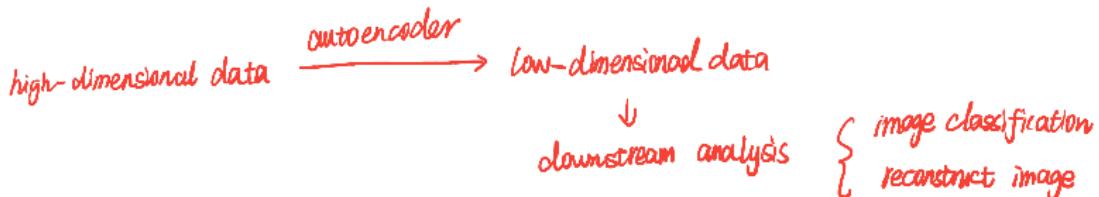


S&DS 365 / 565  
**Data Mining and Machine Learning**

# **Autoencoders**

**Yale**

# Outline



- Autoencoders: Generalization of dictionary learning using deep neural networks
- Variational autoencoders 複合

# Minimal neural network: Recall

- Toy data: 3-class spirals
- First discussed a logistic regression model
- Then a simple 2-layer network, backprop calcs

sparse dictionary learning only has 2 layers { input layer  
{ coding layer

These types of networks are sometimes called *multilayer perceptrons*

variational autoencoder : stack layers up

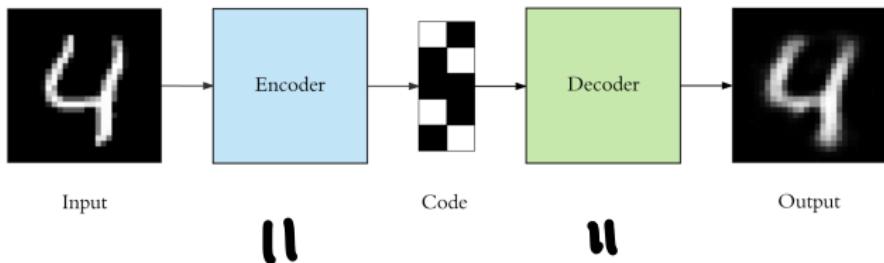
感知机

# Autoencoders

- Unsupervised learning methods
- Squeeze high dimensional data through a “bottleneck” of lower dimension
- Train to minimize reconstruction error

*similar to One of interpretation of PCA*

# Autoencoders



Similar to PCA

PCA

Projection  
Matrix  $A$

Transpose  
of Projection  
Matrix  $A^T$

$$x \rightarrow \boxed{A} \rightarrow Ax \rightarrow A^T(Ax) \rightarrow \hat{x}$$

(v) columns of  $A \in \mathbb{R}^{d \times k}$  are orthogonal

# Important aspects

- Unsupervised: No labels used, (discover useful features of input)
- Compression: Code reduces dimension of data
- Lossy: Input won't be reconstructed exactly
- Trained: The compression algorithm is learned for specific data

