

@author:GauthierLi

@date: 2021/12/20

Supervised Learning

Unsupervised Learning

Auto-Encoder

How to train?

PCA V.S. Auto-Encoders

Other Auto-Encoder

Adversarial AutoEncoders

Another Approach -- Variational Auto-Encoder

Maximize Likelihood

Minimize KL Divergence

sample process is not differential

Result

Supervised Learning

- classification
- regression

but real world exist numerous unlabeled datas

Unsupervised Learning

there exist three usual types of machine learning

- 'pure' reinforcement learning

interaction(交互) with the environment

- need a few bits for some samples

- Supervised learning

- 10 - 10,000 bits per sample

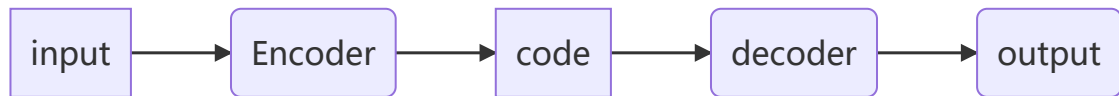
- Unsupervised/predictive learning

- millions bits per sample

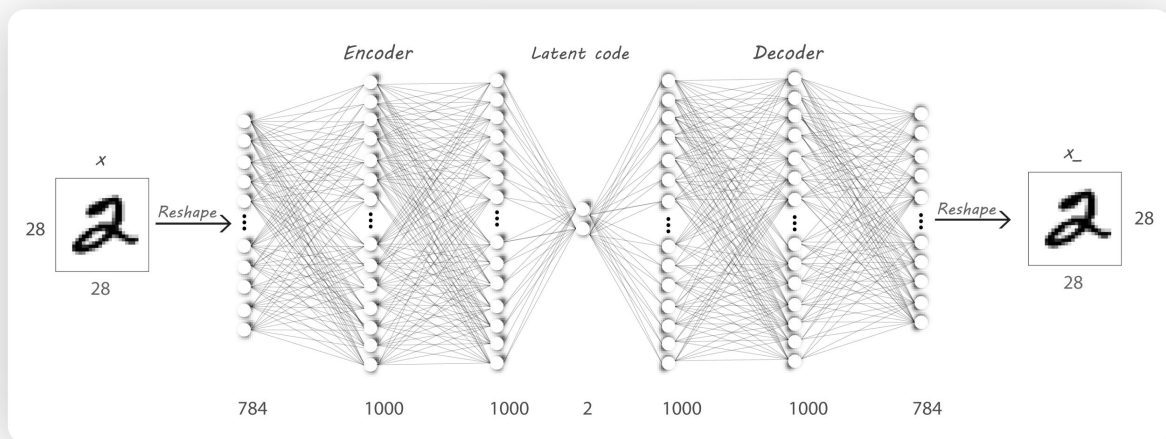
why we need unsupervised learning,

- dimension reduction
- Processing: Huge dimension, say $224 * 224$, is hard to process(将 $224 * 224$ 降为低维)
- Visualization: projector.tensorflow.org
- Compression, denosing, super-resolution(超分辨率)

Auto-Encoder



the aim of auto encoder is reconstruct(重建) the output same as the input via the encoder process and decoder process. Encoder process encode the input into a low dimension, and the decoder can be used for reform the result from the coded feature.



it could use to increase or decrease the feature by the **latent code** (并)

How to train?

- loss function for binary inputs(二进制文件, 如 Mnist 数据集)
 - cross-entropy error function (reconstruct loss) $f(x) \equiv \hat{x}$

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

- loss function for real-valued inputs
 - sum of square differences (reconstruct loss)
 - we used a linear **activation function** at the output

