

# Manifold Learning for Signal and Visual Processing

## Lecture 9: Probabilistic PCA (PPCA), Factor Analysis, Mixtures of PPCA

Radu Horaud

INRIA Grenoble Rhone-Alpes, France

Radu.Horaud@inria.fr

<http://perception.inrialpes.fr/>

# Outline of This Lecture

- A short reminder from Lecture 1
- Probabilistic formulation of PCA (PPCA)
- Maximum-likelihood PPCA
- EM for PPCA
- Mixture of PPCA
- What is Bayesian PCA?
- Factor Analysis

# Material for This Lecture

- C. M. Bishop. Pattern Recognition and Machine Learning. 2006. (Chapter 12)
- More involved readings:
  - S. Roweis. EM algorithms of PCA and SPCA. NIPS 1998.
  - M. E. Tipping and C. M. Bishop. Probabilistic Principal Component Analysis. J. R. Stat. Soc. B. 1999.
  - M. E. Tipping and C. M. Bishop. Mixtures of Probabilistic Principal Component Analysers. Neural Computation. 1999.

## PCA at a Glance

- The input (observation) space (the data are centered):  
 $\mathbf{X} = [\mathbf{x}_1 \dots \mathbf{x}_n \dots \mathbf{x}_N]$ ,  $\mathbf{x}_n \in \mathbb{R}^D$
- The output (latent) space:  $\mathbf{Y} = [\mathbf{y}_1 \dots \mathbf{y}_n \dots \mathbf{y}_N]$ ,  $\mathbf{y}_j \in \mathbb{R}^d$
- **Projection:**  $\mathbf{Y} = \mathbf{W}^\top \mathbf{X}$  with  $\mathbf{W}^\top$  a  $d \times D$  matrix.
- **Reconstruction:**  $\mathbf{X} = \mathbf{W} \mathbf{Y}$  with  $\mathbf{W}$  a  $D \times d$  matrix.
- $\mathbf{W}^\top \mathbf{W} = \mathbf{I}_d$ , i.e.,  $\mathbf{W}^\top$  is a row-orthonormal matrix when both data sets  $\mathbf{X}$  and  $\mathbf{Y}$  are represented in orthonormal bases:  $\mathbf{y}_j = \tilde{\mathbf{U}}^\top (\mathbf{x}_j - \bar{\mathbf{x}})$ . In this case  $\mathbf{W}^\top = \tilde{\mathbf{U}}^\top$ .
- $\mathbf{W}^\top \mathbf{W} = \mathbf{\Lambda}_d^{-1}$ , i.e., this corresponds to the case of *whitening*:  $\mathbf{y}_j = \mathbf{\Lambda}_d^{-1/2} \tilde{\mathbf{U}}^\top (\mathbf{x}_j - \bar{\mathbf{x}})$ , with  $\mathbf{W}^\top = \mathbf{\Lambda}_d^{-1/2} \tilde{\mathbf{U}}^\top$ .
- Remember that  $\mathbf{W}^\top$  was estimated from the  $d$  largest eigenvalue-eigenvector pairs of the data covariance matrix.

# From Lecture #1: Data Projection on a Linear Subspace

- From  $\mathbf{Y} = \mathbf{W}^\top \mathbf{X}$  we have

$$\mathbf{Y}\mathbf{Y}^\top = \mathbf{W}^\top \mathbf{X}\mathbf{X}^\top \mathbf{W} = \mathbf{W}^\top \tilde{\mathbf{U}}\mathbf{\Lambda}_d\tilde{\mathbf{U}}^\top \mathbf{W}$$

- 1 The projected data has a diagonal covariance matrix:  $\mathbf{Y}\mathbf{Y}^\top = \mathbf{\Lambda}_d$ , by identification we obtain

$$\mathbf{W}^\top = \tilde{\mathbf{U}}^\top$$

- 2 The projected data has an identity covariance matrix, this is called *whitening the data*:  $\mathbf{Y}\mathbf{Y}^\top = \mathbf{I}_d$

$$\mathbf{W}^\top = \mathbf{\Lambda}_d^{-\frac{1}{2}} \tilde{\mathbf{U}}^\top$$

- In what follow, we will consider  $\mathbf{W}$  (reconstruction) instead of  $\mathbf{W}^\top$  (projection).