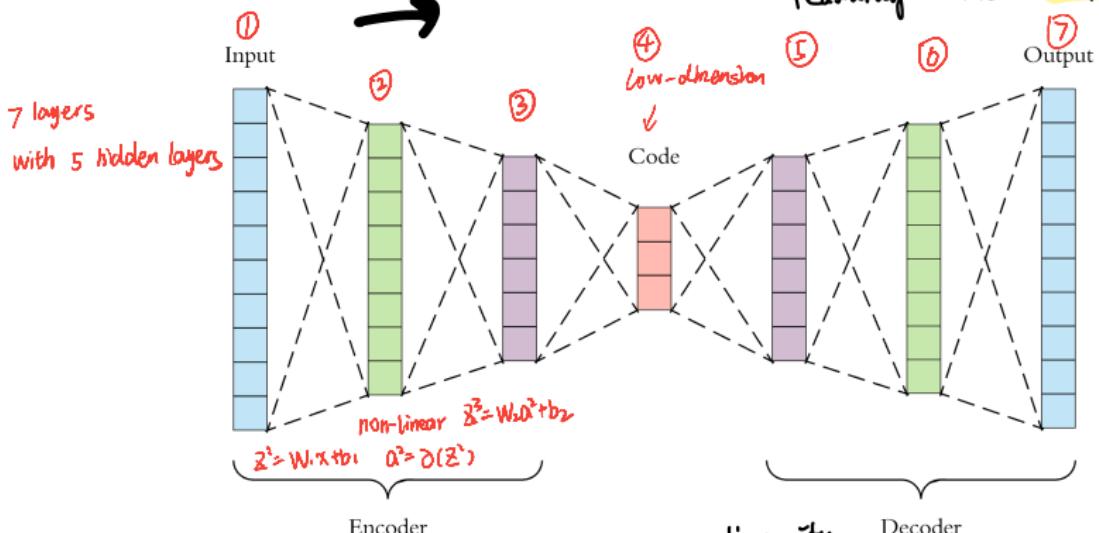
↓  
For PCA just use SVD, but can also train as an optimization problem

# Deep architecture

similar to Neural Network for supervised learning

**Feed** →

**Forward:** Instead of a label, basically "learning" the input itself!



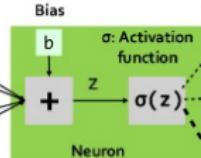
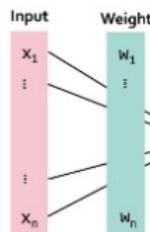
No hidden layer, all linear

$$\begin{aligned}
 \text{PCA} \quad x &\rightarrow W_1 x \in \mathbb{R}^k \rightarrow W_2 (W_1 x) \rightarrow \hat{x} \\
 x \in \mathbb{R}^d \quad W_1 \in \mathbb{R}^{k \times d} \quad W_2 \in \mathbb{R}^{d \times k} &\rightarrow \hat{x} \in \mathbb{R}^d
 \end{aligned}$$

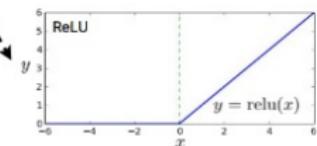
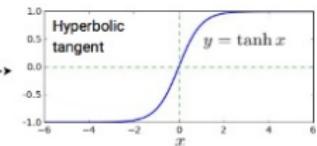
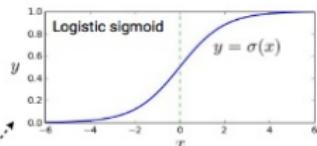
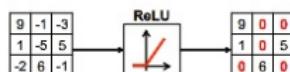
# Activation functions

## Rectified Linear Unit (Relu)

positive part of a number

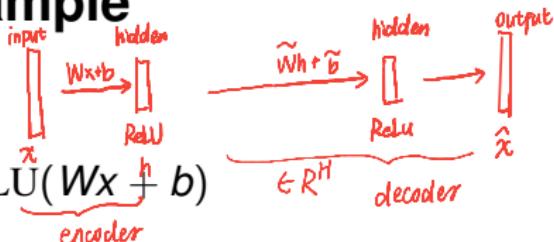


element wise



# Simple autoencoder example

Encoder network of the form



just one hidden layer

$$h = \text{ReLU}(Wx + b)$$

encoder

where  $W \in \mathbb{R}^{H \times D}$  and  $b \in \mathbb{R}^H$ , decoder network is

$$\hat{x} = \text{ReLU}(\tilde{W}h + \tilde{b}) \in \mathbb{R}^D$$

where  $\tilde{W} \in \mathbb{R}^{D \times H}$  and  $\tilde{b} \in \mathbb{R}^D$ .

