



Machine Learning 1

Lecture 10.3 - Unsupervised Learning
Probabilistic Principal Component Analysis

Erik Bekkers

(Bishop 12.2.1)



Probabilistic PCA

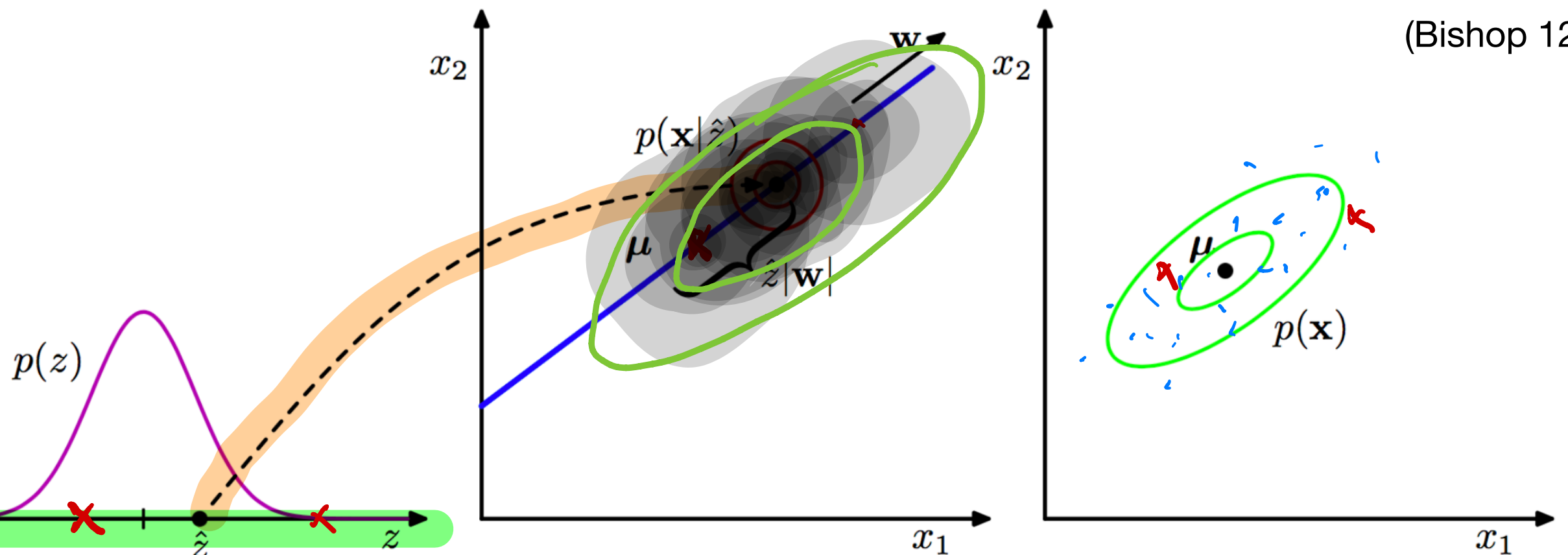
- ▶ Probabilistic view of PCA:
 - ▶ Learn it via maximum likelihood
 - ▶ (Third) alternative view of PCA
 - ▶ Both latent and observed variables are Gaussian

Continuous latent variable model

- ▶ Data: $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}, \mathbf{x}_n \in \mathbb{R}^D$
- ▶ Goal: learn a $M < D$ **continuous latent space** by maximizing the likelihood of the probabilistic model, M given
- ▶ Recall the continuous latent variable model

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

(Bishop 12.9)



PPCA modeling assumptions

- ▶ The **generative model** works as follows

$$\mathbf{x} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\epsilon}$$

- ▶ With $\boldsymbol{\mu} \in \mathbb{R}^D$ and the **continuous latent variable** $\mathbf{z} \in \mathbb{R}^M$ with Gaussian **prior**

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

- ▶ Matrix $\mathbf{W} \in \mathbb{R}^{D \times M}$ transforms the latent variables into observed variables.
- ▶ The independent noise is also Gaussian

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

It follows...

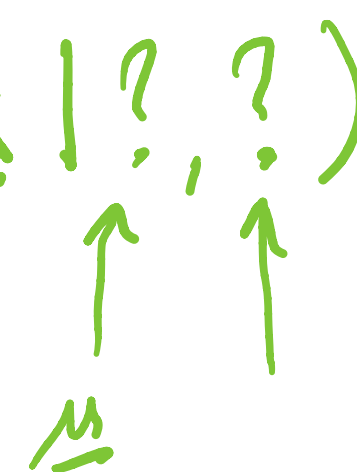
- ▶ The conditional distribution of the observed variable

Gaussian!

