

Autoencoders

"Manifold hypothesis" - data lives on a low-dimensional manifold

Dimensionality reduction

$$x \in \mathbb{R}^D \rightarrow z \in \mathbb{R}^d \quad d < D \quad (d \ll D)$$

$$f(x) = z$$

- f - linear $X \in \mathbb{R}^{N \times D}$ N - num of training objects

$$Z \in \mathbb{R}^{N \times d} \quad Z = X \odot H \quad H \in \mathbb{R}^{D \times d}$$

Reconstruction $\hat{X} = Z \Phi \quad \Phi \in \mathbb{R}^{d \times D}$

$$X \approx \hat{X} = Z \Phi = X \odot H \Phi$$

$$MSE = \|X - \hat{X}\|_F^2 = \|X - \underbrace{X \odot H}_{\text{rank } d} \underbrace{\Phi}_{d \times D}\|_F^2 \rightarrow \min_{H, \Phi}$$

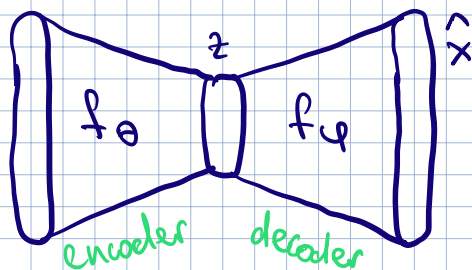
$$\|A\|_F = \sqrt{\sum_{ij} |A_{ij}|^2}$$

• Principal Component Analysis (PCA)

- f - non linear

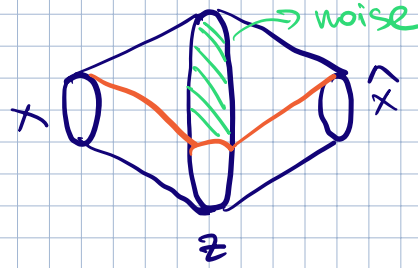
\mathcal{X} - data distribution

$$\mathbb{E}_{x \sim \mathcal{X}} \|x - \hat{x}\|_2^2 =$$



$$= \mathbb{E}_{x \sim \mathcal{X}} \|x - f_\phi(f_\theta(x))\|_2^2 \rightarrow \min_{\phi, \theta}$$

$$d > D$$

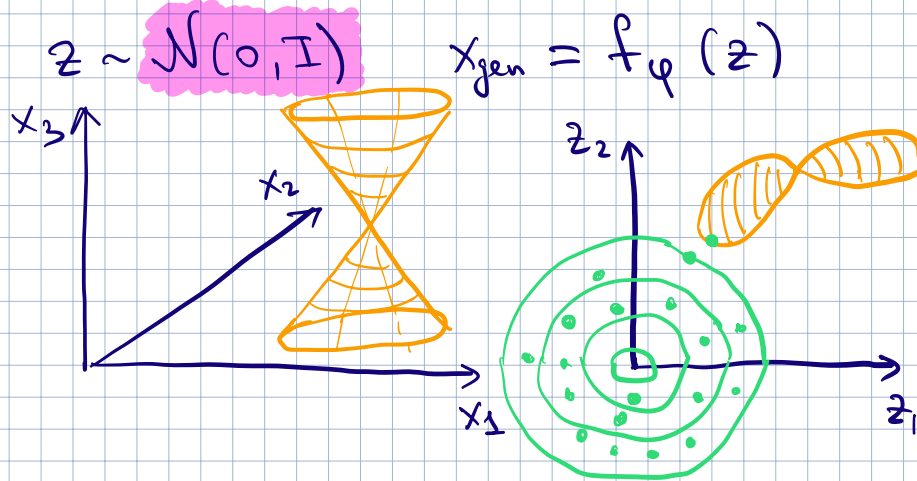


Denoising Autoencoder

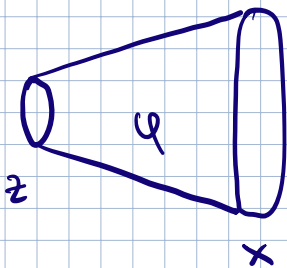
$$\mathbb{E}_{x \sim \mathcal{X}} \|x - f_{\psi}(f_{\theta}(\tilde{x}))\|_2^2 \rightarrow \min_{\psi, \theta}$$

