

VAE understanding

① Generative Model

Recover: P-PCA

a) Model Setting: $z \sim \text{Gaussian}(0, I_d)$
 $x|z \sim \text{Gaussian}(Wz, \sigma^2 I_D)$ $W \in \mathbb{R}^{D \times d}$
 $x \sim \text{Gaussian}(\mu_x, \Sigma_x)$

b) Learning Schema: $(\hat{W}, \hat{\sigma}^2) = \underset{W, \sigma^2}{\operatorname{argmax}} \log P(X|W, \sigma^2)$
(Maximize Likelihood Estimator)
 $= \underset{W, \sigma^2}{\operatorname{argmax}} \sum_{i=1}^N \log P(x_i | W, \sigma^2)$

c) Learning Technique:

(Expectation Maximization)

{ EM Algorithm (tricky)

{ Brute Force (Naive) from $x \sim N(\mu_x, \Sigma_x)$

EM Framework: Primal $\rightarrow (\hat{W}, \hat{\sigma}^2) = \underset{W, \sigma^2}{\operatorname{argmax}} \log P(X|W, \sigma^2)$

$$= \underset{W, \sigma^2}{\operatorname{argmax}} \log \mathbb{E}_Z [P(X, Z | W, \sigma^2)]$$

EM introduced $\rightarrow \theta^{(t+1)} = \underset{\theta}{\operatorname{argmax}} \mathbb{E}_{Z \sim p(\cdot|x, \theta^{(t)})} [\log P(X, Z | \theta)]$

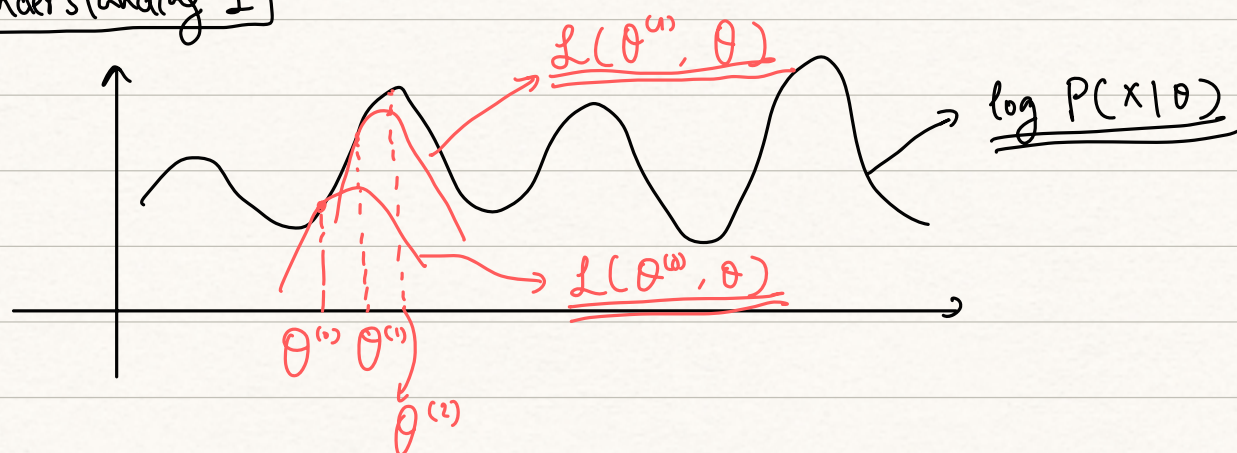
$$= \underset{\theta}{\operatorname{argmax}} \text{ELBO}(P_{z|x, \theta^{(t)}}, P_{(x,z)|\theta})$$

$\mathcal{L}(q, \theta)$

Actually: $\log P(X|\theta) = \text{ELBO}(q, P_{(x,z)|\theta}) + D_{KL}(q \parallel P_{z|x, \theta})$

Specially, $\log P(X|\theta^{(t)}) = \text{ELBO}(P_{z|x, \theta^{(t)}}, P_{(x,z)|\theta})$

Understanding 1



Understanding 2 GEM Framework

Bad Landscape

ideally, we want

$$\hat{\theta} = \operatorname{argmax}_{\theta} \log P(x|\theta)$$

↓ surrogate

consider $\log P(x|\theta) = L(q, \theta) + D_{KL}(q \| p_{z|x, \theta})$

Not bad

$$\begin{cases} q^{(k)} = \operatorname{argmax}_q L(q, \theta^{(k-1)}) = \operatorname{argmin}_q D_{KL}(q \| p_{z|x, \theta^{(k-1)}}) \\ \theta^{(k)} = \operatorname{argmax}_{\theta} L(q^{(k)}, \theta) \end{cases}$$

Why this is good?