# The Probabilistic Framework (I)

- Consider again the *reconstruction* of the observed variables from the latent variables. A point  $x$  is *reconstructed* from  $y$  with:

$$x - \mu = \mathbf{W}y + \varepsilon$$

- $\varepsilon \in \mathbb{R}^D$  is the reconstruction error and let's suppose that it has a Gaussian distribution with zero mean and spherical covariance:

$$p(\varepsilon) = \mathcal{N}(\varepsilon|0, \sigma^2\mathbf{I})$$

## The Probabilistic Framework (II)

- We can now define the conditional distribution of the observed variable  $\mathbf{x}$ , conditioned on the value of the latent variable  $\mathbf{y}$ :

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{y} + \boldsymbol{\mu}, \sigma^2\mathbf{I})$$

- The prior distribution of the latent variable is a Gaussian with zero-mean and unit-covariance:

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

- The marginal (or *predictive*) distribution  $p(\mathbf{x})$  can be obtained from the sum and product rules, supposing continuous latent variables:

$$p(\mathbf{x}) = \int_{\mathbf{y}} p(\mathbf{x}, \mathbf{y})d\mathbf{y} = \int_{\mathbf{y}} p(\mathbf{x}|\mathbf{y})p(\mathbf{y})d\mathbf{y}$$

## The Probabilistic Framework (III)

- This is an instance of the *linear-Gaussian* model, hence it is Gaussian as well:

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \mathbf{C}) \propto \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{C}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

- The posterior distribution can be obtained using the Bayes' theorem for Gaussian variables (see Bishop'06, chapter 2):

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{M}^{-1}\mathbf{W}^\top(\mathbf{x} - \boldsymbol{\mu}), \sigma^{-2}\mathbf{M})$$

- This is the main difference with standard PCA: *the latent variable is in this case a random variable with a Gaussian distribution.*

## The Probabilistic Framework (IV)

- The mean and covariance of this *predictive distribution* can be formally derived from the expression of  $\mathbf{x}$  and from the Gaussian distributions just defined (using the fact that  $\mathbf{y}$  and  $\varepsilon$  are independent random variables):

$$\begin{aligned} E[\mathbf{x}] &= E[\mathbf{W}\mathbf{y} + \boldsymbol{\mu} + \varepsilon] = \mathbf{W}E[\mathbf{y}] + E[\boldsymbol{\mu}] + E[\varepsilon] = \boldsymbol{\mu} \\ \mathbf{C} &= E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top] = E[(\mathbf{W}\mathbf{y} + \varepsilon)(\mathbf{W}\mathbf{y} + \varepsilon)^\top] \\ &= \mathbf{W}E[\mathbf{y}\mathbf{y}^\top]\mathbf{W}^\top + E[\varepsilon\varepsilon^\top] = \mathbf{W}\mathbf{W}^\top + \sigma^2\mathbf{I} \end{aligned}$$

- Gaussian distributions require the inverse of the covariance matrix; Using the *Woodbury identity* (see equation (C.7) in Bishop'06) we have:

$$\begin{aligned} \mathbf{C}^{-1} &= \sigma^{-2}\mathbf{I} - \sigma^{-2}\mathbf{W}\mathbf{M}^{-1}\mathbf{W}^\top \\ \mathbf{M} &= \mathbf{W}^\top\mathbf{W} + \sigma^2\mathbf{I} \end{aligned}$$

where  $\mathbf{M}$  is a  $d \times d$  matrix. *Useful when  $d \ll D$ .*

# Maximum-likelihood PCA (I)

- The observed-data log-likelihood writes:

$$\ln p(\mathbf{x}_1, \dots, \mathbf{x}_N | \boldsymbol{\mu}, \mathbf{W}, \sigma^2) = \sum_{j=1}^N \ln p(\mathbf{x}_j | \boldsymbol{\mu}, \mathbf{W}, \sigma^2)$$

