks. But I'll make it up)

Why needed

- Dimension reduction
 - Preprocessing: Huge dimension, say 224x224, is hard to process
 - Visualization: <https://projector.tensorflow.org/>
 - Taking advantages of unsupervised data
 - Compression, denoising, super-resolution ...
-

Auto-Encoders



<https://towardsdatascience.com/applied-deep-learning-part-3-autoencoders-1c083af4d798>

<https://towardsdatascience.com/a-wizards-guide-to-adversarial-autoencoders-part-1-autoencoder-d9a5f8795af4>

How to Train?

- **Loss function** for binary inputs

$$l(f(\mathbf{x})) = - \sum_k (x_k \log(\hat{x}_k) + (1 - x_k) \log(1 - \hat{x}_k))$$

- Cross-entropy error function (reconstruction loss) $f(\mathbf{x}) \equiv \hat{\mathbf{x}}$

- **Loss function** for real-valued inputs

$$l(f(\mathbf{x})) = \frac{1}{2} \sum_k (\hat{x}_k - x_k)^2$$

- sum of squared differences (reconstruction loss)
- we use a linear activation function at the output

PCA V.S. Auto-Encoders

- PCA, which finds the directions of maximal variance in high-dimensional data, select only those axes that have the largest variance.
 - The linearity of PCA, however, places significant limitations on the kinds of feature dimensions that can be extracted.
-

PCA V.S. Auto-Encoders

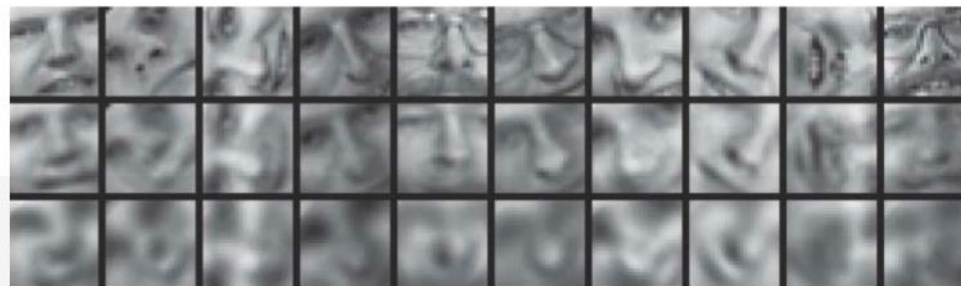


Real data

30-d deep autoencoder

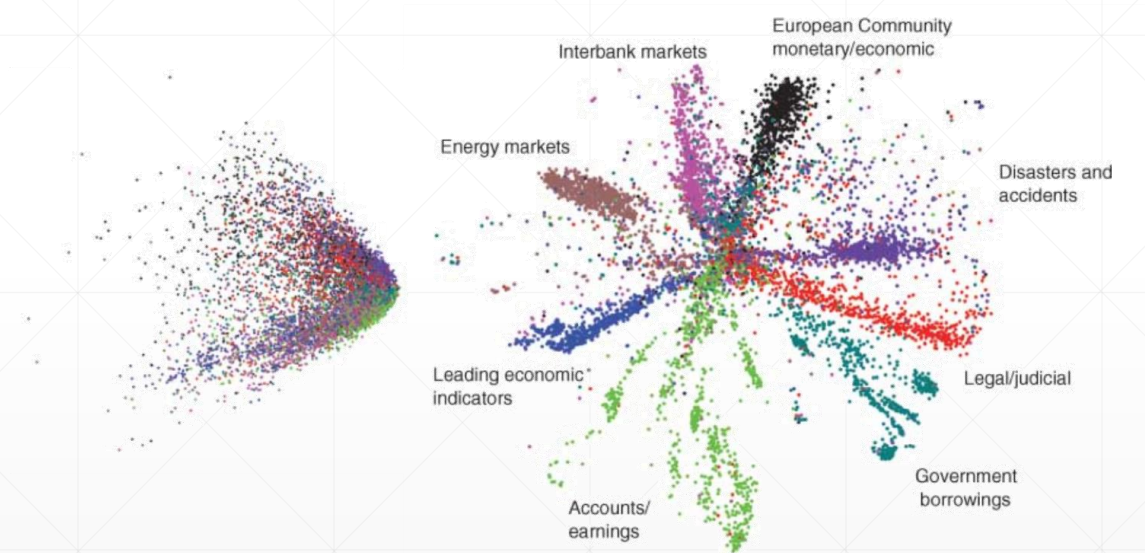
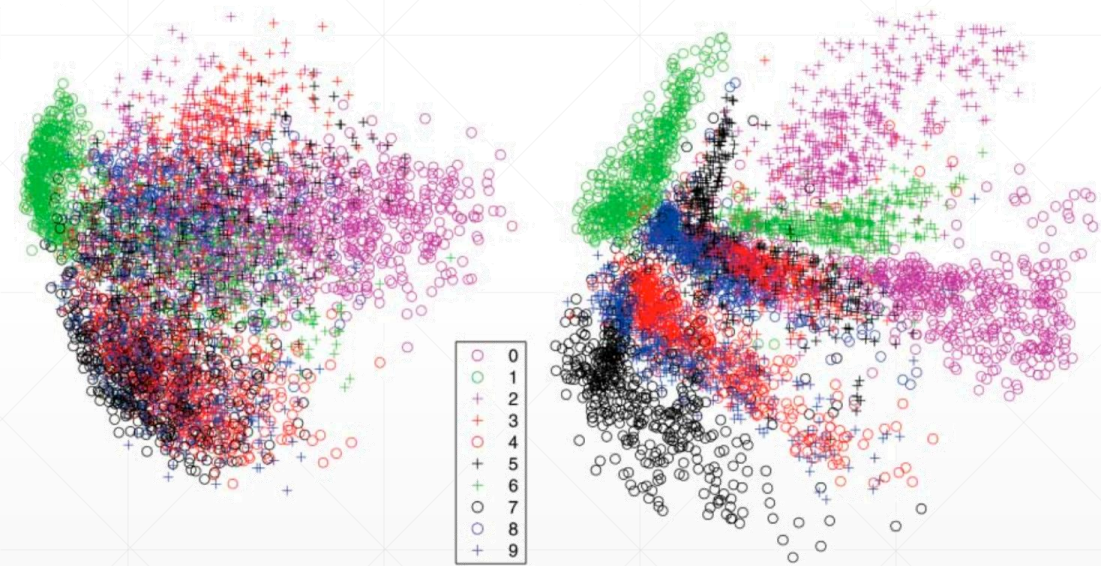
30-d logistic PCA

