

Variational Autoencoder (VAE) [2014]

idea: replace deterministic functions $f(x)$ and $g(z)$ with conditional distributions: encoder $P_E(z|x)$, decoder $P_D(x|z)$

⇒ more flexibility, and new training algorithms

• we want to minimize $KL[P^*(x) \parallel p(x)] \Leftrightarrow \mathbb{E}_{x \sim p^*(x)} [-\log p(x)]$ NLL loss

$$\log p(x) = \log p(x) \cdot \int P_E(z|x) dz = \int P_E(z|x) \log p(x) dz$$

$$= \int P_E(z|x) \log \left[p(x) \cdot \frac{P_E(z|x)}{P_E(z|x)} \cdot \frac{q(z)}{q(z)} \cdot \frac{P_D(x|z)}{P_D(x|z)} \right] dz \quad \begin{array}{l} q(z) \text{ desired} \\ \text{code distrib.} \end{array}$$

$$= \int P_E(z|x) \left[\log \frac{P_E(z|x)}{q(z) P_D(x|z) / p(x)} - \log \frac{P_E(z|x)}{q(z)} + \log P_D(x|z) \right] dz$$

$$\left[p(x) = \int P_D(x|z) dz = \int q(z) P_D(x|z) dz \quad \begin{array}{l} \text{decoder} \\ \text{posterior} \end{array} \quad P_D(z|x) = \frac{q(z) P_D(x|z)}{p(x)} \right]$$

$$\log p(x) = \underbrace{\text{KL}[P_E(z|x) \| P_D(z|x)]}_{\geq 0} - \underbrace{\text{KL}[P_E(z|x) \| q(z)]}_{\text{ELBO}(x)} + \underbrace{\mathbb{E}_{P_E(z|x)} [\log p_D(x|z)]}_{\text{ELBO}(x)}$$

• Observation 1: $\log p(x) \geq \text{ELBO}(x) \Rightarrow$ "evidence lower bound"
 $p(x)$

- maximizing lower bounds also tends to push-up $\log p(x) \Rightarrow$ maximize likelihood

• $\text{KL}[P_E(z|x) \| P_D(z|x)]$ "variational gap" if > 0 , measures the amount of inconsistency between encoder and decoder consistency. joint distributions should be equal

$$P_E(x, z) = p^*(x) \cdot P_E(z|x) \stackrel{!}{=} P_D(x, z) = q(z) \cdot P_D(x|z)$$

$$P_E(z|x) = \frac{q(z) P_D(x|z)}{p^*(x)} = \frac{q(z) P_D(x|z)}{P_D^*(z|x)}$$

- observation 2: $-\log p(x) + \text{KL}[p_E(z|x) \parallel p_D(z|x)] = -\text{ELBO}(x)$

\Rightarrow minimizing $-\text{ELBO}(x)$ has two effects. (i.e. minimizing RHS and LHS)

- minimize $-\log p(x) \Rightarrow$ high quality of generative distribution
(minimize NLL)

minimize $\text{KL}[p_E \parallel p_D] \Rightarrow$ improve self-consistency of model

- observation 3: terms in ELBO optimize trade-off between reconstruction quality ($\hat{x} \approx x$) and distribution fidelity ($p(x) \approx p^*(x)$)

Variational Autoencoders (contd.)

• encoder and decoder are not deterministic functions, but conditional distributions

enc. $p_E(z|x)$ dec. $p_D(x|z)$ + data $p^*(x)$ desired code distrib. $q(z)$

⇒ implies two versions of joint distribution of x and z

$$\text{enc. } p_E(x, z) = p^*(x) \cdot p_E(z|x) \quad \text{dec. } p_D(x, z) = q(z) p_D(x|z)$$

two requirements.

① decoder marginal $p(x) = \int p_D(x, z) dz \approx p^*(x)$ should be a good approx

② enc/dec pair must be self-consistent $p_E(x, z) \stackrel{!}{=} p_D(x, z)$

• maximizing ELBO loss enforces ②:

$$\downarrow \text{KL} [p_E(x, z) \parallel p_D(x, z)] = \mathbb{E}_{x \sim p^*(x)} \left[\mathbb{E}_{z \sim p_E(z|x)} \left[\log \frac{p_E(z|x) p^*(x)}{p_D(x|z) \cdot q(z)} \right] \right]$$

=

$$= \mathbb{E}_{x \sim p^*(x)} \left[\underbrace{\mathbb{E}_{z \sim p_E(z|x)} \left[\log \frac{p_E(z|x)}{q(z)} - \log p_D(x|z) \right]}_{-ELBO(x)} + \underbrace{\log p^*(x)}_{\text{independent of model}} \right]$$

$$\hat{E}, \hat{D} = \underset{E, D}{\text{arg min}} \mathbb{E}_{x \sim p^*(x)} \left[-ELBO(x) \right] \Rightarrow \text{ensures self-consistency upon perfect convergence} \Rightarrow \text{drop}$$

• maximizing the ELBO indirectly maximizes the data likelihood under the model (7)