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

- ▶ The marginal $p(\mathbf{x})$ is also a Gaussian (Bishop Ch. 2.3)

$$p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \int p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) d\mathbf{z} = \mathcal{N}(\mathbf{x} | \mu, \Sigma)$$



- ▶ The expected value for $\mathbf{x} \sim p(\mathbf{x})$ is given by

$$\begin{aligned} \mathbb{E}[\mathbf{x}] &= \mathbb{E}[\mathbf{W} \mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\epsilon}] \\ &= \mathbf{W} \mathbb{E}[\mathbf{z}] + \boldsymbol{\mu} + \mathbb{E}[\boldsymbol{\epsilon}] \\ &= \boldsymbol{\mu} \quad \text{"0} \quad \text{"0} \end{aligned}$$

It follows...

- ▶ The covariance of $\mathbf{x} \sim p(\mathbf{x})$ is given by

$$\text{Cov}[\mathbf{x}] = \mathbb{E}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T]$$

$$\underline{x} = W \underline{z} + \mu + \epsilon$$

$$= \mathbb{E}[(W \mathbf{z} + \boldsymbol{\epsilon})(W \mathbf{z} + \boldsymbol{\epsilon})^T]$$

expand

$$= \mathbb{E}[W \mathbf{z} \mathbf{z}^T W^T + 2W \mathbf{z} \boldsymbol{\epsilon}^T + \boldsymbol{\epsilon} \boldsymbol{\epsilon}^T]$$

linearity

$$= W \mathbb{E}[\mathbf{z} \mathbf{z}^T] W^T + 2W \mathbb{E}[\mathbf{z} \boldsymbol{\epsilon}^T] + \mathbb{E}[\boldsymbol{\epsilon} \boldsymbol{\epsilon}^T]$$

$$= W W^T + \sigma^2 \mathbf{I} = \mathbf{C}$$

$\mathbf{z}, \boldsymbol{\epsilon}$ independent so

$$\text{Cov}[\mathbf{z}, \boldsymbol{\epsilon}] = \mathbb{E}[\mathbf{z} \boldsymbol{\epsilon}^T] - \mathbb{E}[\mathbf{z}] \mathbb{E}[\boldsymbol{\epsilon}]^T = \mathbf{0}$$

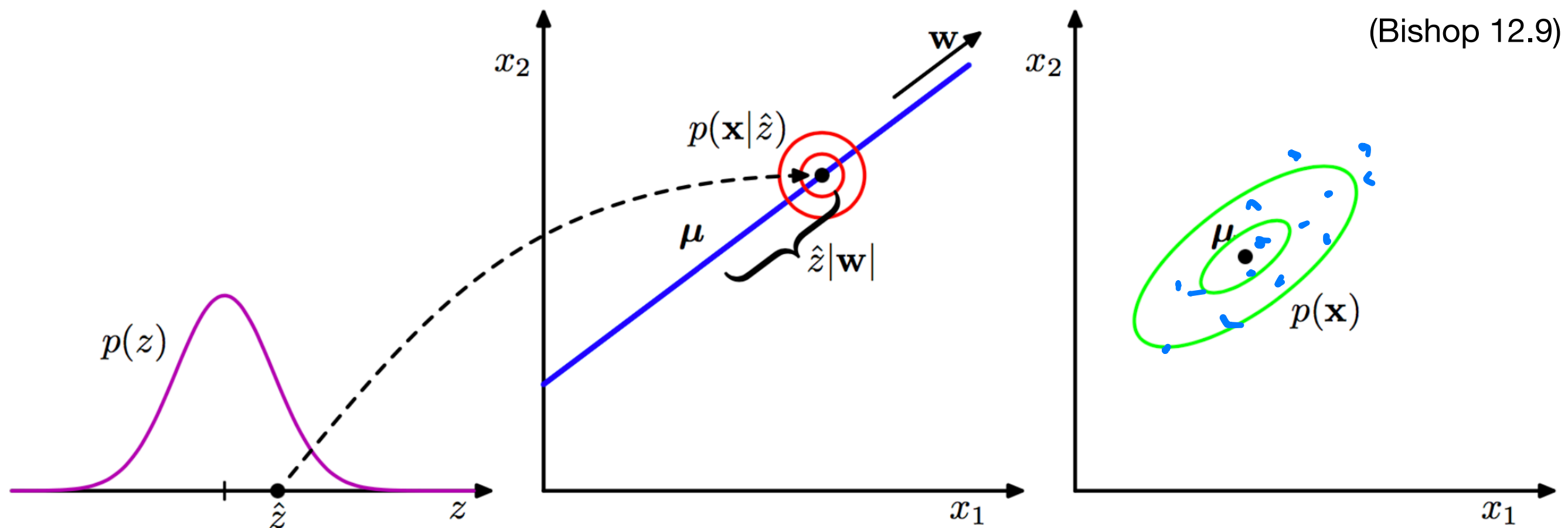
$$\Rightarrow \mathbb{E}[\mathbf{z} \boldsymbol{\epsilon}^T] = \mathbf{0}$$

$$\begin{aligned} \text{Cov}[\boldsymbol{\epsilon}] &= \mathbb{E}[\boldsymbol{\epsilon} \boldsymbol{\epsilon}^T] \\ &\quad + \cancel{\mathbb{E}[\boldsymbol{\epsilon}] \mathbb{E}[\boldsymbol{\epsilon}]^T} \\ &= \sigma^2 \mathbf{I} \end{aligned}$$