30-d PCA

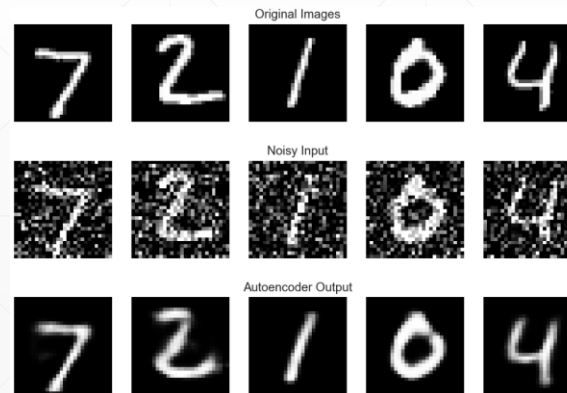
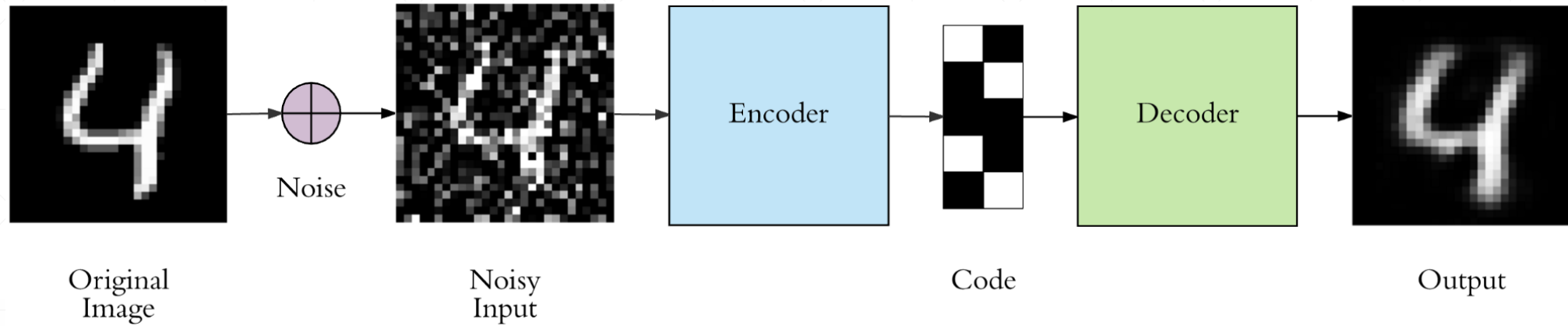


A comparison of reconstruction by an autoencoder (middle) and PCA (bottom) to original image inputs (top)