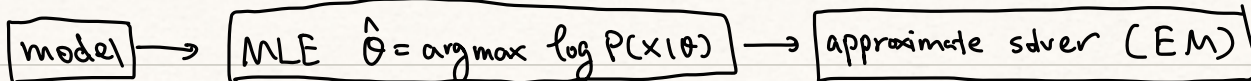
$$\log P(x|\theta^{(k)}) \geq \log P(x|\theta^{(k-1)})$$

Conclusion:

E-step: Given $\theta^{(k)}$, find $p(z|x, \theta^{(k)})$ (posterior)
or approximator of posterior

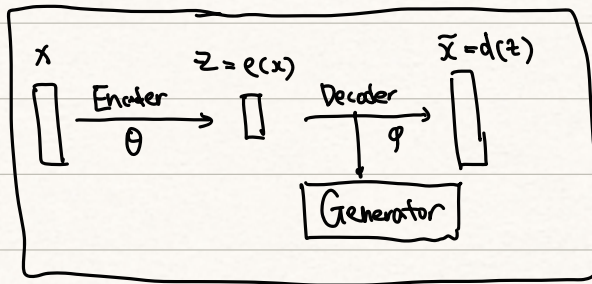
* M-step: Solve $\theta^{(k+1)} = \operatorname{argmax}_{\theta} \mathbb{E}_{z \sim p(z|x, \theta^{(k)})} [\log P(x, z|\theta)]$

Pipeline



② Auto-Encoder (AE model) \leftrightarrow NN discriminative model

not probabilistic model



Model Architecture

\Rightarrow Loss Function $(\hat{\theta}, \hat{\phi}) = \underset{\theta, \phi}{\operatorname{argmin}} \sum_{i=1}^N \|x_i - d(e(x_i; \theta), \phi)\|_2^2$

Reconstruction Loss

Limitation: cannot guarantee the 'regularity' of latent space

cannot

"generate" new meaningful inputs

our interest

What can be modified

{ use probabilistic model to enhance 'Regularity'
add Regularization term to Loss Function

\Downarrow
with respect to hidden (latent) space

③ Variational Auto-Encoder (VAE)

Outline: AE + Variational Inference + probabilistic formulation

\downarrow
explicit generative model

Recap: we want to do Dimensionality Reduction (without re-construction loss)

\nearrow AE's focus

\nearrow VAE's focus
preserve regularity of hidden (latent) space

\hookrightarrow it matters how the latent space is organized

VAE

1. Starting from Generative Model Framework

statistical fashion

model:

$$\begin{cases} z \sim N(0, I_d) \\ x|z \sim N(\underline{\mu}(z; \theta), \underline{\sigma}^2(z; \theta) I_D) \end{cases} \xrightarrow{\text{decoder (generator)}}$$

then $\underline{x} \sim N(\underline{\mu}_x, \underline{\Sigma}_x) \Rightarrow$ MLE estimator \rightarrow Naive

$$\begin{aligned} \hat{\theta} &= \underset{\theta}{\operatorname{argmax}} \log P(X|\theta) \quad \mathbb{E}_z [P(X|Z, \theta)] \\ &= \underset{\theta}{\operatorname{argmax}} \log \int_Z P(X|Z, \theta) \cdot P(Z|\theta) dz \\ &\left\{ \begin{aligned} &\approx \underset{\theta}{\operatorname{argmax}} \log \left\{ \frac{1}{m} \sum_{j=1}^m P(X|z_j, \theta) \right\} \quad \underline{z_1, \dots, z_m} \text{ samples} \\ &\text{EM framework} \end{aligned} \right. \end{aligned}$$