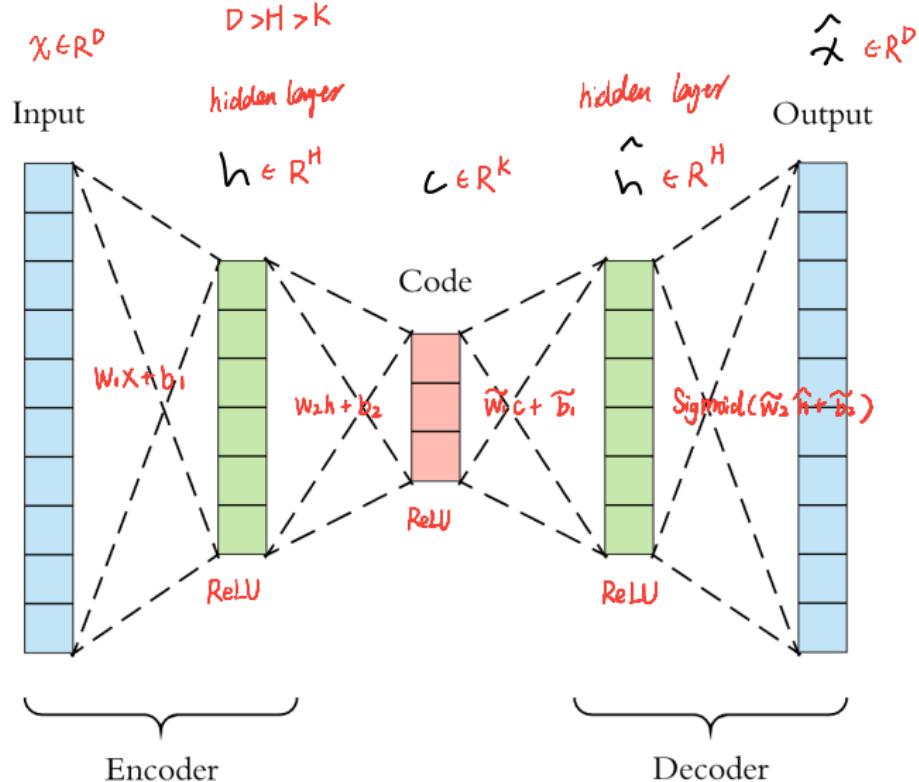
Use back-propagation to train this find  $W, b, \tilde{W}, \tilde{b}$

Objective function:  $\text{Loss} = \frac{1}{n} \sum_{i=1}^n \|x_i - \hat{x}_i\|_2^2$  reconstruction error =  $l_2$ -norm squared

$$\mathcal{L} = \frac{1}{n} \sum_{i=1}^n \|x_i - \text{ReLU}(\tilde{W} \text{ReLU}(Wx_i + b) + \tilde{b})\|_2^2 + \lambda \|h_i\|_1$$

$\text{ReLU}(x) = \max(0, x)$ , applied element-wise component-wise.

# Simple 2-layer architecture



# Simple example – code of dimension $K$

Encoder network of the form

$$h = \text{ReLU}(W_1x + b_1) \in \mathbb{R}^H$$

$$c = \text{ReLU}(W_2h + b_2) \in \mathbb{R}^K$$

where  $W_1 \in \mathbb{R}^{H \times D}$  and  $b_1 \in \mathbb{R}^H$ , and  $W_2 \in \mathbb{R}^{K \times H}$  and  $b_2 \in \mathbb{R}^K$

Decoder network of the form

$$\hat{h} = \text{ReLU}(\tilde{W}_1c + \tilde{b}_1)$$

$$\hat{x} = \text{Sigmoid}(\tilde{W}_2\hat{h} + \tilde{b}_2)$$

where  $\tilde{W}_1 \in \mathbb{R}^{H \times K}$  and  $\tilde{b}_1 \in \mathbb{R}^H$ , and  $\tilde{W}_2 \in \mathbb{R}^{D \times H}$  and  $\tilde{b}_2 \in \mathbb{R}^D$

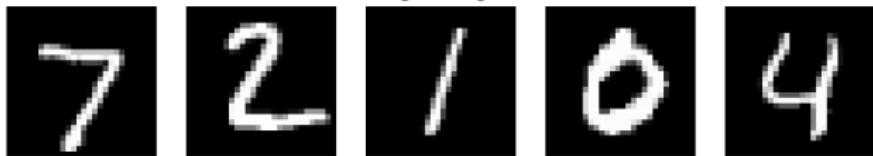
*can use different loss that suit for specific problem*