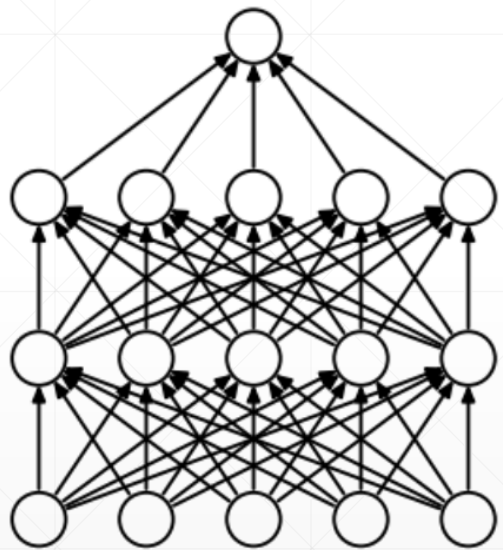
PCA V.S. Auto-Encoders



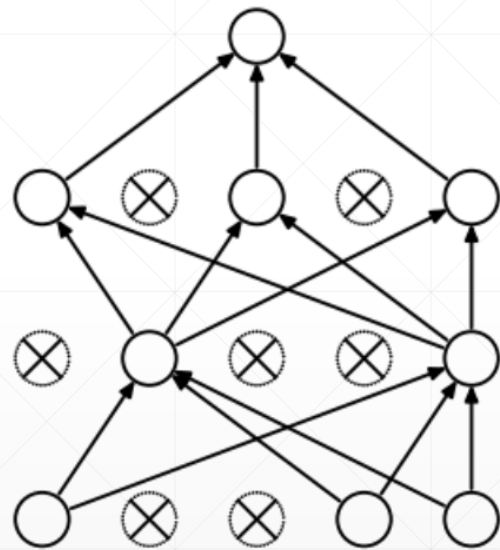
Denoising AutoEncoders



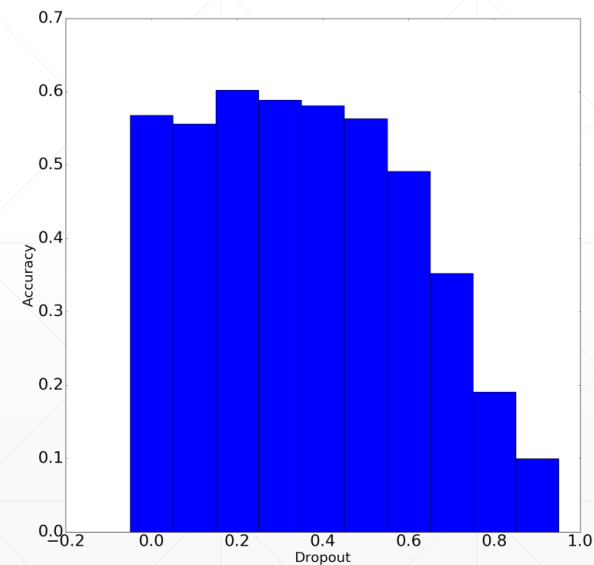
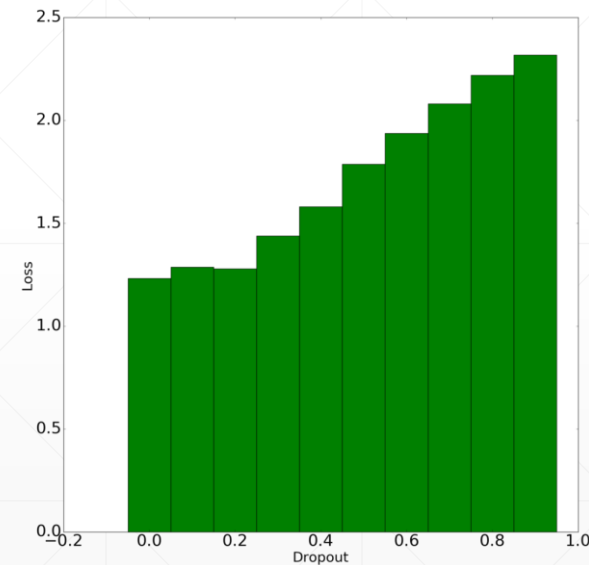
Dropout AutoEncoders



(a) Standard Neural Net



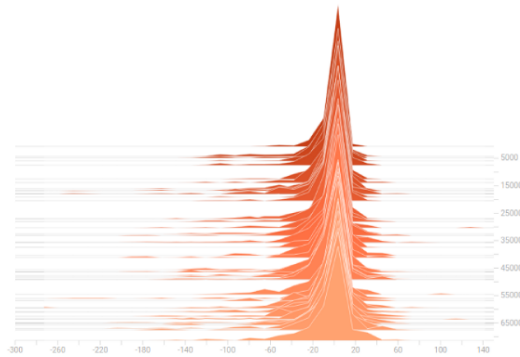
(b) After applying dropout.