Limitation $\left\{ \begin{aligned} &a) \text{ computational cost for } \log \{ \underline{\Sigma}(\cdot) \} \\ &b) \frac{1}{m} \sum_j P(X|z_j, \theta) \approx P(X|\theta) \text{ might be inaccurate} \\ &\quad \downarrow \\ &\text{need } m \text{ large enough} \end{aligned} \right.$

Modification \rightarrow ELBO

$$\underbrace{\log P(X|\theta)}_{\text{likelihood}} = \underbrace{\ell(q, P_{(X,Z)}|\theta)}_{\text{ELBO}} + \underbrace{D_{KL}(q \parallel P_{Z|X, \theta})}_{\text{KL-divergence}}$$

$$\text{Here, } \ell(q, P_{(X,Z)}|\theta) = \int_Z q(z) \log \frac{P(X, Z|\theta)}{q(z)} dz$$

$$= \mathbb{E}_{z \sim q(\cdot)} \left[\log \frac{P(X, Z|\theta)}{q(z)} \right]$$

$$= \underline{\underline{-D_{KL}(q \parallel P_Z)}} + \underline{\underline{\mathbb{E}_{z \sim q(\cdot)} [\log P(X|Z, \theta)]}}$$

Recap, $\begin{cases} z \sim N(0, I_d) \\ x|z \sim N(f(z), g(z) \cdot I_d) \end{cases}$

$$\Rightarrow z|x \sim N(?, ?)$$

↓ Variational Inference

$$z|x \sim N(h(x), k(x) I_d)$$

$$(\hat{f}, \hat{g}) = \arg \max_{f, g} \log P(x|f, g)$$

$$\approx \arg \max_{\theta} \text{ELBO}(\tilde{P}_{z|x, \theta}, P_{(x, z)}|f, g)$$

introduce this setting

How to achieve $\tilde{P}_{z|x, \theta} \approx P_{z|x, \theta}$?

$$\tilde{P}_{z|x, \theta} = \arg \min_{q} D_{KL}(q \parallel P_{z|x, \theta})$$

$$= \arg \max_{q} \text{ELBO}(q, P_{(x, z)}|\theta)$$

$$\text{since } \log P(x|\theta) = \text{ELBO}(q, P_{(x, z)}|\theta) + D_{KL}(q \parallel P_{z|x, \theta}) \\ = \underline{\text{constant}}$$

$$\Rightarrow \{(\hat{f}, \hat{g}), (\hat{h}, \hat{k})\} = \arg \max_{h, k, f, g} \text{ELBO}(q(h, k), P_{(x, z)}(f, g))$$

↓
Decoder

↓
Encoder

$$= \arg \max_{h, k, f, g} \underline{-D_{KL}(q(h, k) \parallel P_z)} + \underline{E_{z \sim q(h, k)} [\log P(x|z, f, g)]}$$

(\hat{h}, \hat{k}) is determined by (\hat{f}, \hat{g})

↓
posterior approx.

generator

can approximately solve:

$$(\hat{f}, \hat{g}) = \arg \max_{f, g} \log P(x|f, g)$$

Here, $D_{KL}(q(h, k) \| p_z)$ $q(h, k) \sim \mathcal{N}(h(x), k(x) \cdot \mathbb{I}_d)$

$$\begin{aligned}
 &= \sum_{j=1}^d \int_z \frac{1}{\sqrt{2\pi k_j(x)}} \exp\left\{-\frac{1}{2k_j(x)}(z - h_j(x))^2\right\} \log \frac{\frac{1}{\sqrt{2\pi k_j(x)}} \exp\left\{-\frac{1}{2k_j(x)}(z - h_j(x))^2\right\}}{\frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}z^2\right\}} dz \\
 &= \sum_{j=1}^d \int_z \left(-\frac{(z - h_j(x))^2}{2k_j(x)^2} + \frac{1}{2}z^2 - \log(k_j(x))\right) \cdot \mathcal{N}(h_j(x), k_j(x)) dz \\
 &= \sum_{j=1}^d \left[-\frac{\mathbb{E}[(z - h_j(x))^2]}{2k_j(x)^2} + \frac{1}{2}\mathbb{E}[z^2] - \log(k_j(x))\right] \\
 &= \sum_{j=1}^d \left(-\frac{1}{2} + \frac{1}{2}(h_j(x)^2 + k_j(x)^2) - \log(k_j(x))\right)
 \end{aligned}$$

