

# Deep Generative Models

## Probabilistic PCA

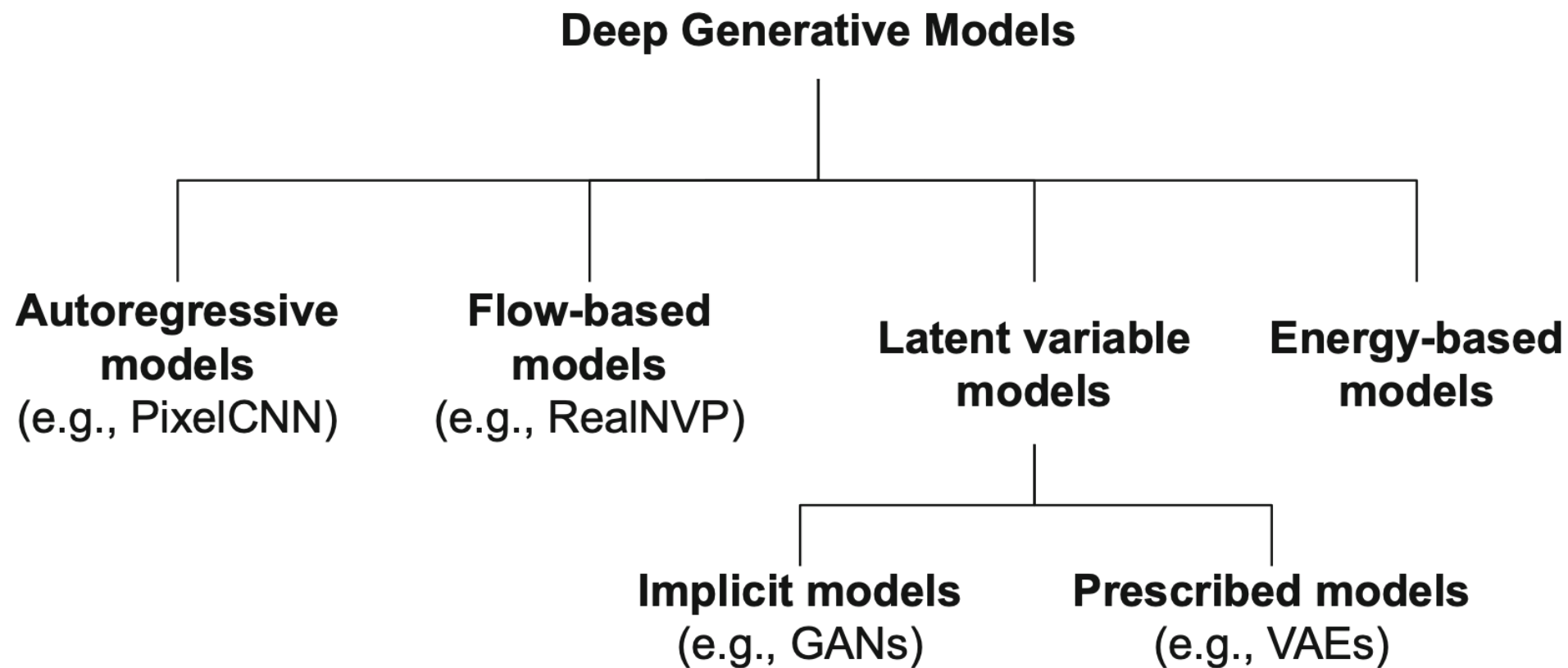
Fall Semester 2025

René Vidal

Director of the Center for Innovation in Data Engineering and Science (IDEAS),  
Rachleff University Professor, University of Pennsylvania  
Amazon Scholar & Chief Scientist at NORCE