Objective function: **binary cross-entropy**      for 0-1 Loss

$$\mathcal{L} = \frac{1}{n} \sum_{i=1}^n (x_i \log \hat{x}_i + (1 - x_i) \log(1 - \hat{x}_i))$$

# MNIST data

Original Images



Dimension

$$28 \times 28 = 784 = D$$

# Adam optimizer

Adaptive gradient descent with momentum

Have package

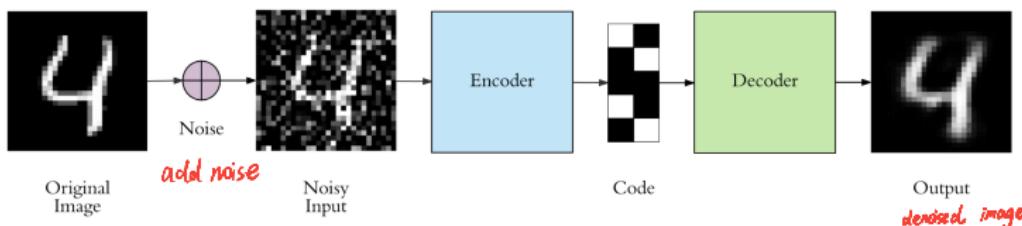
- Variant of stochastic gradient descent where separate learning rate (step size) is maintained for each network weight (parameter) eg  $w^1, w^2, b^1, b^2 \dots$

- Each step size adapted as learning progresses based on second moments of the derivatives

Like Hessian

# Variant: Denoising autoencoder

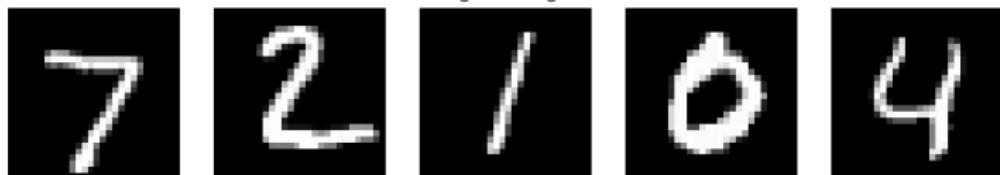
①



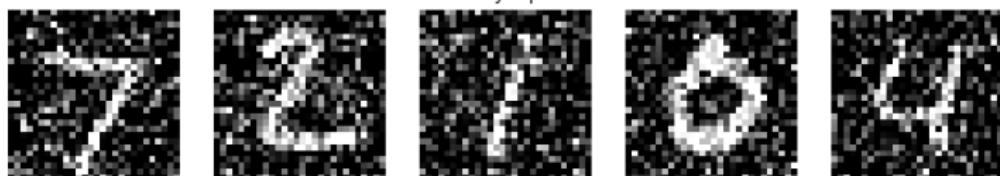
- Feed in noisy data
- Train to match to denoised data  
minimize reconstruction error between output and original input (not noisy input)

# Example result on MNIST

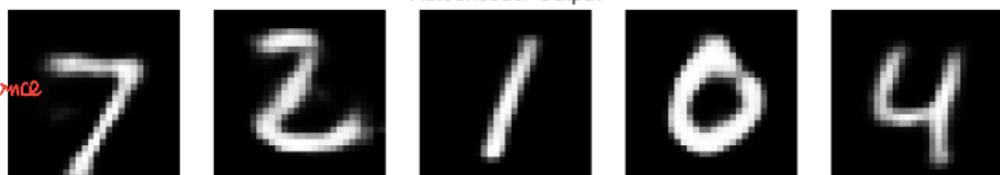
Original Images



Noisy Input



Autoencoder Output



# Variant: Sparse autoencoder

②

Loss function:

$$L = \text{L} + \lambda \| h \|_1$$

e.g. add a  $\ell_1$  norm penalty  
want "code" layer to be sparse

$\begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} \leftarrow \text{l-sparse code}$