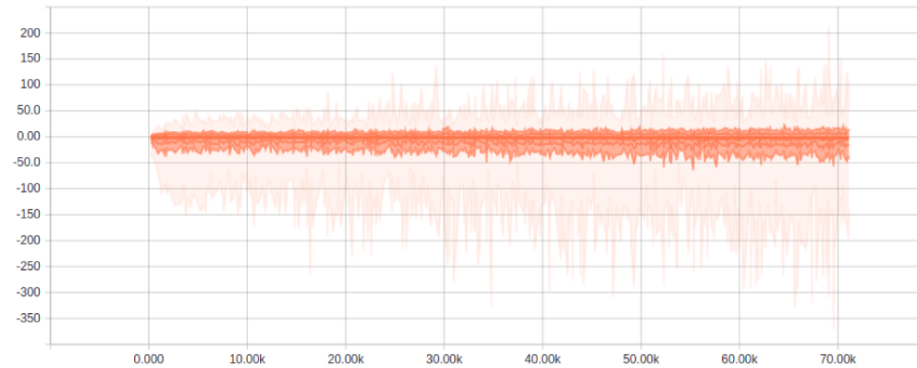
Adversarial AutoEncoders

- Distribution of hidden code

Encoder histogram

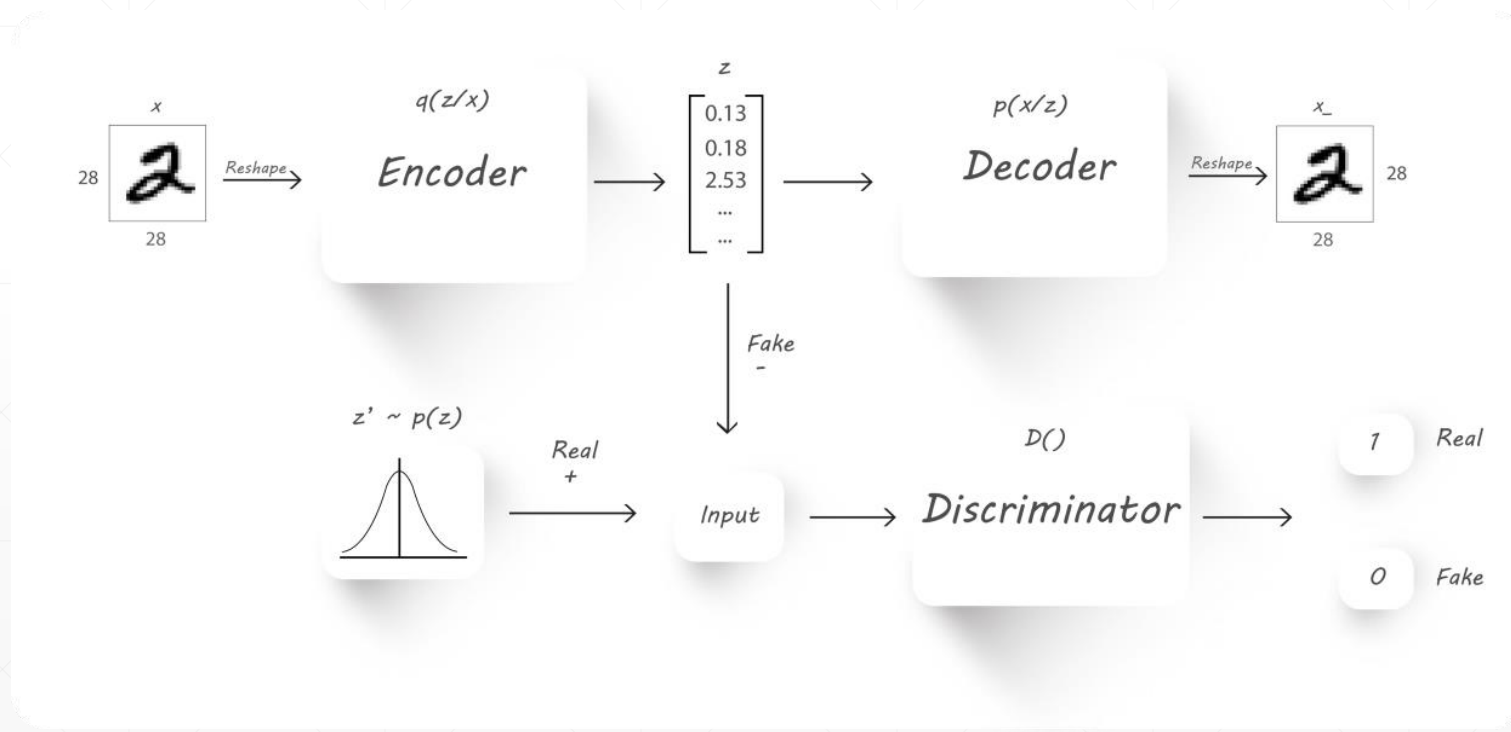


Encoder Distribution



Adversarial AutoEncoders

- Give more details after GAN



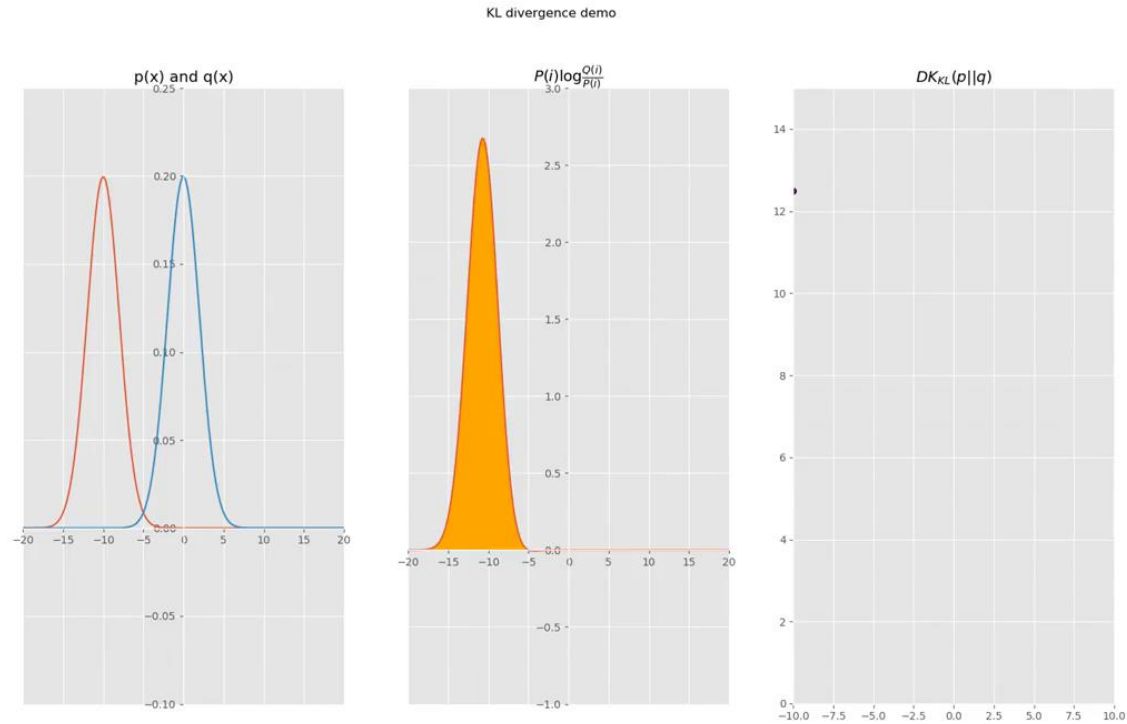
Another Approach: $q(z) \rightarrow p(z)$