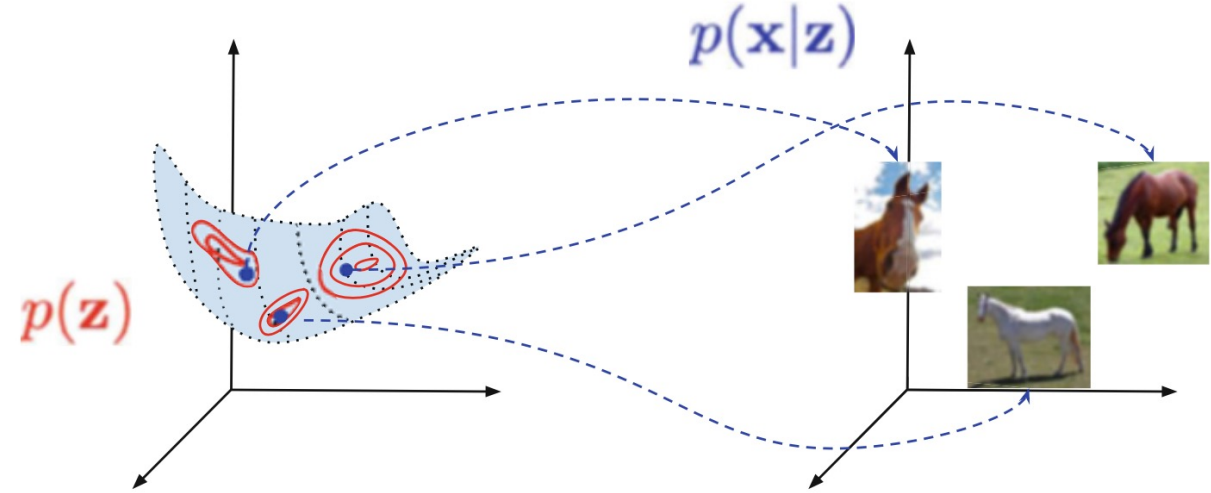


# Taxonomy of Generative Models



# Latent Variable Models

- $X$  = observed variable
- $Z$  = latent variable
- $\mathbf{z} \sim p(\mathbf{z})$
- $\mathbf{x} \sim p(\mathbf{x} | \mathbf{z})$



A latent variable model and a generative process. Note the low-dimensional manifold (here 2D) embedded in the high-dimensional space (here 3D)

- Factorization of the joint model

$$p(\mathbf{x}, \mathbf{z}) = p(\mathbf{x} | \mathbf{z})p(\mathbf{z})$$

- Marginalization of the model

$$p(\mathbf{x}) = \int p(\mathbf{x} | \mathbf{z})p(\mathbf{z})d\mathbf{z}$$

# Probabilistic Principal Component Analysis: Model

- We consider continuous random variables only, i.e.,

$$\mathbf{z} \in \mathbb{R}^d \text{ and } \mathbf{x} \in \mathbb{R}^D \text{ with } d \ll D$$

- The distribution of  $\mathbf{z}$  is the standard Gaussian, i.e.,

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z} \mid \mathbf{0}, \mathbf{I}).$$

- The dependency between  $\mathbf{z}$  and  $\mathbf{x}$  is linear and we assume a Gaussian additive noise:

$$\mathbf{x} = \mathbf{W}\mathbf{z} + \mathbf{b} + \boldsymbol{\varepsilon}$$

- Here  $\boldsymbol{\varepsilon} \sim \mathcal{N}(\boldsymbol{\varepsilon} \mid \mathbf{0}, \sigma^2 \mathbf{I})$  and independent from  $\mathbf{z}$ .

# Probabilistic Principal Component Analysis: Model

- PPCA Model

$$\mathbf{x} = \mathbf{W}\mathbf{z} + \mathbf{b} + \boldsymbol{\epsilon}, \quad \mathbf{z} \sim \mathcal{N}(\mathbf{z} \mid \mathbf{0}, \mathbf{I}), \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\boldsymbol{\epsilon} \mid \mathbf{0}, \sigma^2 \mathbf{I}).$$

- $\mathbf{x}$  is a linear combination of Gaussians, thus  $p(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \mathbf{b}, \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I})$  because

$$\begin{aligned} \mathbb{E}[\mathbf{x}] &= \mathbb{E}[\mathbf{W}\mathbf{z}] + \mathbb{E}[\mathbf{b}] + \mathbb{E}[\boldsymbol{\epsilon}] = \mathbf{W}\mathbb{E}[\mathbf{z}] + \mathbf{b} + \mathbf{0} = \mathbf{b} \\ \mathbb{V}[\mathbf{x}] &= \mathbb{V}[\mathbf{W}\mathbf{z} + \mathbf{b} + \boldsymbol{\epsilon}] = \mathbf{W}\mathbb{V}(\mathbf{z})\mathbf{W}^\top + \mathbb{V}[\boldsymbol{\epsilon}] = \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I} \end{aligned}$$