- This expression can be developed using the previous equations, to obtain:

$$\ln P(\mathbf{X} | \boldsymbol{\mu}, \mathbf{C}) = -\frac{N}{2} ((D \ln(2\pi)) + \ln |\mathbf{C}|) - \frac{1}{2} \sum_{j=1}^N (\mathbf{x}_j - \boldsymbol{\mu})^\top \mathbf{C}^{-1} (\mathbf{x}_j - \boldsymbol{\mu})$$

## Maximum-likelihood PCA (II)

- The log-likelihood is quadratic in  $\boldsymbol{\mu}$ , by setting the derivative with respect to  $\boldsymbol{\mu}$  equal to zero, we obtain the expected result:

$$\boldsymbol{\mu}_{ML} = \sum_{j=1}^N \mathbf{x}_j = \bar{\mathbf{x}}$$

- Maximization with respect to  $\mathbf{W}$  and  $\sigma^2$ , while is more complex, still has a closed-form solution:

$$\begin{aligned}\mathbf{W}_{ML} &= \tilde{\mathbf{U}}(\boldsymbol{\Lambda}_d - \sigma^2 \mathbf{I}_d)^{1/2} \mathbf{R} \\ \sigma_{ML}^2 &= \frac{1}{D-d} \sum_{i=d+1}^D \lambda_i\end{aligned}$$

- With  $\boldsymbol{\Sigma}_X = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^\top \approx \tilde{\mathbf{U}}\boldsymbol{\Lambda}_d\tilde{\mathbf{U}}^\top$ ,  $d < D$ , and  $\mathbf{R}\mathbf{R}^\top = \mathbf{I}$  (a  $d \times d$  matrix).

## Maximum-likelihood PCA (Discussion)

- The covariance of the predictive density,  $\mathbf{C} = \mathbf{W}\mathbf{W}^\top + \sigma^2\mathbf{I}$ , is not affected by the arbitrary orthogonal transformation  $\mathbf{R}$  of the latent space:

$$\mathbf{C} = \tilde{\mathbf{U}}\mathbf{D}[\lambda_i - \sigma^2]\tilde{\mathbf{U}}^\top + \sigma^2\mathbf{I}$$

- The covariance projected onto a unit vector  $\mathbf{v}$  is  $\mathbf{v}^\top\mathbf{C}\mathbf{v}$ . We obtain the following cases:
  - $\mathbf{v}$  is orthogonal to  $\tilde{\mathbf{U}}$ , then  $\mathbf{v}^\top\mathbf{C}\mathbf{v} = \sigma^2$  (noise variance) or the average variance associated with the discarded dimensions.
  - $\mathbf{v} = \mathbf{u}_i$  is one of the column vectors of  $\tilde{\mathbf{U}}$ , then  $\mathbf{v}^\top\mathbf{C}\mathbf{v} = \lambda_i - \sigma^2 + \sigma^2 = \lambda_i$ : the model correctly captures the variance along the principal directions and approximates the variance in the remaining directions with  $\sigma^2$ .
- Matrix  $\mathbf{R}$  introduces an arbitrary orthogonal transformation of the latent space.

# Projecting the Data onto the Latent Space

- Any data point  $\mathbf{x}$  can be summarized by its *posterior mean* and *posterior covariance* in latent space. These are provided by the posterior distribution  $p(\mathbf{y}|\mathbf{x})$ :

$$E[\mathbf{y}|\mathbf{x}] = \mathbf{M}^{-1}\mathbf{W}^\top(\mathbf{x} - \boldsymbol{\mu})$$

$$\mathbf{C}(\mathbf{y}|\mathbf{x}) = \sigma^{-2}\mathbf{M}$$

$$\mathbf{M} = \mathbf{W}^\top\mathbf{W} + \sigma^2\mathbf{I}$$

## From Probabilistic to Standard PCA

- The maximum-likelihood solution allows to estimate the *reconstruction* matrix  $\mathbf{W}$  and the variance  $\sigma$ . The *projection* of the data onto the latent space can be estimated from the posterior mean. We obtain the following projection matrix:

$$(\mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I})^{-1} \mathbf{W}^\top$$

- When  $\sigma^2 = 0$  this corresponds to the standard PCA solution – rotating, projecting and whitening the data:

$$(\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top = \mathbf{\Lambda}^{-1/2} \tilde{\mathbf{U}}^\top$$