- Add **penalty** to encourage code sparsity
- Want to find interesting structure
- Generalizes sparse **dictionary learning** from before

# Autoencoders

Recap

- Unsupervised learning methods
- Squeeze high-dimensional data through a “bottleneck”
- Train to minimize reconstruction error

# Variational autoencoders

- Variational autoencoders offer a more statistical view
- Based on variational inference *(faint)* *More of Bayesian*
- The decoder is a generative model with a latent variable  $\mathbf{z}$
- The encoder approximates the posterior distribution *for the latent variable to find the latent variable*

Recap of Bayes:

Latent Random vector  $\mathbf{z}$

$\mathbf{x} | \mathbf{z} \sim P(\cdot | \mathbf{z})$   $\mathbf{x}$  is observed

$P(\mathbf{z} | \mathbf{x}) \rightarrow$  inference

# Variational autoencoders

Start with a generative model

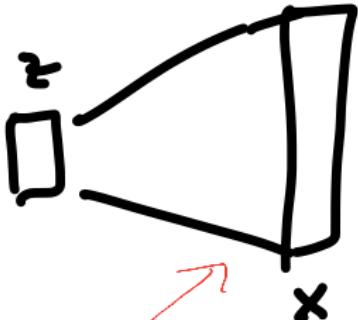
Suppose a latent variable is normal  $Z \sim N(0, I_k)$

$$X | Z = z \sim N(G(z), \gamma^2 I_D)$$

$G(z)$  is the *generator network* or *decoder*

As a simple example:    One hidden layer

$$G(z) = A_2 \text{ReLU}(A_1 z + b_1) + b_2$$



# Transforming Gaussians

Crazy things can happen when you apply activation functions to Gaussians

# Transformed Gaussian

仿射 (缩放变换  $A$  + 平移变换  $b$ )

Let  $Z \sim N(0, I)$  and an affine mapping

hidden layer

$$h = AZ + b$$

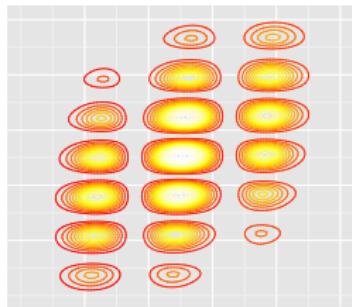
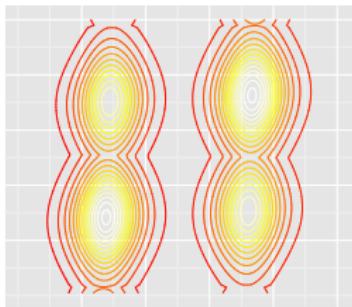
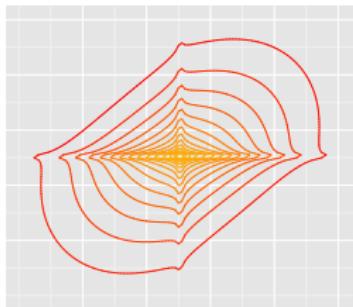
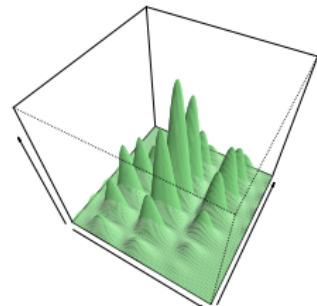
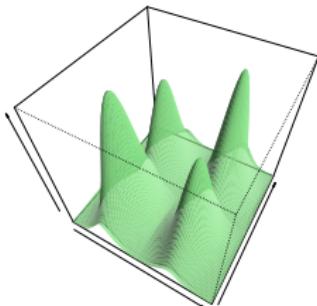
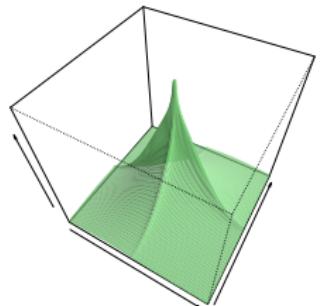
so that  $h \sim N(b, AA^T)$ .