- $\mathbf{x} \mid \mathbf{z}$  is a constant + a Gaussian, thus  $p(\mathbf{x} \mid \mathbf{z}) = \mathcal{N}(\mathbf{x} \mid \mathbf{W}\mathbf{z} + \mathbf{b}, \sigma^2 \mathbf{I})$  because

$$\begin{aligned} \mathbb{E}[\mathbf{x} \mid \mathbf{z}] &= \mathbf{W}\mathbf{z} + \mathbf{b} + \mathbb{E}[\boldsymbol{\epsilon}] = \mathbf{W}\mathbf{z} + \mathbf{b} \\ \mathbb{V}[\mathbf{x} \mid \mathbf{z}] &= \mathbb{V}[\boldsymbol{\epsilon}] = \sigma^2 \mathbf{I} \end{aligned}$$

# Probabilistic Principal Component Analysis: Model

- PPCA model:  $\mathbf{x} = \mathbf{W}\mathbf{z} + \mathbf{b} + \boldsymbol{\epsilon}$ ,  $p(\mathbf{z}) = \mathcal{N}(\mathbf{z} \mid \mathbf{0}, \mathbf{I})$ ,  $p(\boldsymbol{\epsilon}) = \mathcal{N}(\boldsymbol{\epsilon} \mid \mathbf{0}, \sigma^2 \mathbf{I})$ ,

$$p(\mathbf{x} \mid \mathbf{z}) = \mathcal{N}(\mathbf{x} \mid \mathbf{W}\mathbf{z} + \mathbf{b}, \sigma^2 \mathbf{I}), \quad p(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \mathbf{b}, \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I})$$

- Let  $\mathbf{M} = \mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I}$ . We can compute the conditional distribution of  $(\mathbf{z} \mid \mathbf{x})$  as

$$p(\mathbf{z} \mid \mathbf{x}) = \frac{p(\mathbf{x} \mid \mathbf{z})p(\mathbf{z})}{p(\mathbf{x})} \propto e^{-\frac{1}{2\sigma^2}\|\mathbf{x} - \mathbf{W}\mathbf{z} - \mathbf{b}\|^2} e^{-\frac{1}{2}\|\mathbf{z}\|^2}$$

$$p(\mathbf{z} \mid \mathbf{x}) \propto e^{-\frac{1}{2\sigma^2}(\mathbf{z}^\top \mathbf{W}^\top \mathbf{W} \mathbf{z} - 2\mathbf{z}^\top \mathbf{W}^\top (\mathbf{x} - \mathbf{b}) + \sigma^2 \|\mathbf{z}\|^2)} \propto e^{-\frac{1}{2\sigma^2}(\mathbf{z}^\top \mathbf{M} \mathbf{z} - 2\mathbf{z}^\top \mathbf{W}^\top (\mathbf{x} - \mathbf{b}))}$$

$$p(\mathbf{z} \mid \mathbf{x}) = \mathcal{N}(\mathbf{z} \mid \mathbf{M}^{-1} \mathbf{W}^\top (\mathbf{x} - \mathbf{b}), \sigma^2 \mathbf{M}^{-1})$$

# Probabilistic Principal Component Analysis: Learning

- Recall the ML estimators of the parameters of a Gaussian  $\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma})$  are

$$\boldsymbol{\mu}_N = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i, \quad \boldsymbol{\Sigma}_N = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}_N)(\mathbf{x}_i - \boldsymbol{\mu}_N)^T$$

- For PPCA we need to estimate the parameters of a Gaussian with structured covariance  $\boldsymbol{\Sigma} = \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}$ . The estimate of the mean is the same as before  $\boldsymbol{\mu} = \boldsymbol{\mu}_N$ . To estimate  $\mathbf{W}$ , we need to maximize the log-likelihood w.r.t.  $(\mathbf{W}, \sigma)$

$$\ell = -\frac{N}{2} \log(\det(\boldsymbol{\Sigma})) - \frac{N}{2} \text{trace}(\boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma}_N)$$