- ▶ Therefore

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \mathbf{C})$$

Probabilistic PCA in a picture



- ▶ Prior $p(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I})$
- ▶ Likelihood $p(\mathbf{x} | \mathbf{z}) = \mathcal{N}(\mathbf{x} | \mathbf{W} \mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I})$
- ▶ Marginal $p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \mathbf{C})$

The log-likelihood

$$\begin{aligned}\ln p(\mathbf{X} | \boldsymbol{\mu}, \mathbf{W}, \sigma^2) &= \sum_{n=1}^N \ln p(\mathbf{x}_n | \boldsymbol{\mu}, \mathbf{W}, \sigma^2) = \sum_{n=1}^N \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}, \mathbf{C}) \\ &= -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\mathbf{C}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{x}_n - \boldsymbol{\mu})\end{aligned}$$

- Once again, take it's derivative w.r.t. the parameter of interest, set it to zero, and solve it. For the mean:

$$\sum_{n=1}^N \mathbf{C}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) = 0 \quad \Rightarrow \quad \boldsymbol{\mu} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$$

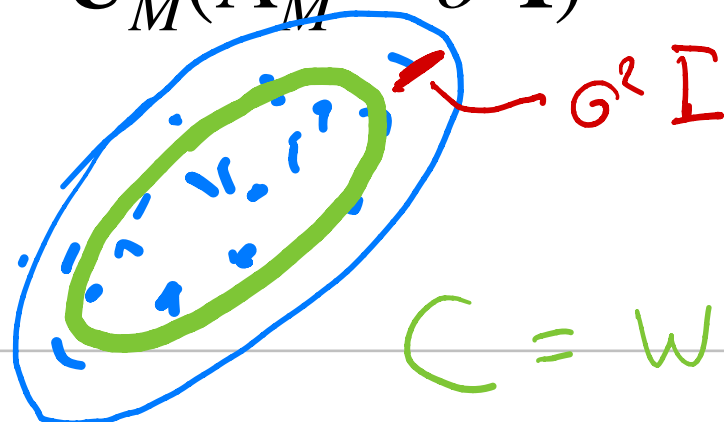
PPCA has closed-form solutions

- Let \mathbf{S} be the sample covariance, as defined in PCA, and let its eigenvectors and eigenvalues be given
- The optimal parameters for the maximal log-likelihood are

$$\boldsymbol{\mu} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$$

$$\sigma^2 = \frac{1}{D - M} \sum_{j=M+1}^D \lambda_j$$

$$\mathbf{W} = \mathbf{U}_M (\boldsymbol{\Lambda}_M - \sigma^2 \mathbf{I})^{1/2}$$



$$\mathbf{C} = \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}$$

Assume $\mathbf{S} = \mathbf{S}_{\text{true}} + \mathbf{S}_{\text{noise}}$

$$\mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^T = \mathbf{U}_M \boldsymbol{\Lambda}_M \mathbf{U}_M^T + \mathbf{U}_- \boldsymbol{\Lambda}_- \mathbf{U}_-^T$$

Covariance of Gaussian was

$$\mathbf{C} = \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}$$

\mathbf{C} should match \mathbf{S}_{true}

$$\mathbf{U}_M \boldsymbol{\Lambda}_M \mathbf{U}_M^T = \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}$$

$$\mathbf{W}\mathbf{W}^T = \mathbf{U}_M (\boldsymbol{\Lambda}_M + \sigma^2 \mathbf{I}) \mathbf{U}_M^T$$

PPCA

- ▶ PPCA is the probabilistic **generative** version of PCA: we can also draw samples from it
- ▶ PPCA is a form of Gaussian distribution with number of parameters **restricted** (by the latent space)
- ▶ PPCA is the basis of **Bayesian PCA** (Bishop 12.2.3) in which the dimension of the latent space can be found from the data
- ▶ Like with the Gaussian Mixture Model, also the PPCA can be done via an **EM algorithm** (Bishop 12.2.2) - *Though not necessary because we have close form solutions*

PCA: Summary

- **Three views**

- Max variance, min reconstruction error, probabilistic

- **Applications**

- Dimensionality reduction
 - 2D/3D visualization
 - Compression
 - Whitening (de-correlating features)
 - (not mentioned) De-noising: discard the smallest variance features = the noise components (hopefully!)

- **Limitation**

- Only linear transformations