- Explicitly enforce

$$l_i(\theta, \phi) = -E_{z \sim q_\theta(z|x_i)} [\log p_\phi(x_i|z)] + KL(q_\theta(z|x_i) || p(z))$$

$$KL(P||Q) = \int_{-\infty}^{\infty} p(x) \log \frac{p(x)}{q(x)} dx$$

Intuitively comprehend $KL(p \circledast q)$



Maximize Likelihood

$$E_{z \sim q_{\theta}(z|x_i)} [\log p_{\phi}(x_i|z)]$$

- **Loss function** for binary inputs

$$l(f(\mathbf{x})) = - \sum_k (x_k \log(\hat{x}_k) + (1 - x_k) \log(1 - \hat{x}_k))$$

- Cross-entropy error function (reconstruction loss) $f(\mathbf{x}) \equiv \hat{\mathbf{x}}$

- **Loss function** for real-valued inputs

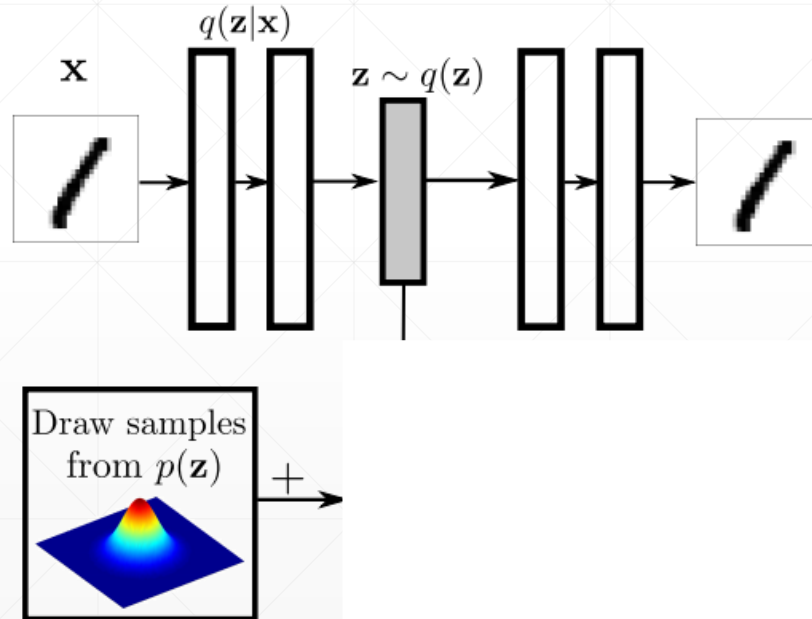
$$l(f(\mathbf{x})) = \frac{1}{2} \sum_k (\hat{x}_k - x_k)^2$$

- sum of squared differences (reconstruction loss)
- we use a linear activation function at the output

