

Deep Generative Models

Probabilistic PCA

Fall Semester 2025

René Vidal

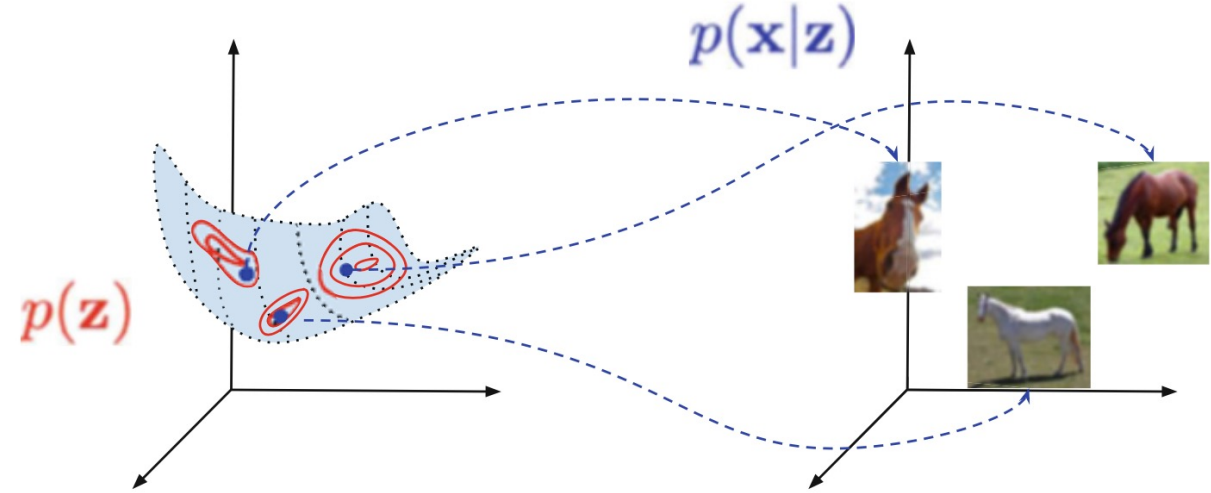
Director of the Center for Innovation in Data Engineering and Science (IDEAS),
Rachleff University Professor, University of Pennsylvania
Amazon Scholar & Chief Scientist at NORCE

Outline

- Probabilistic PCA: Model
 - Joint, Marginal and Conditional Distributions
 - PPCA as an Encoder-Decoder Architecture
- Linear Algebra Background
 - Singular Value Decomposition (SVD)
 - SVD Properties
- Probabilistic PCA: Learning
 - Maximum Likelihood Estimation
- Applications
 - Application of PPCA to Generating Face Images

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, plus 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

