

2024-11-02

Unsupervised Deep learning

* L11: Variational Autoencoders

* L12: Generative Adversarial Networks (GAN)

* L13: Diffusion models

Generative AI

Introduction

- So far, **supervised learning**

- Discriminative methods:

$$\{\mathbf{x}_i, y_i\}_{i=1}^N \rightarrow p(y|\mathbf{x})$$

- Generative methods:

$$\{\mathbf{x}_i, y_i\}_{i=1}^N \rightarrow P(y), p(\mathbf{x}|y)$$

class
size

class
characteristics

Gaussian Discriminant
analysis

linear regression
logistic / softmax regression
FNN
CNN
ViT

labelled data

PCA

- Next, **unsupervised learning**:

- Finite mixture models for clustering [Skipped]

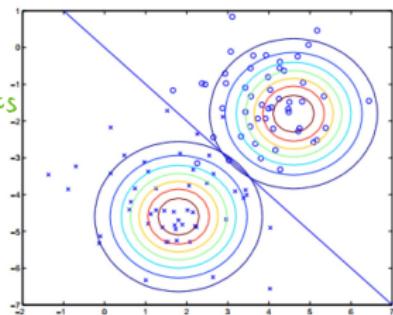
$$\{\mathbf{x}_i\}_{i=1}^N \rightarrow P(z), p(\mathbf{x}|z)$$

cluster
size

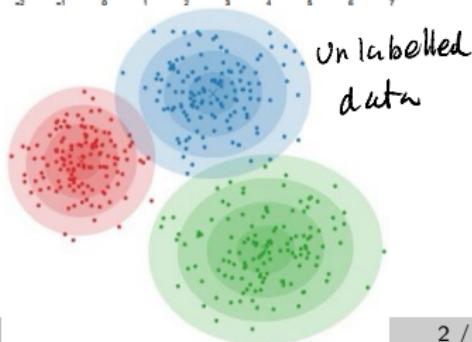
cluster
characteristics

z : latent variable

Intelligence
Math grade Physics grade History grade



Unlabelled
data



Gaussian Mixture
Models

Observed Variable

- Next, **unsupervised learning**:

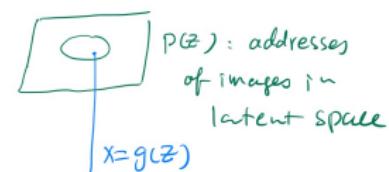
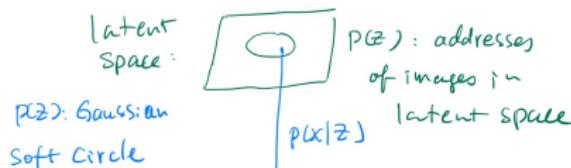
- *Variational autoencoder* for data generation and representation learning

$$\{x_i\}_{i=1}^N, p(z) \rightarrow p(x|z) \quad q(z|x) \text{ used in inference}$$

* Diffusion Model

- *Generative adversarial networks* for data generation

$$\{x_i\}_{i=1}^N, p(z) \rightarrow x = g(z)$$



VAE

Model used to

generate new data:

VAE : $z \sim p(z)$

GAN : $z \sim p(z)$

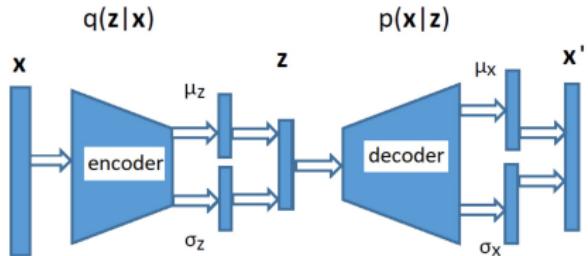
$x \sim p(x|z)$

$x = g(z)$

New
image

Lecture 11: Variational Autoencoder (VAE)

$$\Sigma = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}$$



* Want : $p(x|z)$

Decoder : $z \rightarrow x$

* Introduce encoder $q(z|x)$
in order to train
the decoder.