Minimize KL Divergence

- Evidence Lower BOund

$$KL(q_{\theta}(z|x_i)||p(z))$$



How to compute KL between $q(z)$ and $p(z)$

$$p(z_i) \sim N(\mu_1, \sigma_1^2)$$

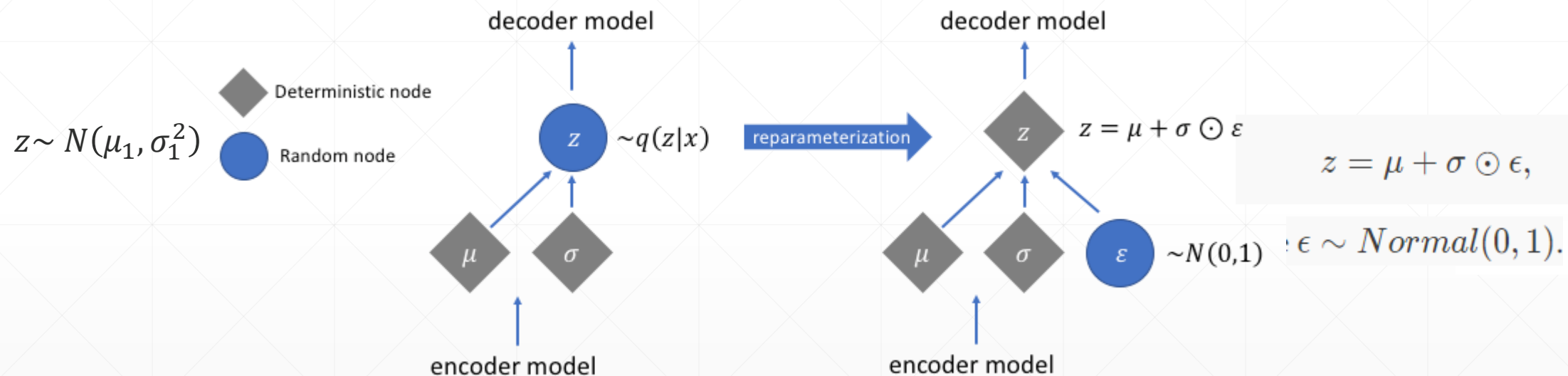
$$q(z_i) \sim N(\mu_2, \sigma_2^2)$$

$$KL(p, q) = - \int p(x) \log q(x) dx + \int p(x) \log p(x) dx$$

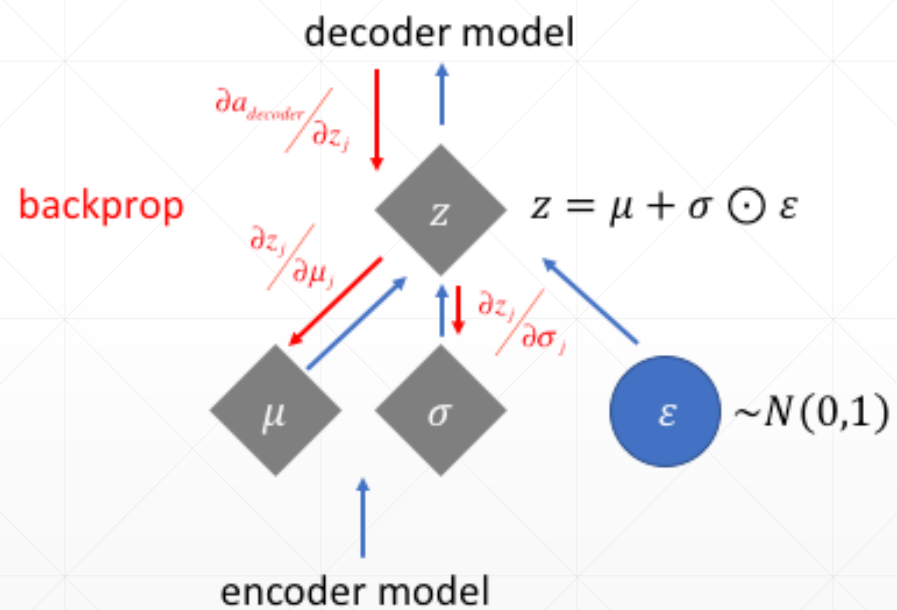
$$= \frac{1}{2} \log(2\pi\sigma_2^2) + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2} (1 + \log 2\pi\sigma_1^2)$$

$$= \log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}$$

Sample() is not differentiable



Reparameterization trick



Penalizing reconstruction loss encourages the distribution to describe the input



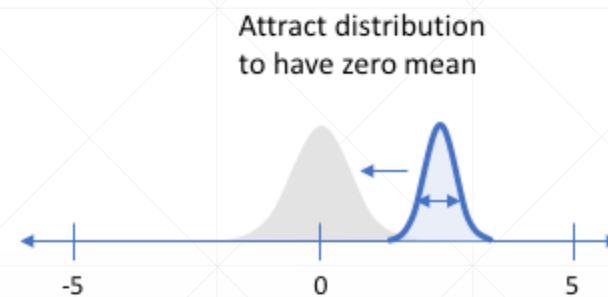
Our distribution deviates from the prior to describe some characteristic of the data

Without regularization, our network can “cheat” by learning narrow distributions



With a small enough variance, this distribution is effectively only representing a single value

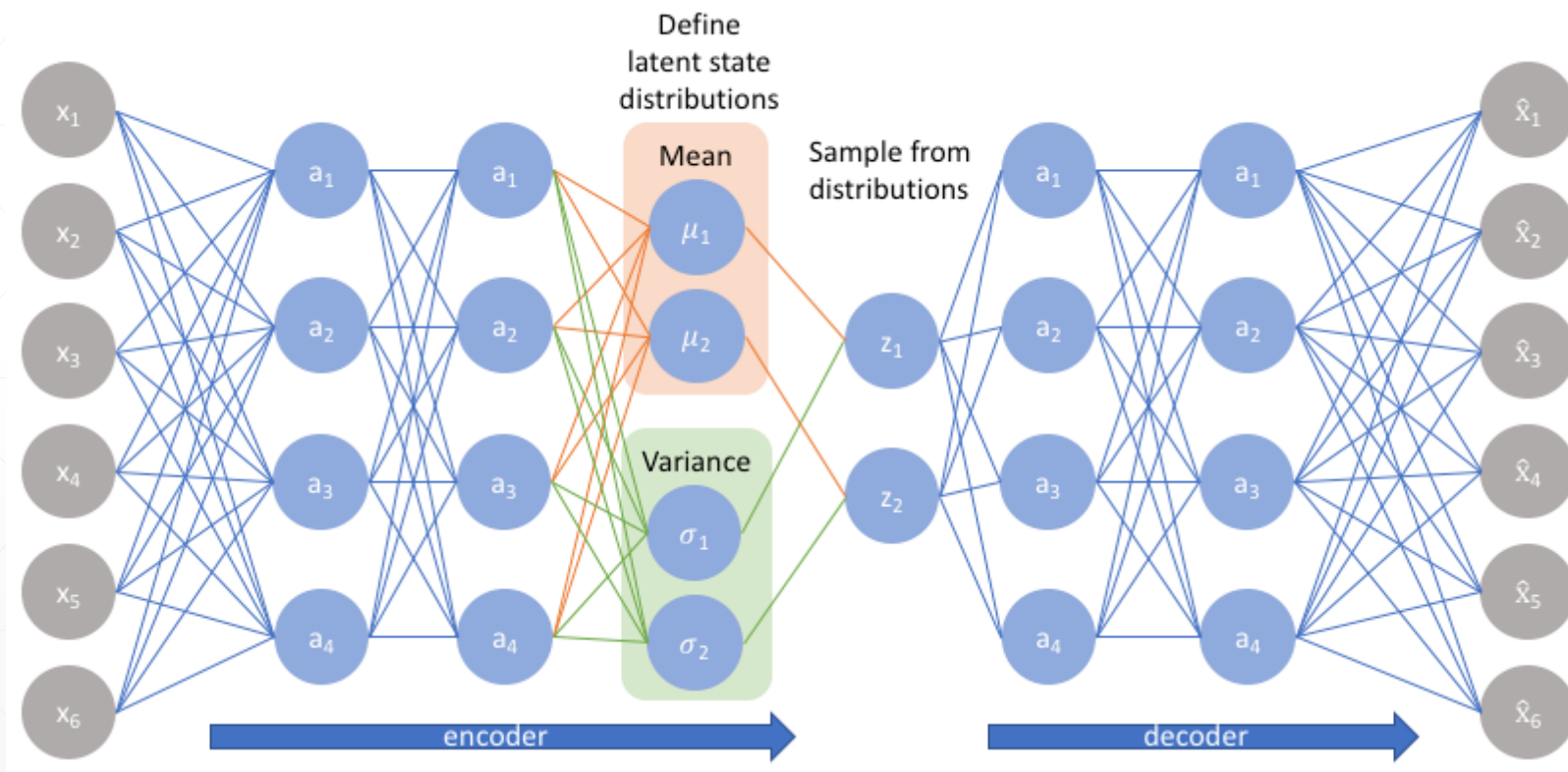
Penalizing KL divergence acts as a regularizing force