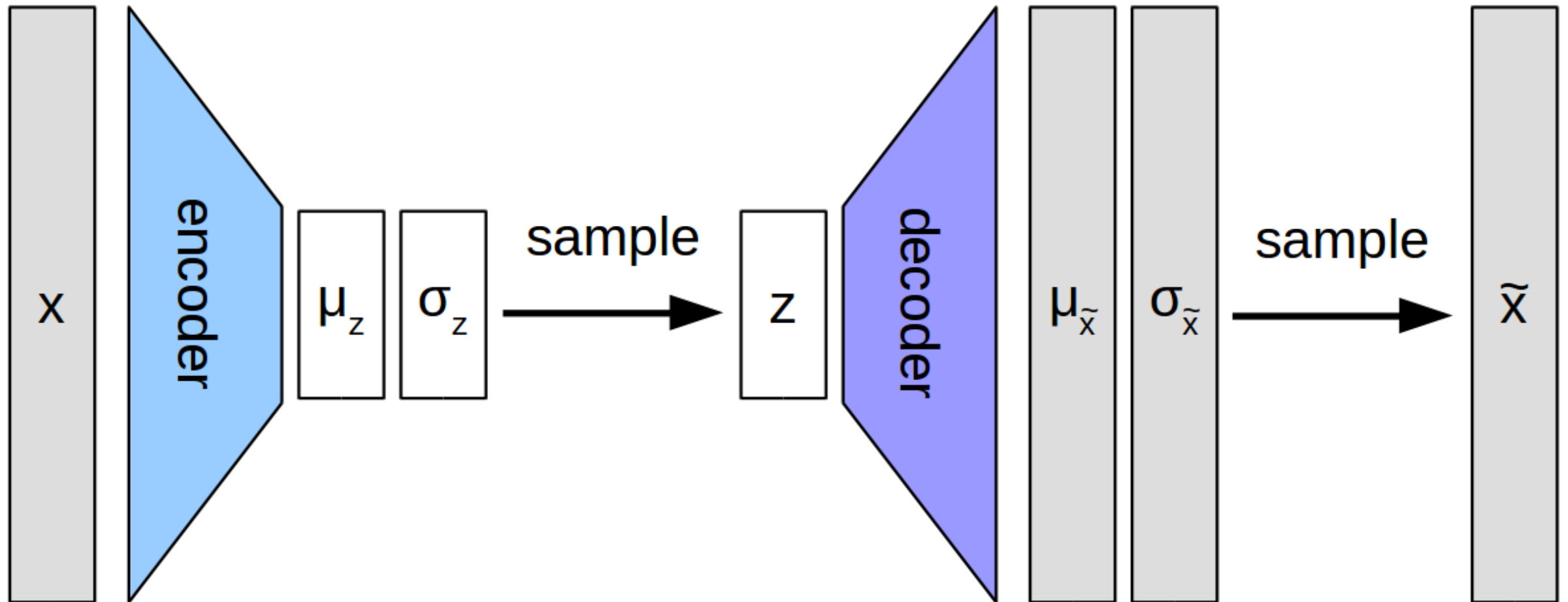
$$p(\mathbf{z} \mid \mathbf{x}) = \frac{p(\mathbf{x} \mid \mathbf{z})p(\mathbf{z})}{p(\mathbf{x})} \propto e^{-\frac{1}{2\sigma^2}\|\mathbf{x} - \mathbf{W}\mathbf{z} - \mathbf{b}\|^2} e^{-\frac{1}{2}\|\mathbf{z}\|^2}$$

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

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

PPCA as an Encoder Decoder Architecture

- PPCA model: $\mathbf{x} = \mathbf{W}\mathbf{z} + \mathbf{b} + \epsilon$, $p(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I})$, $p(\epsilon) = \mathcal{N}(\epsilon | \mathbf{0}, \sigma^2 \mathbf{I})$,
 $p(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\mathbf{z} | M^{-1}W^\top(\mathbf{x} - \mathbf{b}), \sigma^{-2}M)$ $p(\mathbf{x} | \mathbf{z}) = \mathcal{N}(\mathbf{x} | W\mathbf{z} + \mathbf{b}, \sigma^2 \mathbf{I})$



Singular Value Decomposition (SVD)

- A matrix $X \in \mathbb{R}^{m \times n}$ of rank r can be decomposed as

$$X = U\Sigma V^\top = [u_1 \cdots u_r u_{r+1} \cdots u_m] \begin{bmatrix} \sigma_1 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_r & 0 & \cdots & 0 \\ 0 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & 0 & \cdots & 0 \end{bmatrix} [v_1 \cdots v_r v_{r+1} \cdots v_m]^\top$$

- $\Sigma \in \mathbb{R}^{m \times n}$ is a “diagonal” matrix with r singular values $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_r > 0$
- $U \in \mathbb{R}^{m \times m}$ and $V \in \mathbb{R}^{n \times n}$ are orthonormal matrices of left/right singular vectors
- Orthonormal means that columns are unit norm and orthogonal to each other:

$$U^\top U = UU^\top = I_{m \times m},$$

$$V^\top V = VV^\top = I_{n \times n}$$

Compact Singular Value Decomposition (SVD)