$$\tilde{x} = x + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I)$$

Sampling from Autoencoder



Variational Autoencoder (VAE)



$$p(x, z | \psi) = p(x | z, \psi) p(z)$$

$$p(x | z, \psi) = \mathcal{N}(x | \underbrace{\mu_{\psi}(z)}_{\in \mathbb{R}^D}, \underbrace{\Sigma_{\psi}(z)}_{\in \mathbb{R}^{D \times D}})$$

$$p(z) = \mathcal{N}(z | 0, I)$$

x - observe z - latent

$$\log p(x|\varphi) = \log \int p(x, z|\varphi) dz \rightarrow \max_{\varphi}$$

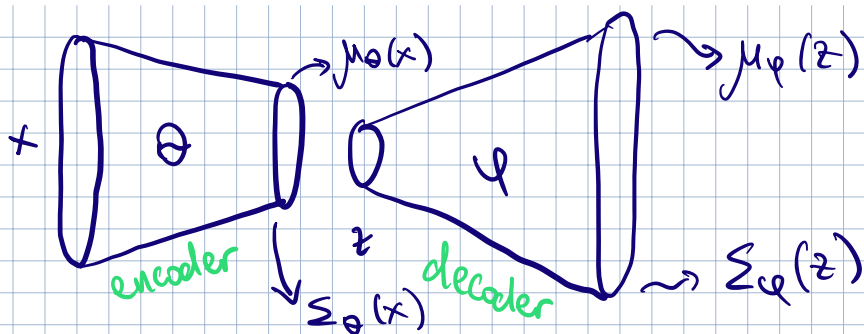
Evidence lower bound (ELBO)

$$\begin{aligned} \log p(x|\varphi) &= \log p(x|\varphi) \cdot 1 = \int q(z) \log p(x|\varphi) dz \quad \text{arbitrary distribution over } z \\ &= \log p(x|\varphi) \cdot \int q(z) dz = \int q(z) \log p(x|\varphi) dz = \\ &= \int p(x, z|\varphi) = p(z|x, \varphi) \cdot p(x|\varphi) \} = \\ &= \int q(z) \log \frac{p(x, z|\varphi)}{p(z|x, \varphi)} dz = \\ &= \int q(z) \log \frac{p(x, z|\varphi)}{p(z|x, \varphi)} \cdot \frac{q(z)}{q(z)} dz = \underbrace{\int q(z) \log \frac{p(x, z|\varphi)}{q(z)} dz}_{\text{ELBO}} + \underbrace{\int q(z) \log \frac{q(z)}{p(z|x, \varphi)} dz}_{\text{KL}(q||p(z|x, \varphi))} = \\ &= \underbrace{\mathcal{L}(q, \varphi)}_{\text{ELBO}} + \text{KL}(q||p(z|x, \varphi)) \\ &\quad q = p(z|x, \varphi) \quad \text{KL} = 0 \\ &\quad \mathcal{L}(q, \varphi) = \log p(x|\varphi) \end{aligned}$$

q - arbitrary $q(z) := q(z|x, \theta)$

$$\mathcal{L}(q, \varphi) = \int q(z|x, \theta) \log \frac{p(x, z|\varphi)}{q(z|x, \theta)} dz \rightarrow \max_{\varphi, \theta}$$

$$q(z|x, \theta) = \mathcal{N}(z|\mu_{\theta}(x), \Sigma_{\theta}(x))$$



$$\mu_\theta(x) \in \mathbb{R}^d$$

$$\mu_\psi(z) \in \mathbb{R}^D$$

$$\Sigma_\theta(x) = \text{diag}(v_{\theta 1}^2(x), \dots, v_{\theta d}^2(x)) \quad \Sigma_\psi(z) = \text{diag}(v_{\psi 1}(z), \dots, v_{\psi D}(z))$$