- Taking derivatives w.r.t.  $\mathbf{W}$  we get

$$\frac{\partial \ell}{\partial \mathbf{W}} = \frac{\partial \ell}{\partial \boldsymbol{\Sigma}} \frac{\partial \boldsymbol{\Sigma}}{\partial \mathbf{W}} = -\frac{N}{2} (\boldsymbol{\Sigma}^{-1} - \boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma}_N \boldsymbol{\Sigma}^{-1}) 2\mathbf{W} = 0 \Rightarrow \boldsymbol{\Sigma}_N \boldsymbol{\Sigma}^{-1} \mathbf{W} = \mathbf{W}$$

# Probabilistic Principal Component Analysis: Learning

- We thus need to solve the nonlinear equations

$$\mathbf{\Sigma}_N \mathbf{\Sigma}^{-1} \mathbf{W} = \mathbf{W} \quad \text{and} \quad \mathbf{\Sigma} = \mathbf{W} \mathbf{W}^T + \sigma^2 \mathbf{I}$$

- A trivial solution is  $\mathbf{W} = 0$ , but this is a minimum of the log-likelihood.
- Another solution is  $\mathbf{\Sigma} = \mathbf{\Sigma}_N$ , but this would require the structure of the sample covariance  $\mathbf{\Sigma}_N$  to match the structure of  $\mathbf{\Sigma} = \mathbf{W} \mathbf{W}^T + \sigma^2 \mathbf{I}$ , i.e., the smallest eigenvalues would need to be all equal to each other and equal to  $\sigma^2$ .
- Alternatively, let

$$\mathbf{W} = [\mathbf{Z}_1 \quad \mathbf{Z}_2] \begin{bmatrix} \mathbf{\Gamma}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} [\mathbf{V}_1 \quad \mathbf{V}_2]^T = \mathbf{Z}_1 \mathbf{\Gamma}_1 \mathbf{V}_1^T$$

- Then

$$\mathbf{\Sigma} = \mathbf{W} \mathbf{W}^T + \sigma^2 \mathbf{I} = \mathbf{Z}_1 \mathbf{\Gamma}_1^2 \mathbf{Z}_1^T + \sigma^2 (\mathbf{Z}_1 \mathbf{Z}_1^T + \mathbf{Z}_2 \mathbf{Z}_2^T) = \mathbf{Z}_1 (\mathbf{\Gamma}_1^2 + \sigma^2 \mathbf{I}) \mathbf{Z}_1^T + \sigma^2 \mathbf{Z}_2 \mathbf{Z}_2^T$$

$$\mathbf{\Sigma}^{-1} \mathbf{W} = (\mathbf{Z}_1 (\mathbf{\Gamma}_1^2 + \sigma^2 \mathbf{I})^{-1} \mathbf{Z}_1^T + \sigma^{-2} \mathbf{Z}_2 \mathbf{Z}_2^T) \mathbf{Z}_1 \mathbf{\Gamma}_1 \mathbf{V}_1^T = \mathbf{Z}_1 (\mathbf{\Gamma}_1^2 + \sigma^2 \mathbf{I})^{-1} \mathbf{\Gamma}_1 \mathbf{V}_1^T$$



# Probabilistic Principal Component Analysis: Learning

- Therefore,

$$\begin{aligned}\Sigma_N \Sigma^{-1} \mathbf{W} = \mathbf{W} &\Rightarrow \Sigma_N \mathbf{Z}_1 (\Gamma_1^2 + \sigma^2 \mathbf{I})^{-1} \Gamma_1 \mathbf{V}_1^T = \mathbf{Z}_1 \Gamma_1 \mathbf{V}_1^T \Rightarrow \\ \Sigma_N \mathbf{Z}_1 (\Gamma_1^2 + \sigma^2 \mathbf{I})^{-1} &= \mathbf{Z}_1 \Rightarrow \Sigma_N \mathbf{Z}_1 = \mathbf{Z}_1 (\Gamma_1^2 + \sigma^2 \mathbf{I}) \Rightarrow \Sigma_N \mathbf{z}_i = (\gamma_i^2 + \sigma^2) \mathbf{z}_i\end{aligned}$$

- In other words,  $\mathbf{z}_i$  is an eigenvector of  $\Sigma_N$  with eigenvalue  $\gamma_i^2 + \sigma^2$ .