- A matrix $X \in \mathbb{R}^{m \times n}$ of rank r can be decomposed as

$$\begin{aligned} X = U\Sigma V^\top &= [u_1 \cdots u_r u_{r+1} \cdots u_m] \begin{bmatrix} \sigma_1 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_r & 0 & \cdots & 0 \\ 0 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & 0 & \cdots & 0 \end{bmatrix} [v_1 \cdots v_r v_{r+1} \cdots v_m]^\top \\ &= U_1 \Sigma_1 V_1^\top = [u_1 \cdots u_r] \begin{bmatrix} \sigma_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_r \end{bmatrix} [v_1 \cdots v_r]^\top \end{aligned}$$

- $U_1 \in \mathbb{R}^{m \times r}$ and $V \in \mathbb{R}^{n \times r}$ are orthonormal matrices of left/right singular vectors
- $\Sigma \in \mathbb{R}^{r \times r}$ is a diagonal matrix with r singular values $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_r > 0$

SVD Properties

- We can use any partition of the SVD of a matrix to expand it as

$$X = U\Sigma V^{\top} = [U_1 \quad U_2] \begin{bmatrix} \Sigma_1 & 0 \\ 0 & \Sigma_2 \end{bmatrix} [V_1 \quad V_2]^{\top} = U_1 \Sigma_1 V_1^{\top} + U_2 \Sigma_2 V_2^{\top}$$

- **Example:** $I = U_1 U_1^{\top} + U_2 U_2^{\top}$, where $U_1^{\top} U_1 = I$, $U_1^{\top} U_2 = 0$, and $U_2^{\top} U_2 = I$
- If X is invertible, we can use its SVD to compute the inverse of a matrix

$$X = U\Sigma V^{\top} \Rightarrow X^{-1} = V\Sigma^{-1} U^{\top}$$

- We can use any partition of the SVD of a matrix to compute its inverse as

$$X = U_1 \Sigma_1 V_1^{\top} + U_2 \Sigma_2 V_2^{\top} \Rightarrow X^{-1} = V_1 \Sigma_1^{-1} U_1^{\top} + V_2 \Sigma_2^{-1} U_2^{\top}$$

SVD Properties

- We can use the SVD of a square matrix to compute its trace and determinant

$$X = U\Sigma V^{\top} \Rightarrow \text{trace}(X) = \text{trace}(\Sigma) = \sum_i \sigma_i \quad \text{and} \quad \det(X) = \det(\Sigma) = \prod_i \sigma_i$$

- We can use any partition of the SVD to compute its trace and determinant

$$X = U_1\Sigma_1V_1^{\top} + U_2\Sigma_2V_2^{\top} \Rightarrow \text{trace}(X) = \text{trace}(\Sigma_1) + \text{trace}(\Sigma_2)$$

$$X = U_1\Sigma_1V_1^{\top} + U_2\Sigma_2V_2^{\top} \Rightarrow \det(X) = \det(\Sigma_1)\det(\Sigma_2)$$

- Example

$$X = \begin{bmatrix} I_m & 0 \\ 0 & \sigma^2 I_n \end{bmatrix} \Rightarrow \text{trace}(X) = m + \sigma^2 n \quad \text{and} \quad \det(X) = \sigma^{2n}$$

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}, σ)

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

- 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 $(\mathbf{W}, \sigma) = (\mathbf{0}, 0)$, 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 σ^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 Σ_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 Σ_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 σ .

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 that Λ_1 corresponds to **top** d eigenvalues of Σ_N
 - Exercise 2.13 in GPCA textbook

- **Theorem.** The ML estimates for the parameters of the PPCA model \mathbf{b} , \mathbf{W} , and σ can be obtained from the ML estimates of the mean and covariance of the data, μ_N and Σ_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 Σ_N , Λ_1 is the matrix with the corresponding top d eigenvalues, $R \in \mathbb{R}^{d \times d}$ is an arbitrary orthogonal matrix, and λ_i is the i th largest eigenvalue of Σ_N .

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 192×168 .

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 $y_i = -1 : \frac{1}{3} : 1, i = 1, 2$.