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

Real



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

Fake

Training
Target : $x' \approx x$

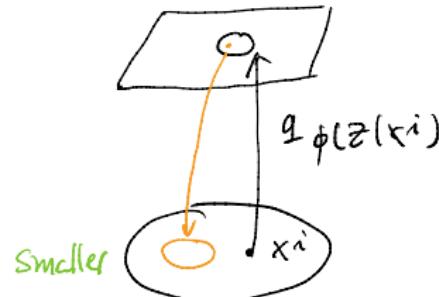
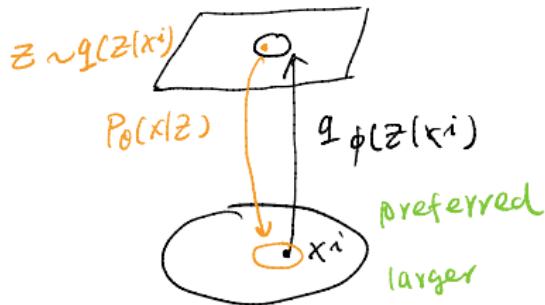
* General idea

Encoder : $x \rightarrow z$
real data latent vector

Decoder : $z \rightarrow x'$
latent vector fake data

The VAE Objective function

$$\{x^i\}_{i=1}^N, p(z)$$



$$L(x^i; \theta, \phi) = E_{z \sim q(z|x^i)} [\log P(x^i|z)] - KL(q(z|x^i) || p(z))$$

reconstruction error

regularization

$$\max_{\theta, \phi} L(\theta, \phi) = \frac{1}{N} \sum_{i=1}^N L(x^i; \theta, \phi)$$

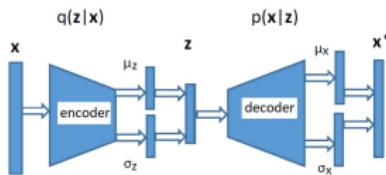


VAE

Given : $\{x^i\}_{i=1}^N \quad p(z)$

$$\max_{\theta, \phi} \frac{1}{N} \sum_{i=1}^N L(x^i, \theta, \phi)$$

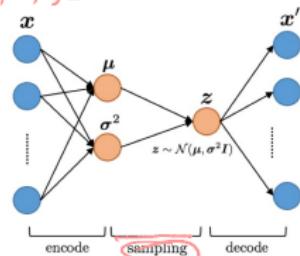
$$\mathcal{L}(x^{(i)}, \theta, \phi) = E_{z \sim q_\phi(z|x^{(i)})} [\log p_\theta(x^{(i)}|z)] - \mathcal{D}_{KL}[q_\phi(z|x^{(i)}) || p_\theta(z)]$$



The Need For Reparameterization

Forward

[$x, +, g$]



[$x, +, g$]

- The computation of the first term \mathcal{L}_1 of \mathcal{L} requires **sampling**:

$$\mathcal{L}_1 = E_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x}^{(i)})} [\log p_\theta(\mathbf{x}^{(i)}|\mathbf{z})] \approx \frac{1}{L} \sum_{l=1}^L \log p_\theta(\mathbf{x}^{(i)}|\mathbf{z}^{(i,l)})$$

where $\mathbf{z}^{(i,l)} \sim q_\phi(\mathbf{z}|\mathbf{x}^{(i)})$. $\underbrace{\text{depends on } \theta, \phi}_{\text{Not depend on } \phi}$

The Reparameterization Trick

$$\sigma = x^2$$

scalar x, z

$$q(z|x) = N(x^2+2, x^4)$$

$$z \sim q(z|x) \quad \frac{dz}{dx} = 0$$

$$z = x^2 + 2 + x^2 \varepsilon \quad \varepsilon \sim N(0, 1)$$

Gaussian distribution

$$E[z] = x^2 + 2$$

$$\sigma(z) = x^2$$

$$\frac{dz}{dx} = 2x + 2x\varepsilon$$

The Reparameterization Trick

$$\mathcal{L}_1 = E_{z \sim q_\phi(z|x^{(i)})} \left[\log p_\theta(x^{(i)}|z) \right]$$