non-linear 单调变换

Then, apply a monotonic transformation

$$X = g(h)$$

# Example densities: $k = 2 \rightarrow d = 2$



# Posterior inference

Back to training our generative network:

$$Z \sim N(0, I_K)$$

$$X | Z = z \sim N(G(z), \gamma^2 I_D)$$

We want the posterior distribution  $p(z | x)$ .

This is generally hard to compute, because  $G(\cdot)$  is nonlinear

Approach: Use variational inference

powerful method for solving all Bayesian inference

# Mean field variational inference

One method ↑

In mean field variational inference we take

a proposal distribution  $q(z | x) = N(\mu(x), \text{diag}(\sigma^2(x)))$

where now  $\mu_j(x)$  and  $\sigma_j^2(x)$  are the *variational parameters* for  $j = 1 \dots, K$ .  $z \in \mathbb{R}^K$   $\mu$  and  $\sigma^2$  are independent

This is a very popular approach for topic models!

Want APPROXIMATE

the posterior distribution  
complex  $p(z|x)$  with  $q(z|x)!$

⇒ final manifold<sup>3</sup>  
that minimize discrepancy  
between  $p(z|x)$  and  $q(z|x)$

# Using neural networks

Build a “recognition” or **encoder neural network** that outputs the mean and variance.  $\sigma^2$   $\mu$

For example: *analogous to encoder network of Adam autoencoder  
encoding  $x$  into mean param  $\mu(x)$*

$$\mu(x) = B_2 \text{ReLU}(B_1 x + d_1) + d_2$$

and similarly for  $\log \sigma^2(x)$ . ( $\sigma^2 = e^{\log \sigma^2}$ )

# Using neural networks

For more details  
David Blei tutorial.

Loss function involves  $\mathbb{E}_q(\log p(x|Z))$ . Approximate by sampling  $z_s$

comes trying to minimize KL divergence between 2 distributions  $q(x|z)$  and  $p(x|z)$  true

$$\text{argmin}_{\mu, \sigma^2} \text{KL}(q || p) = \underset{q}{\text{argmin}} \mathbb{E}_q(\log p(x|Z)) \approx \frac{1}{N} \sum_{s=1}^N \log p(x|Z_s)$$

key: expectation over  $q$  distribution for an observed  $x$

has order. not symmetric: with one takes expectation with respect to  $Z$  is only random variable in equation,  $x$  is not random, it's observed

But the parameters of the recognition network have disappeared!

$\mu, \sigma$

"Reparameterization trick": Re-express the samples by  $z_s$

$$Z_i = \mu(x) + \sigma(x)\epsilon_i \text{ where } \epsilon_i \sim N(0, I_k).$$

# Simple example

Suppose  $x | z \sim N(G(z), I)$  where generator network is

$$G(z) = \text{ReLU}(Az + b).$$

Negative log conditional probability (< NLL?)

Then  $-\log p(x | z)$  is  $\frac{1}{2} \text{l}_2\text{-norm squared}$  see in linear regression

$$\begin{aligned} & \exp(-\|x - G(z)\|_2^2) \\ & \text{take -log} \quad = \quad \frac{1}{2} \|x - \text{ReLU}(Az + b)\|_2^2 \\ & \Rightarrow \|x - G(z)\|_2^2 \quad \text{Then How to deal with } z? \end{aligned}$$

# Simple example

Next, suppose the approximate posterior is  $q(z|x) = N(\mu(x), I_K)$

remove  $\sigma^2$  in simple e.g  
means  $z = \mu(x) + \epsilon$

$\epsilon \sim N(0, I_K)$   
reparametrization trick:  
instead of randomly sample  $z$

where recognition network is  $\mu(x) = \text{ReLU}(Bx + d)$ . Then we randomly sample  $E$