$$\Sigma_\psi(z) = \mathbf{I}$$

$$\mathcal{L}(\theta, \psi) = \int q(z|x, \theta) \log \frac{p(x, z|\psi)}{q(z|x, \theta)} dz =$$

$$= \int p(x, z|\psi) = p(x|z, \psi) p(z) \} =$$

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

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

$$\underbrace{\mathbb{E}_{q(z|x, \theta)} \log p(x|z, \psi)}_{\text{(JSE) reconstruction loss}} - \underbrace{\text{KL}(q(z|x, \theta) || p(z))}_{\text{regularization}} \rightarrow \max_{\psi, \theta}$$

$$\mathbb{E}_{q(z|x, \theta)} \log p(x|z, \psi) \rightarrow \nabla_\psi, \nabla_\theta$$

$$\nabla_\psi \mathbb{E}_{q(z|x, \theta)} \log p(x|z, \psi)$$

$$\nabla_{\theta} \mathbb{E}_{q(z|x, \theta)} \log p(x|z, \psi)$$

Reparametrization trick

$$\xi \sim \mathcal{N}(\xi | m, s^2) \quad \varepsilon \sim \mathcal{N}(\varepsilon | 0, 1) \quad \xi = m + s \cdot \varepsilon$$

$$z \sim \mathcal{N}(z | \mu_{\theta}(x), \Sigma_{\theta}(x))$$

$$\varepsilon \sim \mathcal{N}(\varepsilon | 0, I) \quad \underline{z = \mu_{\theta}(x) + \Sigma_{\theta}(x)^{1/2} \cdot \varepsilon}$$

$$\nabla_{\theta} \mathbb{E}_{q(z|x, \theta)} \log p(x|z, \psi) = \nabla_{\theta} \mathbb{E}_{\varepsilon \sim \mathcal{N}(0, I)} \log p(x | \underset{z(\mu_{\theta}(x), \Sigma_{\theta}(x), \varepsilon)}{\parallel} z, \psi)$$

$$\begin{aligned} \log p(x|z, \psi) &= \log \left[\frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp \left(\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right) \right] = \\ &= -\frac{D}{2} \log 2\pi - \frac{1}{2} \log |\Sigma| + \frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) = \\ &= -\frac{1}{2} \sum_{i=1}^D \log \sigma_i^2 + \frac{1}{2} \sum_{i=1}^D \frac{(x_i - \mu_i)^2}{\sigma_i^2} \propto \{ \Sigma_{\psi}(z) = I \} \propto \\ &\propto \frac{1}{2} \sum_{i=1}^D (x_i - \mu_i)^2 \end{aligned}$$

↑ object
↑ decoder output

$$q(z|x, \theta) = \mathcal{N}(z | \mu_{\theta}, \Sigma_{\theta}) \quad p(z) = \mathcal{N}(z | 0, I)$$

$$\xi_1 \sim \mathcal{N}(\xi_1 | \mu_1, \Sigma_1) \quad \xi_2 \sim \mathcal{N}(\xi_2 | \mu_2, \Sigma_2)$$

$$KL(\xi_1 \parallel \xi_2) = \frac{1}{2} \left[\text{tr}(\Sigma_2^{-1} \Sigma_1) - d + (\mu_2 - \mu_1)^T \Sigma_2^{-1} (\mu_2 - \mu_1) + \log \frac{|\Sigma_2|}{|\Sigma_1|} \right]$$

$$\Sigma_0 = \text{diag}(v_1^2, \dots, v_d^2)$$

$$\begin{aligned} \text{KL}(q||p) &= \frac{1}{2} [\text{tr}(\Sigma_0) - d + \|\mu_0\|_2^2 - \log |\Sigma_0|] = \\ &= \frac{1}{2} \times \sum_{i=1}^d (v_i^2 + \mu_i^2 - \log v_i^2) - \frac{d}{2} \end{aligned}$$

$\log v_i^2$ - output of linear layer

$$v_i^2 = \exp(\log v_i^2)$$