regularization term

$$\begin{aligned}
 &\mathbb{E}_{z \sim q(h, k)} [\log P(X|Z, f, g)] \quad \begin{cases} X|Z, f, g \sim \mathcal{N}(f(z), g(z)^2 \mathbb{I}_D) \\ z \sim \mathcal{N}(h(x), k(x)^2 \mathbb{I}_d) \end{cases} \\
 &= \mathbb{E}_{z \sim q(h, k)} \left[\log \prod_{j=1}^D P(x_j | z, f, g) \right] \\
 &= \mathbb{E}_{z \sim q(h, k)} \left[\sum_{j=1}^D \left\{ -\frac{1}{2} \log(2\pi \cdot g_j^2(z)) - \frac{1}{2} \frac{(x_j - f_j(z))^2}{g_j^2(z)} \right\} \right]
 \end{aligned}$$

Lastly, conclusion is

$$\begin{aligned}
 \{(\hat{f}, \hat{g}), (\hat{h}, \hat{k})\} &= \arg\max \frac{1}{n} \sum_{i=1}^n \log P(x_i | \theta) \\
 &\approx \arg\max \frac{1}{n} \sum_{i=1}^n \text{ELBO}_i(q(h, k), p(x, z)(f, g)) \\
 &= \arg\max \frac{1}{n} \sum_{i=1}^n \text{ELBO}(\mathcal{N}(h(x_i), k^2(x_i) \mathbb{I}_d), p(x_i, z_i)(f, g)) \\
 &= \arg\max \frac{1}{n} \sum_{i=1}^n \left\{ -D_{KL}(\mathcal{N}(h(x_i), k^2(x_i) \mathbb{I}_d), \mathcal{N}(0, \mathbb{I}_d)) \right. \\
 &\quad \left. + \mathbb{E}_{z \sim \mathcal{N}(h(x_i), k^2(x_i) \mathbb{I}_d)} [\log \mathcal{N}(x_i; f(z), g(z)^2 \mathbb{I}_D)] \right\} \\
 &= \arg\max -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^d \left(-\frac{1}{2} + \frac{1}{2}(h_j(x_i)^2 + k_j(x_i)^2) - \log(k_j(x_i))\right)
 \end{aligned}$$

Gaussian

$$+ \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^d \mathbb{E}_{z \sim N(h, k)} \left[-\frac{1}{2} \log(q_j^2(z)) - \frac{1}{2} \frac{(x_{ij} - f_j(z))^2}{q_j^2(z)} \right]$$

$$(\hat{h}, \hat{k}, \hat{f}) \stackrel{\boxed{g=1}}{=} \underset{f, h, k}{\operatorname{argmax}} - \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^d \left(-\frac{1}{2} + \frac{1}{2} (k_j(x_i)^2 + h_j(x_i)^2) - \log(k_j(x_i)) \right) + \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{z \sim N(h, k)} \left[-\frac{1}{2} \|x_i - f(z)\|^2 \right]$$

$$= \underset{f, h, k}{\operatorname{argmin}} \underbrace{\frac{1}{n} \sum_{i=1}^n \mathbb{E}_{z \sim N(h, k)} \left[\frac{1}{2} \|x_i - f(z)\|_2^2 \right]}_{\text{reconstruction error}}$$

$$+ \underbrace{\frac{1}{n} \sum_{i=1}^n \frac{1}{2} \left(\|h(x_i)\|_2^2 + \|k(x_i)\|_2^2 - \log(k(x_i)^2) - 1 \right)}_{\text{regularization on latent layer (space)}}$$

regularization on latent layer (space)

Highlight: Up to now, we haven't introduced "MN" yet!