For a fixed  $x$ , sample a bunch of  $Z_s$  according to distribution  $q$ :  $Z = \mu(x) + \epsilon$ ,  $\epsilon \sim N(0, I_K)$

$$-\mathbb{E}_q \log p(x|Z) \approx \frac{1}{N} \sum_{s=1}^N \frac{1}{2} \|x - \text{ReLU}(AZ_s + b)\|^2 = \text{ReLU}(Bx + d) + \epsilon$$

sometimes  $N = 1$

notation  $\stackrel{d}{=}$  means equal in distribution

$$\stackrel{d}{=} \frac{1}{N} \sum_{s=1}^N \frac{1}{2} \|x - \text{ReLU}(A(\text{ReLU}(Bx + d) + \epsilon_s) + b)\|^2$$

$$\hat{\mu} = \arg \min_{\mu} (-\mathbb{E}_{z \sim q_{\mu}} \log p(x|z))$$

train by gradient descent,  
compute  $\frac{\partial E}{\partial \mu}$ , but inside of  $E$  don't have  $\mu$  ↑ nice thing:  $\epsilon_s$  in independent of param  $\mu$

$$\begin{aligned} -\mathbb{E}_{z \sim q_{\mu}} \log(p(x|z)) &= \mathbb{E}_{z \sim q_{\mu}} \frac{1}{2} \|x - \text{ReLU}(Az + b)\|^2 \\ &\quad \text{use reparametrization trick, can compute } \frac{\partial \| \cdot \|_2^2}{\partial \mu}, \text{ then compute } \frac{\partial E}{\partial \mu} \text{ 整个式子} \\ &= \mathbb{E}_{\epsilon \sim N(0, I_K)} \frac{1}{2} \|x - \text{ReLU}(A(\mu(x) + \epsilon) + b)\|^2 \end{aligned}$$

$z$ 服从  $q$  分布,  $q$  取决于参数  $\mu$   
we want optimize  $\mu$

# Déjà vu all over again

似曾相识

We saw this before—almost.

Simple autoencoder

$$\|x - \text{ReLU}(A \text{ReLU}(Bx + d) + b)\|^2$$

Variational autoencoder      Sum over  $N \uparrow$  random sample of  $\epsilon_s$

$$\frac{1}{N} \sum_{s=1}^N \frac{1}{2} \|x - \text{ReLU}(A(\text{ReLU}(Bx + d) + \epsilon_s) + b)\|^2$$

① Perspective 2

Similar to simple autoencoder, can think as add noise  $\epsilon_s$  to code  
want reconstruction to be robust against adding noise

Perspective and thinking have changed—back to statistics!

② Perspective 2 : Bayesian

want code have latent variable interpretation: given  $x$ , the approximate posterior distribution of  $z$  is

$$\text{Gaussian: } q(z|x) = N(\mu(x), I_k)$$

$$z = \text{ReLU}(Bx + d) + \epsilon$$

# Generating data

Use generator  $G(z) = \text{ReLU}(Az + b)$  to generate random  $z$ s from Gaussian distribution  $z \sim (0, I_d)$