# EM for PCA

- We can derive an EM algorithm for PCA, by following the EM framework: derive the complete-data log-likelihood conditioned by the observed data, and take its expectation.
- Complete-data log-likelihood for observed-latent pairs  $\mathbf{x}_j, \mathbf{y}_j$ :

$$\ln P(\mathbf{X}, \mathbf{Y} | \boldsymbol{\mu}, \mathbf{W}, \sigma^2) = \sum_{j=1}^n (\ln P(\mathbf{x}_j | \mathbf{y}_j) + \ln P(\mathbf{y}_j))$$

- Then we take the expectation with respect to the posterior distribution over the latent variables,  $E[\ln P(\mathbf{X}, \mathbf{Y} | \boldsymbol{\mu}, \mathbf{W}, \sigma^2)]$ , which depends on the current model parameters  $\boldsymbol{\mu} = \bar{\mathbf{x}}$ ,  $\mathbf{W}$ , and  $\sigma^2$ , as well as on (these are the posterior statistics):

$$\begin{aligned} E[\mathbf{y}_j] &= \mathbf{M}^{-1} \mathbf{W}^\top (\mathbf{x}_j - \bar{\mathbf{x}}) \\ E[\mathbf{y}_j \mathbf{y}_j^\top] &= \sigma^2 \mathbf{M}^{-1} + E[\mathbf{y}_j] E[\mathbf{y}_j]^\top \end{aligned}$$

# The EM Algorithm (I)

- *Initialize* the parameter values  $\mathbf{W}$  and  $\sigma^2$ .
- *E-step*: Estimate the posterior statistics  $E[\mathbf{y}_j]$  and  $E[\mathbf{y}_j\mathbf{y}_j^\top]$  using the current parameter values.
- *M-step*: Maximize with respect to  $\mathbf{W}$  and  $\sigma^2$  while keeping the posterior statistics fixed. The equations are:

$$\mathbf{W}_{new} = \left( \sum_{j=1}^N (\mathbf{x}_j - \bar{\mathbf{x}}) E[\mathbf{y}_j]^\top \right) \left( \sum_{j=1}^N E[\mathbf{y}_j\mathbf{y}_j^\top] \right)^{-1}$$

$$\sigma_{new}^2 = \frac{1}{ND} \sum_{j=1}^N \left( \|\mathbf{x}_j - \bar{\mathbf{x}}\|^2 - 2E[\mathbf{y}_j]^\top \mathbf{W}_{new}^\top (\mathbf{x}_j - \bar{\mathbf{x}}) \right. \\ \left. + \text{tr}(E[\mathbf{y}_j\mathbf{y}_j^\top] \mathbf{W}_{new}^\top \mathbf{W}_{new}) \right)$$

## The EM Algorithm (II)

- By substitution of  $E[\mathbf{y}_j]$  and  $E[\mathbf{y}_j \mathbf{y}_j^\top]$  (the E-step) into the expressions of  $\mathbf{W}_{new}$  and  $\sigma_{new}^2$  (the M-step), we get:

$$\mathbf{W}_{new} = \Sigma_X \mathbf{W}_{old} (\sigma_{old}^2 \mathbf{I} + \mathbf{M}_{old}^{-1} \mathbf{W}_{old}^\top \Sigma_X \mathbf{W}_{old})^{-1}$$
$$\sigma_{new}^2 = \frac{1}{D} \text{tr}(\Sigma_X - \Sigma_X \mathbf{W}_{old} \mathbf{M}_{old}^{-1} \mathbf{W}_{old}^\top)$$

## EM for PCA (Discussion)

- Computational efficiency for high-dimensional spaces. EM is iterative, but each iteration can be quite efficient. The covariance matrix is never estimated explicitly.
- The case of  $\sigma^2 = 0$  corresponds to a valid EM algorithm: *S. Roweis. EM algorithms of PCA and SPCA. NIPS 1998.*
- More details can be found in *M. E. Tipping and C. M. Bishop. Probabilistic Principal Component Analysis. J. R. Stat. Soc. B. 1999*