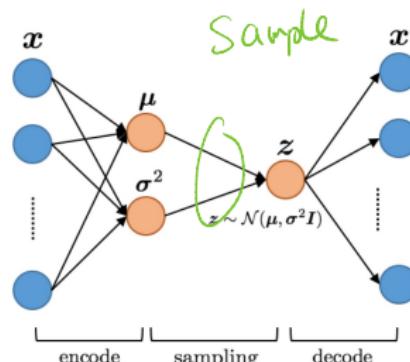
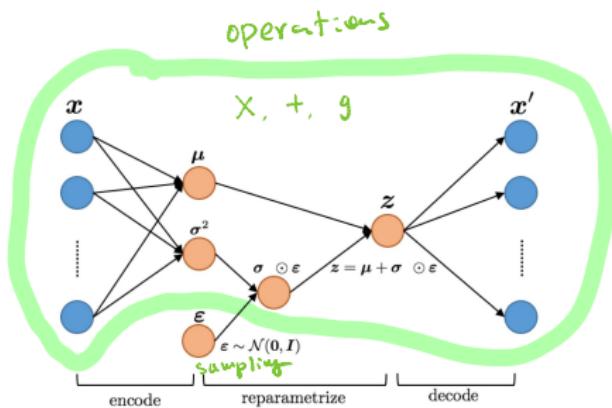
$z(x, \phi, \varepsilon)$

$$\approx \frac{1}{N} \sum_{l=1}^L \log p_\theta(x^l | z(x^l; \phi, \varepsilon^l))$$

$z = \mu_z(x, \phi) + \sigma_z(x, \phi) \odot \epsilon, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

$\varepsilon^1, \dots, \varepsilon^L \sim \mathcal{N}(0, I)$

Depends on θ, ϕ



The Second Term

Gaussian Distribution

$$\text{max } \mathcal{L}(\mathbf{x}^{(i)}, \theta, \phi) = E_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x}^{(i)})} \left[\log p_\theta(\mathbf{x}^{(i)}|\mathbf{z}) \right] - \mathcal{D}_{KL}[q_\phi(\mathbf{z}|\mathbf{x}^{(i)}) || p_\theta(\mathbf{z})]$$



$$\text{max } \mathcal{L}_2 = -\mathcal{D}_{KL}[q_\phi(\mathbf{z}|\mathbf{x}^{(i)}) || p_\theta(\mathbf{z})]$$

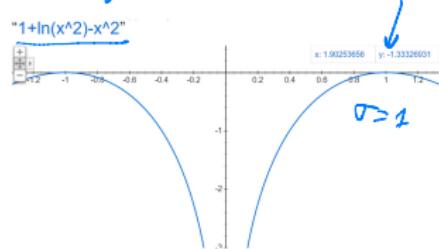
$$= \frac{1}{2} \sum_{j=1}^J \left(1 + \log((\sigma_j^{(i)})^2) - \underbrace{(\mu_j^{(i)})^2}_{- (\sigma_j^{(i)})^2} \right)$$

$$\begin{array}{c} q_\phi(\mathbf{z}|\mathbf{x}^{(i)}) \\ \hline \begin{bmatrix} \mu_1^{(i)} \\ \vdots \\ \mu_J^{(i)} \end{bmatrix} \xrightarrow{\text{max } L_2} \begin{bmatrix} \sigma \\ \vdots \\ \sigma \end{bmatrix} \end{array}$$

$$\text{max } \underline{1 + \log \sigma^2 - \sigma^2}$$

$\bar{\mu}$

$$\Sigma \begin{bmatrix} (\sigma_1^{(i)})^2 & & & \\ & \ddots & & 0 \\ & & 0 & (\sigma_J^{(i)})^2 \end{bmatrix} \xrightarrow{\text{max } L_2} \begin{bmatrix} 1 & & & \\ & \ddots & & 0 \\ & & 0 & 1 \end{bmatrix}$$