(in literature, ELBO is usually maximized, conceptually, minimizing $-ELBO$ is simpler)

maximum likelihood principle. TS should be a typical (= high prob.) model outcome

\Leftrightarrow minimize expected negative log-likelihood (NLL)

$$\begin{aligned} \mathbb{E}_{x \sim p^*(x)} \left[-\log p(x) \right] &= \mathbb{E}_{x \sim p^*(x)} \left[-\log \int p_D(x|z) q(z) dz \right] \\ &= \mathbb{E}_{x \sim p^*(x)} \left[-\log \int p_E(z|x) \frac{p_D(x|z) q(z)}{p_E(z|x)} dz \right] \\ &= \mathbb{E}_{z \sim p_E(z|x)} \left[\frac{p_D(x|z) q(z)}{p_E(z|x)} \right] \end{aligned}$$

Jensen's inequality: if ψ is a convex function: $\psi(\mathbb{E}[t]) \leq \mathbb{E}[\psi(t)]$

$$\leq \mathbb{E}_{x \sim p^*(x)} \left[\underbrace{\mathbb{E}_{z \sim p_E(z|x)} \left[-\log \frac{p_D(x|z) q(z)}{p_E(z|x)} \right]}_{-ELBO(x)} \right]$$

$\Rightarrow \mathbb{E}_{x \sim p^*(x)} [-ELBO(x)]$ is upper bound for NLL loss
 minimization here leads to minimize this as well

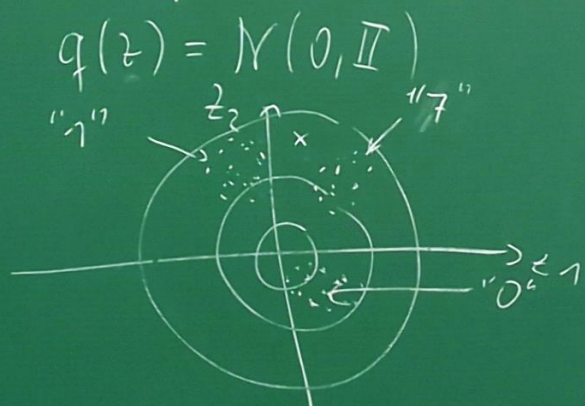
Interpretation of ELBO(x) terms: there is a trade-off between two objectives

$$-ELBO(x) = \underbrace{KL[p_E(z|x) \parallel q(z)]}_{\text{how close is conditional mode distribution } p_E(z|x) \text{ to desired } q(z)} + \underbrace{\mathbb{E}_{z \sim p_E(z|x)} [-\log p_D(x|z)]}_{\text{reconstruction error: original } x \text{ should be a very probable reconstruction}}$$

minimized if $p_E(z|x) = q(z)$
 ie encoder is independent of x ("ignores x ")
 \Rightarrow output $p(x)$ is unrelated to input, no reconstr.

minimized if encoder & decoder are deterministic with perfect reconstruction:
 $p_E(z|x) = \delta(z - f(x))$, $p_D(x|z) = \delta(x - g(z))$, $x = g(f(x))$

typical compromise in practice: MNIST (10 classes)



$$p_E(z) = \int p^*(x) p_E(z|x) dx \quad \text{has visible clusters}$$

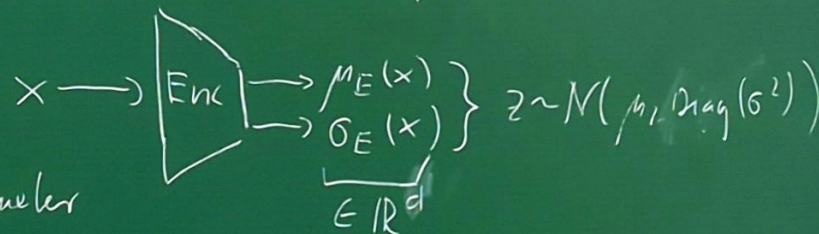
clusters for each digit with visible gaps

\Rightarrow sampling $z \sim q(z)$ can land in a gap

\Rightarrow bad synthetic point $p_D(x|z)$

most common implementation: encoder and decoder are diagonal Gaussians

$$p_E(z|x) = N(z | \mu_E(x), \text{Diag}(\sigma_E(x)^2))$$



$$p_D(x|z) = N(x | \mu_D(z), \beta^2 \mathbb{I})$$

hyperparameter

if $q(z) = N(0, \mathbb{I})$, $KL[p_E \| q]$ can be calculated analytically

$$KL[N(z | \mu_E(x), \text{Diag}(\sigma(x)^2)) \| N(0, \mathbb{I})] = \frac{1}{2} \sum_{j=1}^d (\mu_j(x)^2 + \sigma_j(x)^2 - \log \sigma_j(x)^2 - 1)$$

if $p_D(x|z) = N(x | \mu_D(z), \beta^2 \mathbb{I})$, reconstruction term becomes

$$\mathbb{E}_{z \sim p_E(z|x)} \{-\log p_D(x|z)\} = \mathbb{E}_{z \sim p_E(z|x)} \left[\frac{(x - \mu_D(z))^2}{2\beta^2} + \frac{\log(2\pi\beta^2)^{D/2}}{2\beta^2} \right]$$