# Mixture of PPCA

- The log-likelihood of a mixture of PPCA:

$$\sum_{i=1}^N \ln(p(\mathbf{x}_i)) = \sum_{i=1}^N \ln \left( \sum_{j=1}^M \pi_j p(\mathbf{x}_i|j) \right)$$

- We seek  $\boldsymbol{\mu}_j$ ,  $\mathbf{W}_j$ , and  $\sigma_j^2$  for each mixture component  $j$ .
- For a given data point  $\mathbf{x}$  there is a posterior distribution associated with each latent space  $j$ , the mean of which is  $\mathbf{M}_j^{-1} \mathbf{W}_j^\top (\mathbf{x} - \boldsymbol{\mu}_j)$ .
- It is also possible to define *the posterior*, or responsibility, of mixture component  $j$  for generating a data point:

$$r_{ij} = \frac{p(\mathbf{x}_i|j)\pi_j}{p(\mathbf{x}_i)}$$

# EM for Mixtures of PPCA

- Initialization of the model parameters
- E-step: estimate the posteriors  $r_{ij}$
- M-step: Use the maximum-likelihood formulation of PPCA to estimate  $\mathbf{W}_j$ , and  $\sigma_j^2$  from the *local responsibility-weighted* covariance matrix:

$$\Sigma_j = \frac{1}{\pi_j N} \sum_{i=1}^N r_{ij} (\mathbf{x}_i - \boldsymbol{\mu}_j)(\mathbf{x}_i - \boldsymbol{\mu}_j)^\top$$

with  $\pi_j = 1/N \sum_{i=1}^N r_{ij}$  and  $\boldsymbol{\mu}_j = \frac{\sum_{i=1}^N r_{ij} \mathbf{x}_i}{\sum_{i=1}^N r_{ij}}$

## Bayesian PCA (I)

- Select the dimension  $d$  of the latent space.
- The generative model just introduced (well defined likelihood function) allows to address the problem in a principled way.
- The idea is to consider each column in  $\mathbf{W}$  as having an independent Gaussian prior:

$$P(\mathbf{W}|\boldsymbol{\alpha}) = \prod_{i=1}^d \left(\frac{\alpha_i}{2\pi}\right)^{D/2} \exp\left(-\frac{1}{2}\alpha_i \mathbf{w}_i^\top \mathbf{w}\right)$$

- where  $\alpha_i = 1/\sigma_i^2$  is called the precision parameter. The objective is to estimate these parameters, one for each principal direction, and select only a subset of these directions.
- We need to select directions of maximum variance, hence directions with *infinite precision* will be disregarded.

## Bayesian PCA (II)

- The approach is based on *evidence approximation* or *empirical Bayes*.
- The marginal likelihood function (the latent space  $\mathbf{W}$  is *integrated out*):

$$P(\mathbf{X}|\boldsymbol{\alpha}, \boldsymbol{\mu}, \sigma_2) = \int \underbrace{P(\mathbf{X}|\boldsymbol{\mu}, \mathbf{W}, \sigma_2)}_{\text{ML PCA}} P(\mathbf{W}|\boldsymbol{\alpha}) d\mathbf{W}$$

- The formal derivation is quite involved. The maximization with respect to the precision parameters yields a simple form:

$$\alpha_i^{new} = \frac{D}{\mathbf{w}_i^\top \mathbf{w}}$$

- This estimation is interleaved with the EM updates for estimating  $\mathbf{W}$  and  $\sigma^2$ .

# Factor Analysis

- Probabilistic PCA so far (the predictive covariance is isotropic):

$$P(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{y} + \boldsymbol{\mu}, \sigma^2\mathbf{I})$$

- In factor analysis, the covariance is diagonal rather than isotropic:

$$P(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{y} + \boldsymbol{\mu}, \boldsymbol{\Psi})$$

- the columns of  $\mathbf{W}$  are called *factor loadings* and the diagonal entries of  $\boldsymbol{\Psi}$  are called *uniquenesses*.
- The factor analysis point of view: one form of latent-variable density model, the form of the latent space is of interest but not the particular choice of coordinates (up to an orthogonal transformation).
- The factor analysis parameters,  $\mathbf{W}$ , and  $\boldsymbol{\Psi}$  are estimated via the maximum likelihood and EM frameworks.