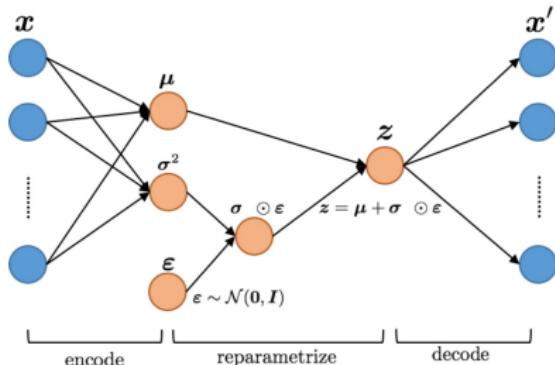


The Final Objective Function

- Putting together, this is the objective function that we maximize using gradient ascent

$$\mathcal{L} \approx \frac{1}{L} \sum_{l=1}^L \log p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}^{(i,l)}) + \frac{1}{2} \sum_{j=1}^J \left(1 + \log((\sigma_j^{(i)})^2) - (\mu_j^{(i)})^2 - (\sigma_j^{(i)})^2 \right)$$

where $\mathbf{z}^{(i,l)} = \mu_z(\mathbf{x}^{(i)}, \phi) + \sigma_z^2(\mathbf{x}^{(i)}, \phi) \odot \boldsymbol{\epsilon}^{(i)}$, and $\boldsymbol{\epsilon}^{(i)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