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

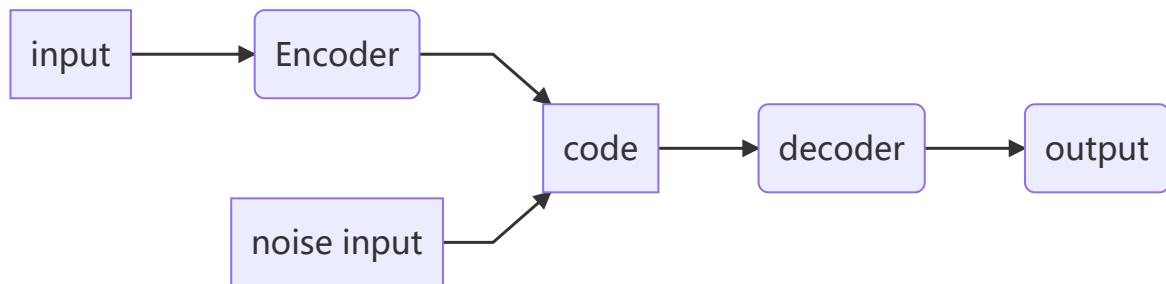
PCA V.S. Auto-Encoders

comparing to Auto-Encoders:

- PCA, which finds the directions of maximal variance in high-dimension 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. (丢失很多信息)

Other Auto-Encoder

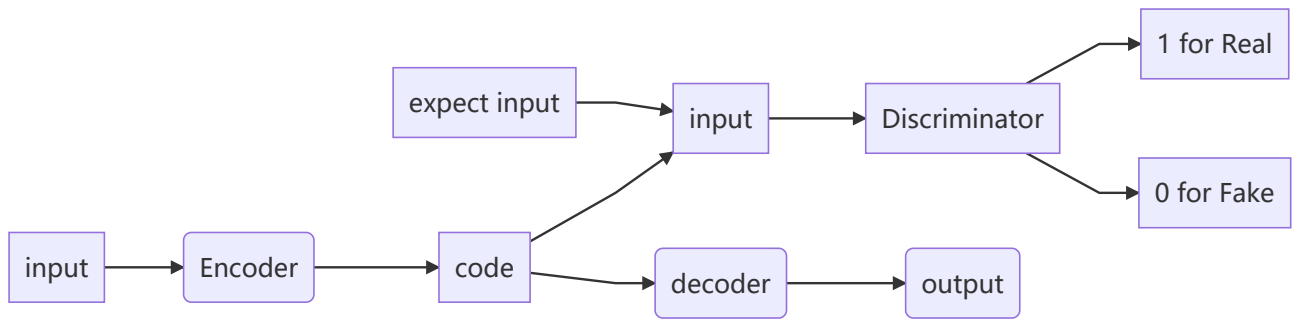
- Denosing AutoEncoder



- Dropout AutoEncoders
- Adversarial AutoEncoders

Adversarial AutoEncoders

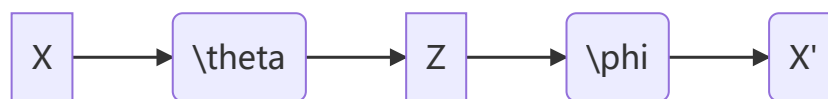
- distribution of hidden code



make the code submit to a specific distribution.

Another Approach -- Variational Auto-Encoder

- Maximum likelihood similarity



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

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

for the first part, to reduce the reconstruction loss of auto encoders, KL use to describe the difference between two distributions, means to make the feature similar to the specific distribution.

Maximize Likelihood

realise the $E_{Z \sim q_\theta}(\log[p_\phi(x_i|z)])$ part by:

- loss function for binary inputs(二进制文件, 如 Mnist 数据集)
 - cross-entropy error function (reconstruct loss) $f(x) \equiv \hat{x}$

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

- loss function for real-valued inputs
 - sum of square differences (reconstruct loss)

- we used a linear **activation function** at the output

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

Minimize KL Divergence

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

to calculate the $KL(q, p)$, we assume that

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

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

then,

$$\begin{aligned} KL(p, q) &= - \int p(x) \log(q(x)) dx + \int p(x) \log(p(x)) dx \\ &= \frac{1}{2} \log(2\pi\sigma_2^2) + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2} (1 + \log(2\pi\sigma_1^2)) \\ &= \log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2} \end{aligned}$$

Though the loss function, what we get is not a determinate feature, rather a **distribution**, so when we make a reconstruction, we have to make a 'sample' process to get a determinate feature vector. So, sample process is not differential.