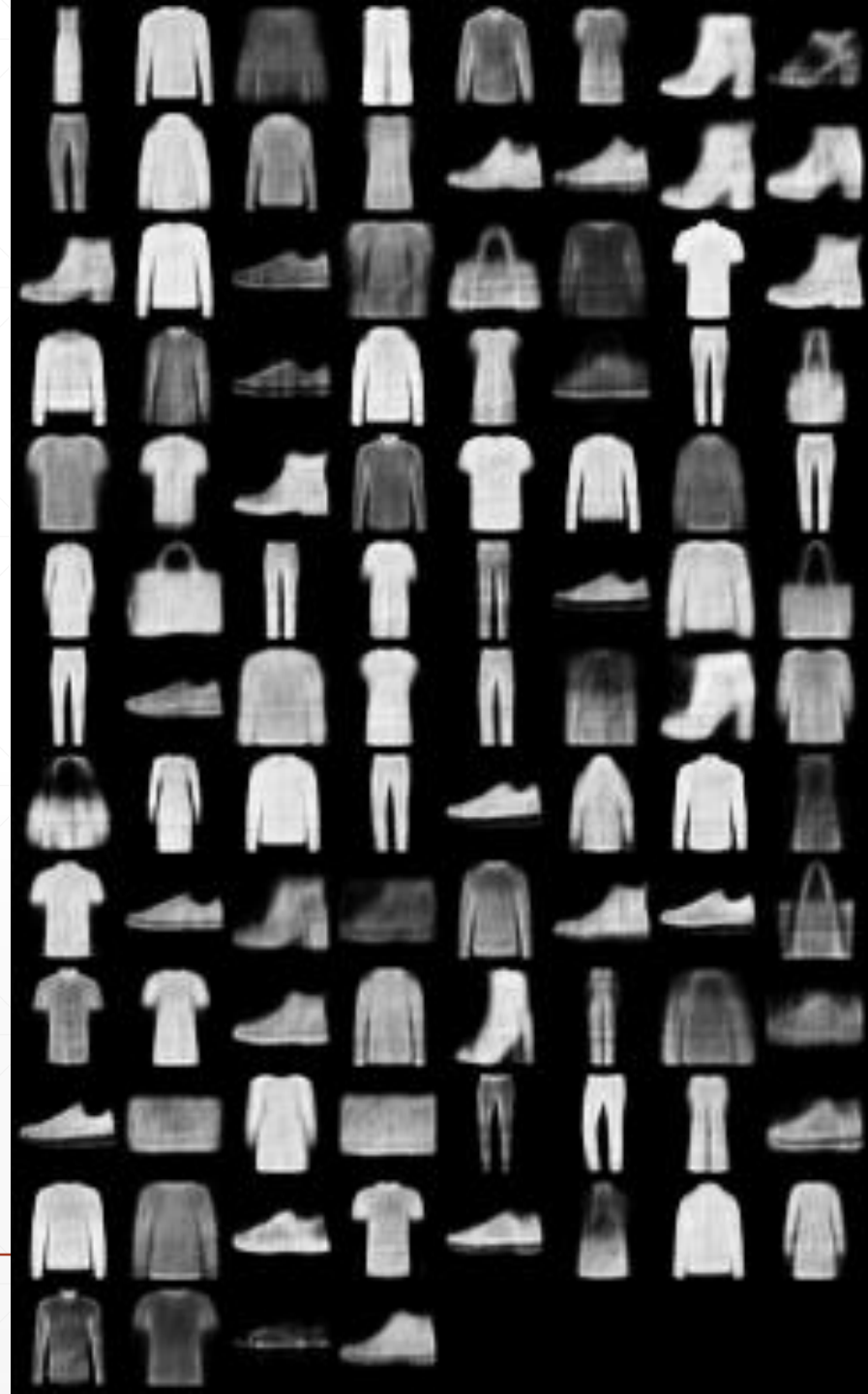
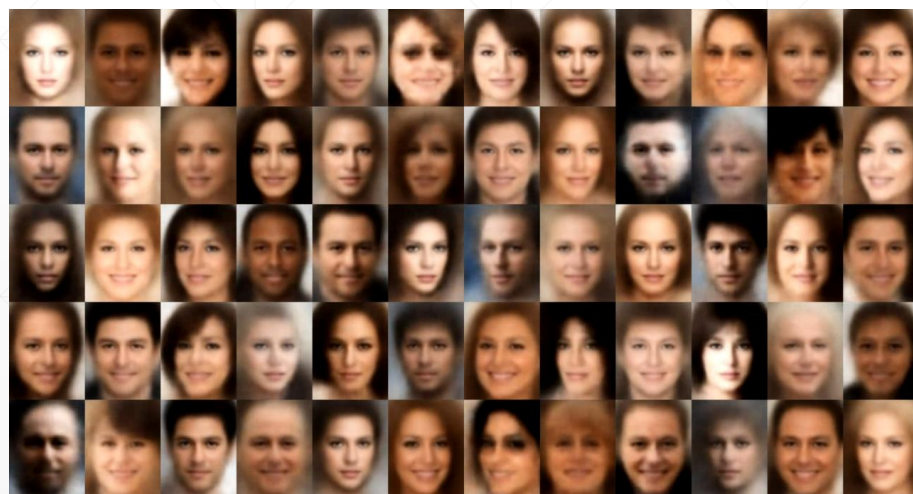
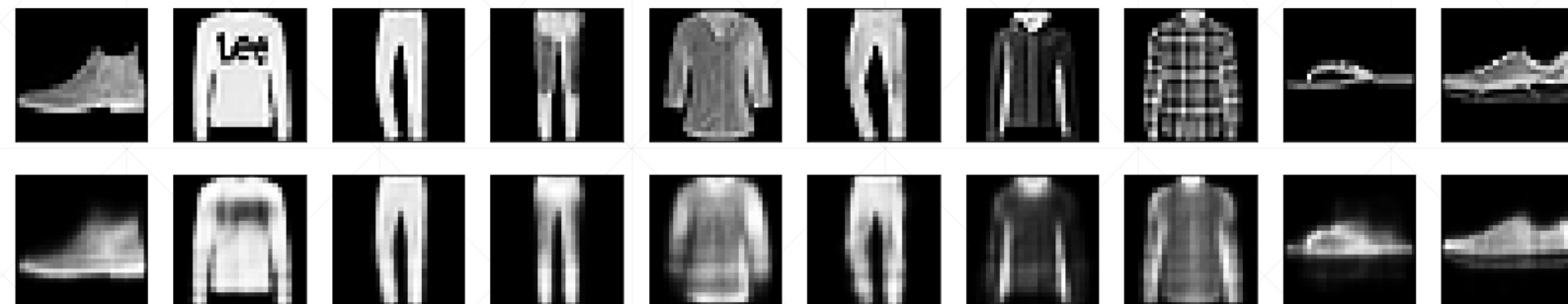
Attract distribution to have zero mean