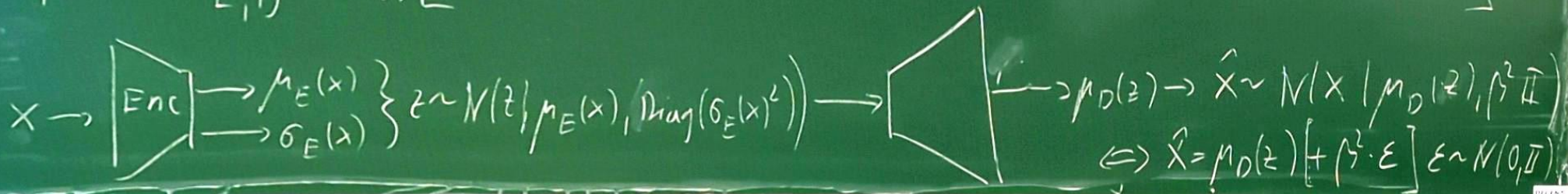
independent of model \Rightarrow drop

$$\left\{ z_k \sim p_E(z|x) \right\}_{k=1}^M \approx \frac{1}{M} \sum_{k=1}^M \frac{(x - \mu_D(z_k))^2}{2\beta^2}$$

$\beta^2 \gg 1$ downscale squared loss \Rightarrow reconstruction error unimportant
 $\beta^2 \ll 1$ upscale $-||-$ \Rightarrow $-||-$ dominates

full loss:

$$\hat{D}, \hat{E} = \arg \min_{E, D} \frac{1}{N} \sum_{i=1}^N \left[\sum_{j=1}^d \frac{1}{2} (\mu_{E_{ij}}(x_i)^2 + \sigma_{E_{ij}}(x_i)^2 - \log \sigma_{E_{ij}}(x_i)^2 - 1) + \frac{1}{M} \sum_{k=1}^M \frac{(x_i - \mu_D(z_{i,k}))^2}{2\beta^2} \right]$$



Paper [Kingma & Welling 2013] > 30000 citations

- generation $z \sim q(z)$, $x \sim p_D(x|z) \Leftrightarrow \mu_D(z) + \beta^2 \epsilon$
- inference if $p(x) = p^*(x)$ and $p_E(x,z) = p_D(x,z)$

$$\Rightarrow p_E(x,z) = p(x) p_E(z|x) = q(z) p_D(x|z) = p_D(x,z)$$

$$p(x) = \frac{q(z) \cdot p_D(x|z)}{p_E(z|x)} \quad \text{must give same value for all } z \sim p_E(z|x)$$

[measure of self-consistency: variance of $\log q(z) + \log p_D(x|z) - \log p_E(z|x)$]

conditional VAE

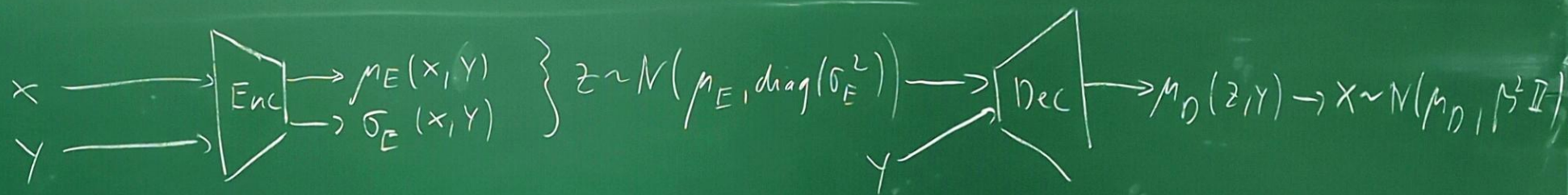
- plain VAE learns $p(x)$, conditional VAE learns $p(x|Y)$ for some variable Y
ex $Y \in \{0, \dots, 9\}$ MNIST label

$$p(x|Y) = \int q(z) p_D(x|Y,z) dz \quad \Rightarrow \quad p(x) = \sum_Y p(x|Y) p^*(Y) \quad \begin{array}{l} \text{hopefully} \\ \text{better quality} \end{array}$$

model change add Y is an extra condition to enc. & dec.

$$p_E(z|x, Y) \quad p_D(x|z, Y)$$

if enc. & dec. are Gaussian, add Y as input to the networks



application: style transfer $Y = \{ \text{painting, foto, carbon, drawing} \}$

encode x with true Y (e.g. painting of Newton)

decode with different Y (e.g. foto of Newton)

generative classification: test x with unknown label
 \Rightarrow try encoding with every label and calculate the corresponding probs $p(x|Y)$
 return the label maximizing the probs.

$$\hat{Y} = \arg \max_Y \int p_E(z|x, Y) p_D(x|z, Y) dz$$