- Thus if  $\Sigma_N = [\mathbf{U}_1 \ \mathbf{U}_2] \begin{bmatrix} \Lambda_1 & \mathbf{0} \\ \mathbf{0} & \Lambda_2 \end{bmatrix} [\mathbf{U}_1 \ \mathbf{U}_2]^T$ , then  $\mathbf{Z}_1 = \mathbf{U}_1$ ,  $\Gamma_1^2 + \sigma^2 \mathbf{I} = \Lambda_1$ .

- In other words,  $\mathbf{U}_1$  is a matrix whose  $d$  columns correspond to  $d$  singular vectors of  $\Sigma_N$

- Therefore,  $\mathbf{W} = \mathbf{Z}_1 \Gamma_1 \mathbf{V}_1^T = \mathbf{U}_1 (\Lambda_1 - \sigma^2 \mathbf{I})^{1/2} \mathbf{V}_1^T$

- Having “almost” found  $\mathbf{W}$  (we don’t know which  $d$  columns), we now turn to finding  $\sigma$ .

# Probabilistic Principal Component Analysis: Learning

- Recall the log-likelihood

$$\ell = -\frac{N}{2} \log(\det(\mathbf{\Sigma})) - \frac{N}{2} \text{trace}(\mathbf{\Sigma}^{-1} \mathbf{\Sigma}_N)$$

- We have  $\mathbf{\Sigma} = \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}$  and  $\mathbf{W} = \mathbf{U}_1 \mathbf{\Gamma}_1 \mathbf{V}_1^T$ . Thus,  $\mathbf{W}\mathbf{W}^T = \mathbf{U}_1 \mathbf{\Gamma}_1^2 \mathbf{U}_1^T$  and

$$\mathbf{\Sigma} = \mathbf{U}_1 (\mathbf{\Gamma}_1^2 + \sigma^2 \mathbf{I}) \mathbf{U}_1^T + \sigma^2 \mathbf{U}_2 \mathbf{U}_2^T = \mathbf{U}_1 \mathbf{\Lambda}_1 \mathbf{U}_1^T + \sigma^2 \mathbf{U}_2 \mathbf{U}_2^T$$

$$\mathbf{\Sigma}^{-1} \mathbf{\Sigma}_N = (\mathbf{U}_1 \mathbf{\Lambda}_1^{-1} \mathbf{U}_1^T + \sigma^{-2} \mathbf{U}_2 \mathbf{U}_2^T) (\mathbf{U}_1 \mathbf{\Lambda}_1 \mathbf{U}_1^T + \mathbf{U}_2 \mathbf{\Lambda}_2 \mathbf{U}_2^T) = \mathbf{U}_1 \mathbf{U}_1^T + \sigma^{-2} \mathbf{U}_2 \mathbf{\Lambda}_2 \mathbf{U}_2^T$$

- Substituting into the log-likelihood, we get

$$\ell = -\frac{N}{2} \log(\det(\mathbf{\Lambda}_1) \sigma^{2(D-d)}) - \frac{N}{2} (d + \sigma^{-2} \text{trace}(\mathbf{\Lambda}_2))$$

# Probabilistic Principal Component Analysis: Learning

- Taking the derivative yields  $\frac{\partial \ell}{\partial \sigma^2} = -\frac{N}{2} \left( \frac{D-d}{\sigma^2} - \frac{\text{trace}(\Lambda_2)}{\sigma^4} \right) = 0 \Rightarrow \sigma^2 = \frac{\text{trace}(\Lambda_2)}{D-d}$
- Final piece: we have not shown  $\Lambda_1$  corresponds to **top**  $d$  eigenvalues of  $\Sigma_N$ 
  - Exercise 2.13 in GPCA textbook