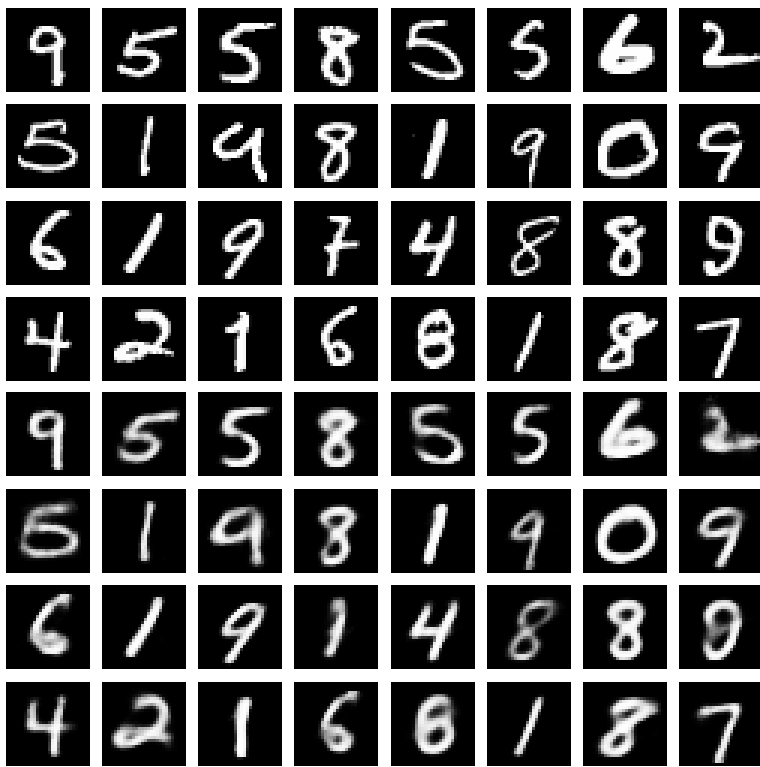
sample process is not differential

sample process is not differential, it hard to realize a back propagation. so, we used a **reparameterization trick**

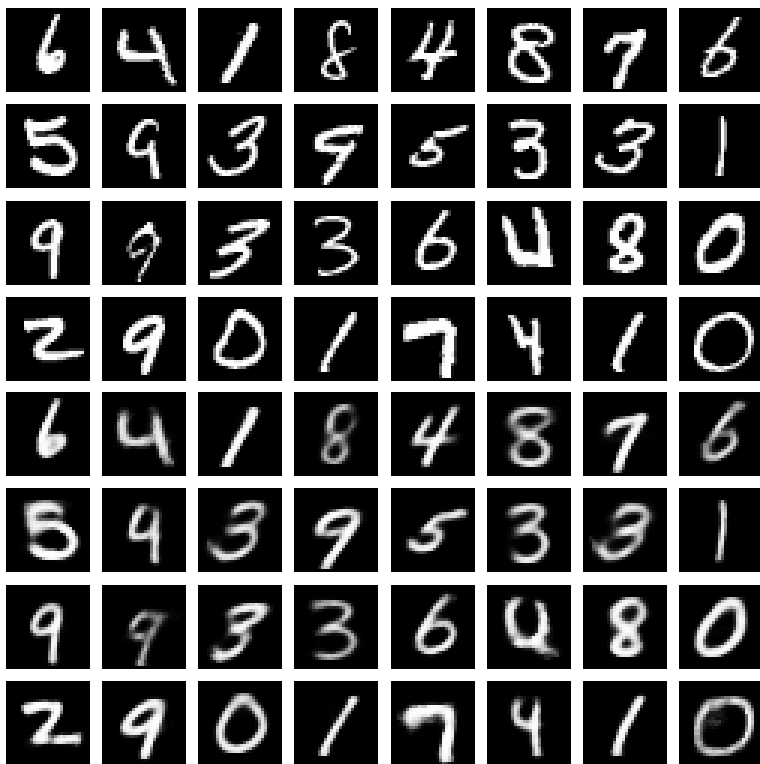
$$z \sim N(\mu, \sigma^2) \rightarrow z = \mu + \sigma \odot \epsilon, \quad \epsilon \sim N(0, 1)$$

\odot means dot product

Result



the two pictures above are the origin pic and after AE pic



the two pictures above are the origin pic and after VAE pic