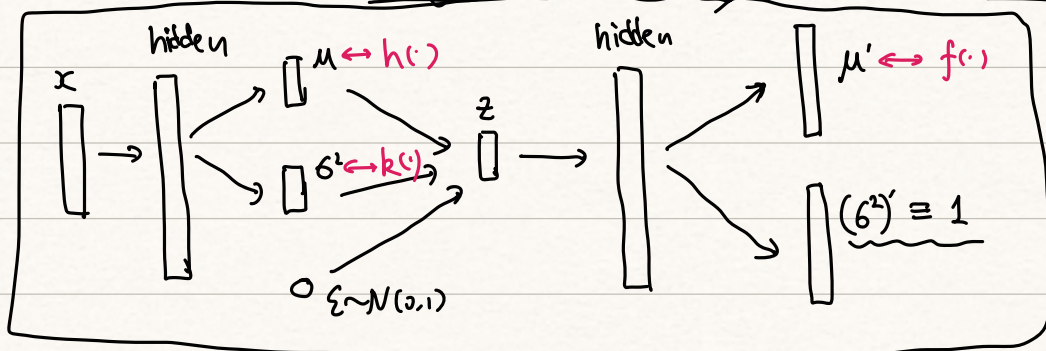
Idea is: After we achieve the simplified objective function

↳ which is achieved by

- ① MLE $P(x|\theta)$
- ② ELBO with approximate posterior
- ③ simplify the expression
- ④ achieve the form that,

$$\text{obj} = \text{reconstruction error} + \text{regularization term}$$

We can re-organize the formulation as NN

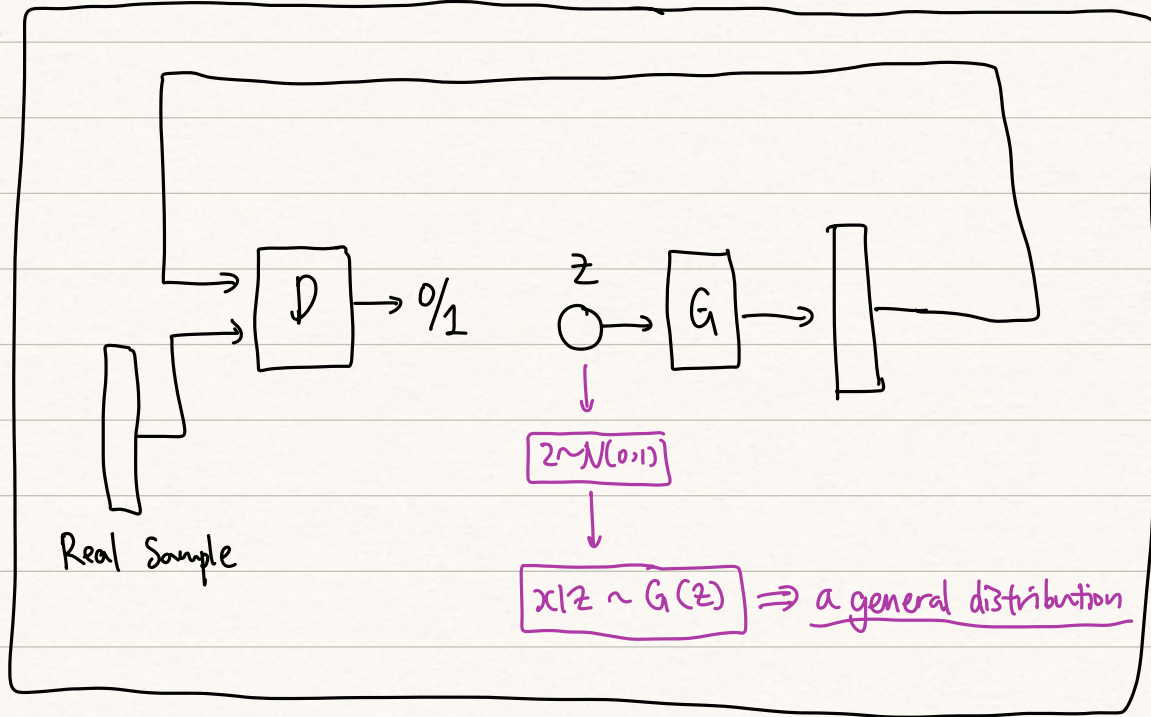


VAE (NN) architecture

encoder

decoder (generator)

→ Compare with GAN (VAE can do inference $P(z|x)$)



VAE $\begin{cases} z \sim \mathcal{N}(0,1) \\ x|z \sim \mathcal{N}(f(z), g(z)^2 I_D) \end{cases}$ (Strong Assumption)
↳ restrict the model power