data by

random  $z$ s from Gaussian distribution  $z \sim (0, I_d)$

Generated data

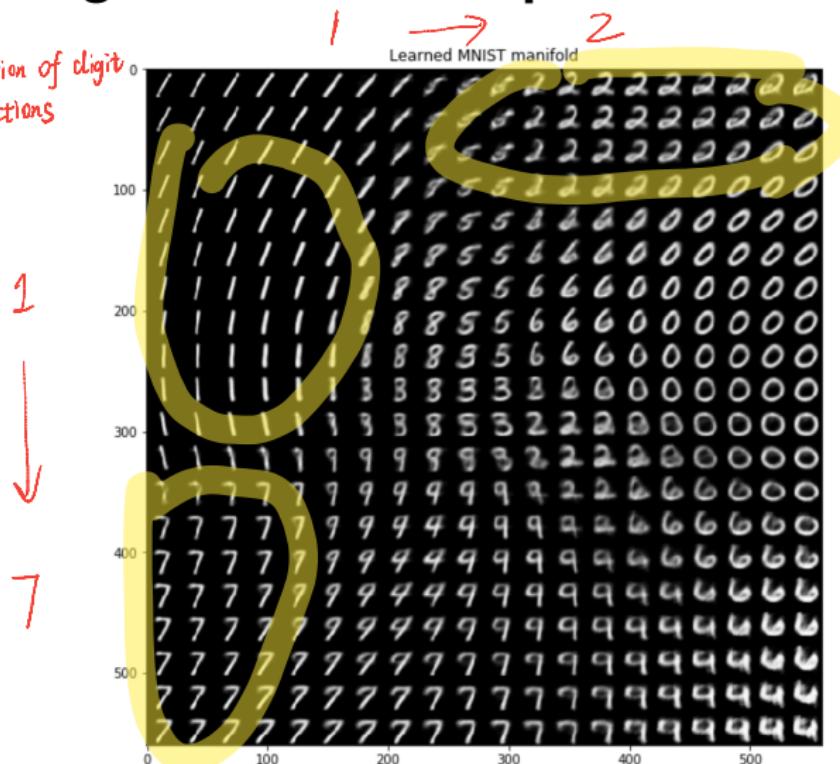
looks very like real numbers



# Visualizing a 2-dim latent space

2D code

can see transition of digit  
along 2 directions



# Interpolating in the latent space

插值

Topic Modelling

Transition

Topic Gradation	Generated Equation (Greedyly decoded)
Astrophysics (100%)	$G_{\text{eff}} = \frac{1}{2} \left( \frac{M_{\text{eff}}}{M_{\odot}} \right)^{-1} \left( \frac{M_{\text{eff}}}{M_{\odot}} \right)^{-1}$
$\vdots$	$G_{\text{eff}} = \frac{1}{2} \left( \frac{M_{\text{eff}}}{M_{\odot}} \right)^{-1}$
$\vdots$	$G_{\text{eff}} = \frac{1}{2} \left( \frac{1}{2} + \frac{1}{2} \right)$
50% — 50%	$G_s = \frac{1}{2} \left( 1 - \frac{1}{2} \right)$
$\vdots$	$G_s(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{x} + \mathbf{x}^T \mathbf{x}$
$\vdots$	$G_s(\mathcal{C}) = \mathcal{C}(\mathcal{C}) + \mathcal{C}(\mathcal{C}).$
Graph theory (100%)	$G_i = \{(x, y) \in \mathbb{R}^n : x_i = x_i\}$
Optimization (100%)	$L = \frac{1}{2} \ \mathbf{x} - \mathbf{x}\ _2^2 + \lambda \ \mathbf{x}\ _2^2 + \lambda \ \mathbf{x}\ _2^2 + \lambda \ \mathbf{x}\ _2^2$
$\vdots$	$L = \frac{1}{2} \ \mathbf{x} - \mathbf{x}\ _2^2 + \lambda \ \mathbf{x}\ _2^2 + \lambda \ \mathbf{x}\ _2^2$
$\vdots$	$L = \frac{1}{2} \sum_{i=1}^n \sum_{i=1}^n (x_i - x_i)^2 + \sum_{i=1}^n (x_i - x_i)^2$
50% — 50%	$L = \frac{1}{2} \sum_{i=1}^n \sum_{i=1}^n (x_i - x_i)^2 + \sum_{i=1}^n x_i^2$
$\vdots$	$L = \frac{1}{N} \sum_{i=1}^N \sum_{i=1}^N (x_i - x_i)^2$
$\vdots$	$L = \frac{1}{N} \sum_{i=1}^N \mathbb{E}[\mathbf{x}_i^T \mathbf{x}_i],$
Statistics (100%)	$L = \frac{1}{N} \sum_{i=1}^N \mathbb{E}[\mathbf{x}_i^T \mathbf{x}_i]$

Generator function

Loss function