Ensure sufficient variance to yield a smooth latent space

Too Complex!

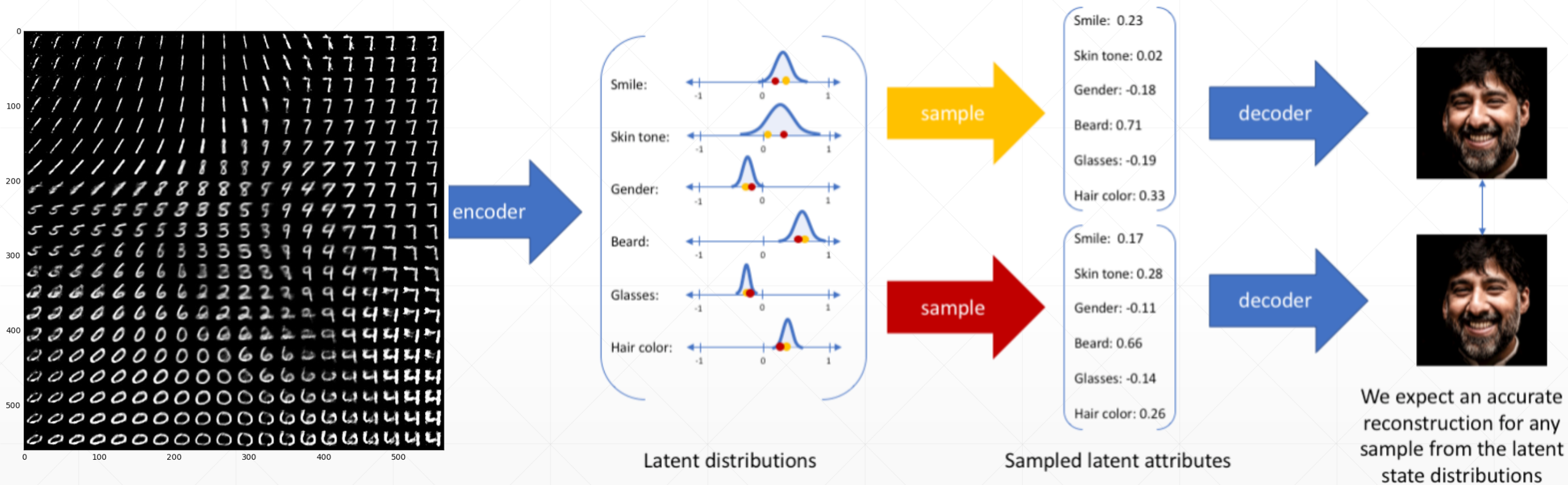


AE V.S. VAE



https://github.com/cryer/Variational_Auto-Encoder/blob/master/image/reconst_images_7.png

Generative model



VAE V.S. GAN



下一课时

Thank You.