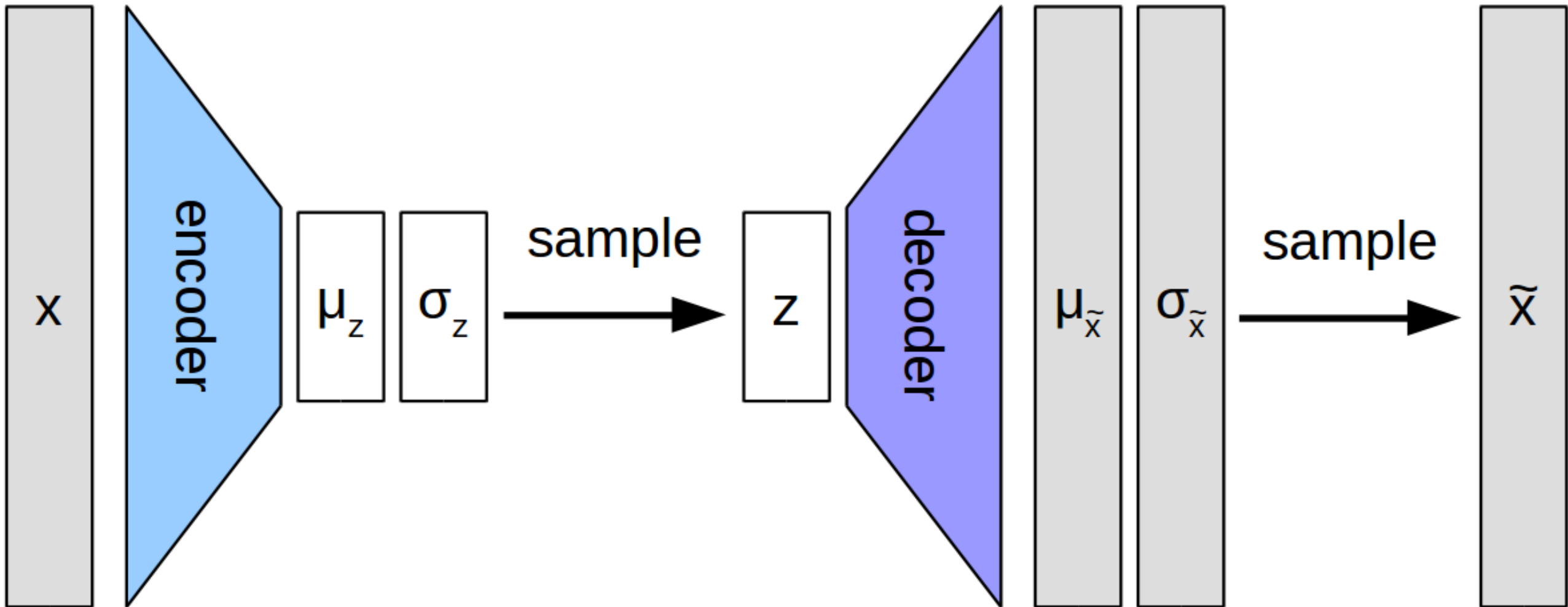
- **Theorem.** The ML estimates for the parameters of the PPCA model  $\mathbf{b}$ ,  $\mathbf{W}$ , and  $\sigma$  can be obtained from the ML estimates of the mean and covariance of the data,  $\mu_N$  and  $\Sigma_N$ , respectively, as

$$\mathbf{b} = \mu_N, \mathbf{W} = \mathbf{U}_1(\Lambda_1 - \sigma^2 I)^{1/2} R \text{ and } \sigma^2 = \frac{1}{D-d} \sum_{i=d+1}^D \lambda_i$$

- where  $\mathbf{U}_1$  is the matrix with the top  $d$  eigenvectors of  $\Sigma_N$ ,  $\Lambda_1$  is the matrix with the corresponding top  $d$  eigenvalues,  $R \in \mathbb{R}^{d \times d}$  is an arbitrary orthogonal matrix, and  $\lambda_i$  is the  $i$ th largest eigenvalue of  $\Sigma_N$ .

# PPCA as an Encoder Decoder Architecture

•  $p(z|x) = \mathcal{N}(z \mid M^{-1}W^{\top}(x - b), \sigma^{-2}M)$        $p(x \mid z) = \mathcal{N}(x \mid Wz + b, \sigma^2I)$



# Application of PPCA to Generating Face Images



**Fig. 2.2** Face images of subject 20 under 10 different illumination conditions in the extended Yale B data set. All images are frontal faces cropped to size  $192 \times 168$ .

# Application of PPCA to Generating Face Images



(a) mean face



(b) first eigenface



(c) second eigenface

**Fig. 2.5** Mean face and the first two eigenfaces by applying PPCA to the ten images in Figure 2.2.

# Application of PPCA to Generating Face Images



(a) Variation along the first eigenface



(b) Variation along the second eigenface

**Fig. 2.6** Variation of the face images along the two eigenfaces given by PPCA. Each row plots  $\mu + y_i \mathbf{u}_i$  for  $y_i = -1 : \frac{1}{3} : 1, i = 1, 2$ .