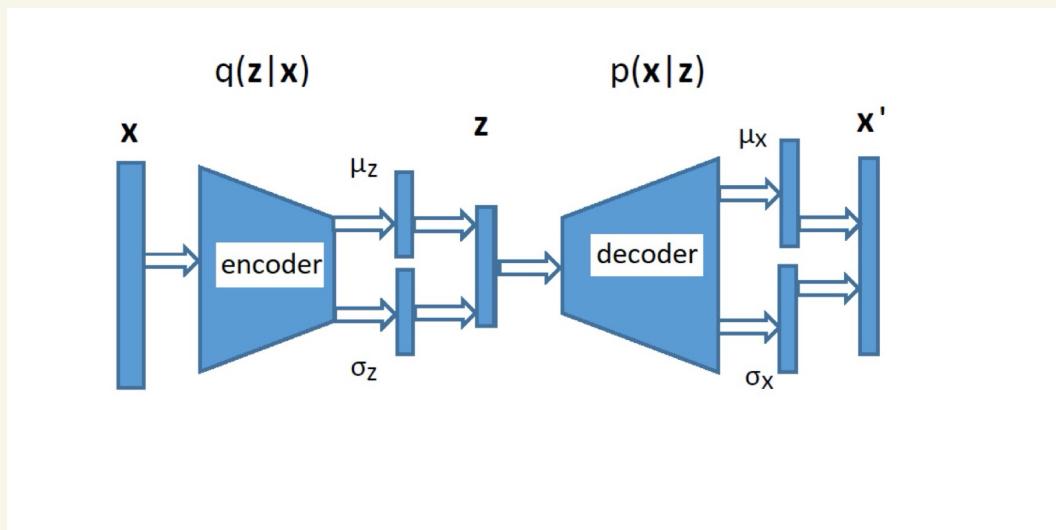


Tutorial 7 VAE

MNIST: images with binary pixels

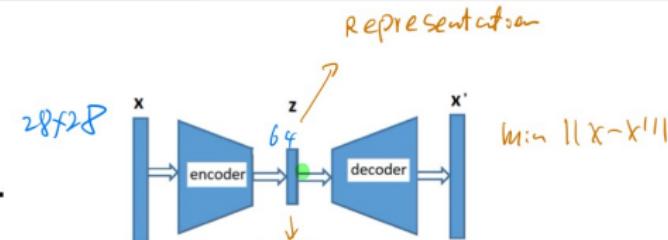
each pixel: 0, 1

$p(x|z)$: prob of each pixel being 1

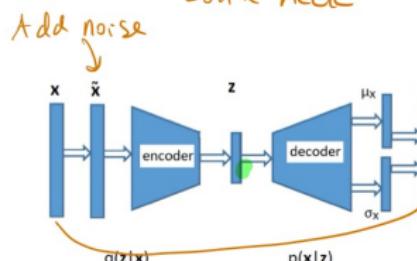


other Autoencoder models

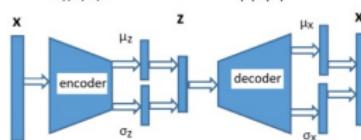
Auto encoder
AE



Denoising AE
DAG



VAE

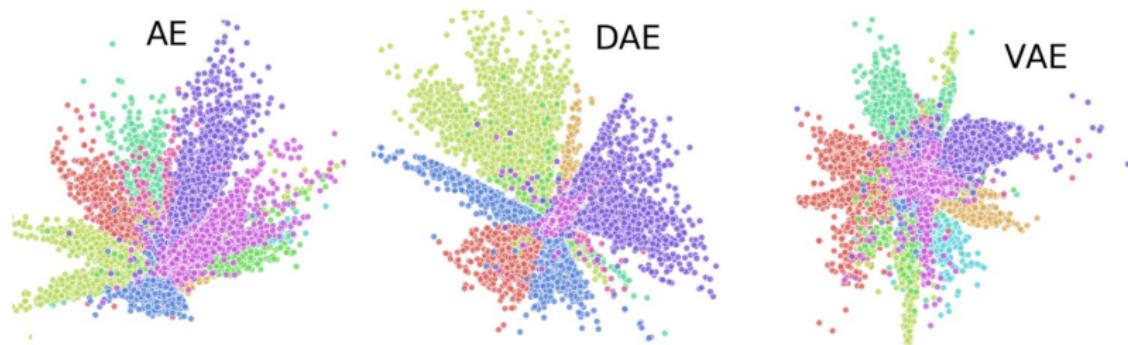


$$P(z): \text{Mo, I}$$

} Not regularize $p(z)$

$$\hat{x} \approx x'$$

Data Distribution in Latent Space (MNIST)



- VAE: Forces data into a normal distribution in the latent space.
- DAE: Preserves class separation better.



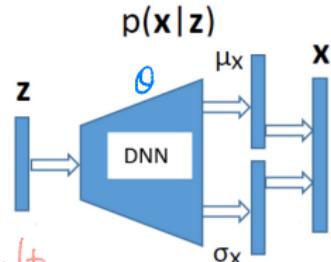
we know where the real images live in the latent space
so we can generate fake images that look like them

Theory of VAE (shared by Diffusion models)

$$\{x^i\}_{i=1}^N, p(z) \quad \Longrightarrow$$

$$p_{\theta}(x) = \int_{\text{learner}} p_{\theta}(x|z) \frac{p(z)}{\text{prior}} dz$$

$$\max \ell(\theta) = \frac{1}{N} \sum_{i=1}^N \log p_{\theta}(x^i) : \text{Difficult.}$$



$$E_{q(z|x)} [\log \frac{p(x|z)}{q(z|x)}]$$

$$= \int q(z|x) \log \frac{p(x|z)}{q(z|x)} dz$$

$$\leq \log \int q(z|x) \frac{p(x|z)}{q(z|x)} dz$$

Jensen's inequality

$$= \log \int p(x, z) dz$$

$$= \log p(x)$$

Therefore:

$$\log p(x) \geq E_{q(z|x)} [\log p(x|z)] - KL(q(z|x) || p(z))$$

$$E_{q(z|x)} [\log \frac{p(x|z)}{q(z|x)}]$$

$$= E_{q(z|x)} [\log \frac{p(z)p(x|z)}{q(z|x)}]$$

$$= E_{q(z|x)} [\log p(x|z)]$$

$$- E_{q(z|x)} [\log \frac{q(z|x)}{p(z)}]$$

$$\stackrel{\downarrow}{KL}(q(z|x) || p(z))$$

Variational lower bound,

↙ evidence lower bound
(ELBO)

