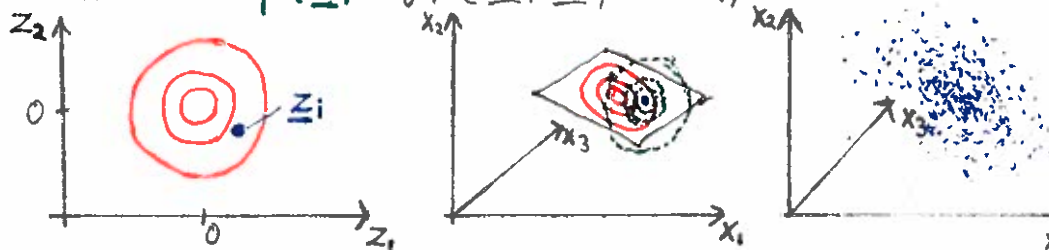


# Probabilistic PCA

- Latent  $\underline{z}$  with  $p(\underline{z}) = \mathcal{N}(\underline{z} | \underline{0}, I_k)$

- Observable  $\underline{x}$  with  $\underline{x} = W\underline{z} + \underline{\mu} + \underline{\varepsilon}$

where  $p(\underline{\varepsilon}) = \mathcal{N}(\underline{\varepsilon} | \underline{0}, \sigma^2 I_d)$



Joint distribution:

$$p(\underline{z}, \underline{x}) = \frac{1}{\text{const}} \exp\left(-\frac{1}{2}[(\underline{x} - W\underline{z} - \underline{\mu})^T (\frac{1}{\sigma^2} I_d) \cdot (\underline{x} - W\underline{z} - \underline{\mu}) + \underline{z}^T \underline{z}]\right)$$

Completing the square:

For a given  $-\frac{1}{2} \underline{x}^T A \underline{x} + \underline{x}^T \underline{y} + \text{const}$ :

1) Isolate 2<sup>nd</sup> order term  $-\frac{1}{2} \underline{x}^T A \underline{x}$

and obtain  $\Sigma = A^{-1}$

2) Isolate linear term  $\underline{x}^T \underline{y}$  and use  $\Sigma$

to obtain  $\underline{\mu} = \Sigma \underline{y}$

Applied to the PPCA model:

- Take  $\begin{pmatrix} \underline{z} \\ \underline{x} \end{pmatrix}$  as a single random vector and isolate

2<sup>nd</sup> order terms of  $p(\underline{z}, \underline{x})$ :

$$-\frac{1}{2} \left[ \underline{z}^T (\mathbf{I}_k + \mathbf{W}^T \frac{1}{\sigma^2} \mathbf{W}) \underline{z} + \underline{x}^T \frac{1}{\sigma^2} \underline{x} - \underline{x}^T \frac{1}{\sigma^2} \mathbf{W} \underline{z} - \underline{z}^T \mathbf{W}^T \frac{1}{\sigma^2} \underline{x} \right]$$

$$= -\frac{1}{2} \begin{pmatrix} \underline{z} \\ \underline{x} \end{pmatrix}^T \begin{pmatrix} \mathbf{I}_k + \mathbf{W}^T \frac{1}{\sigma^2} \mathbf{W} & -\mathbf{W}^T \frac{1}{\sigma^2} \\ -\frac{1}{\sigma^2} \mathbf{W} & \frac{1}{\sigma^2} \mathbf{I}_d \end{pmatrix} \begin{pmatrix} \underline{z} \\ \underline{x} \end{pmatrix}$$

$$\Rightarrow \text{Cov} \left[ \begin{pmatrix} \underline{z} \\ \underline{x} \end{pmatrix} \right] = \left( \begin{array}{cc} & \\ & \end{array} \right)^{-1} \stackrel{\text{lin}}{\text{alg}} \left( \begin{array}{cc} \mathbf{I}_k & \mathbf{W}^T \\ \mathbf{W} & \mathbf{W} \mathbf{W}^T + \sigma^2 \mathbf{I}_d \end{array} \right)$$

- Linear terms:

$$-\underline{z}^T \mathbf{W}^T \frac{1}{\sigma^2} \underline{\mu} + \underline{x}^T \frac{1}{\sigma^2} \underline{\mu} = \begin{pmatrix} \underline{z} \\ \underline{x} \end{pmatrix}^T \begin{pmatrix} -\mathbf{W}^T \frac{1}{\sigma^2} \underline{\mu} \\ \frac{1}{\sigma^2} \underline{\mu} \end{pmatrix}$$

$$\rightarrow \mathbb{E} \left[ \begin{pmatrix} \underline{z} \\ \underline{x} \end{pmatrix} \right] = \begin{pmatrix} \mathbf{I}_k & \mathbf{W}^T \\ \mathbf{W} & \mathbf{W} \mathbf{W}^T + \sigma^2 \mathbf{I}_d \end{pmatrix} \begin{pmatrix} -\mathbf{W}^T \frac{1}{\sigma^2} \underline{\mu} \\ \frac{1}{\sigma^2} \underline{\mu} \end{pmatrix} = \begin{pmatrix} \underline{0} \\ \underline{\mu} \end{pmatrix}$$

$$\Rightarrow p(\underline{z}, \underline{x}) = \mathcal{N} \left( \begin{pmatrix} \underline{z} \\ \underline{x} \end{pmatrix} \middle| \begin{pmatrix} \underline{0} \\ \underline{\mu} \end{pmatrix}, \begin{pmatrix} \mathbf{I}_k & \mathbf{W}^T \\ \mathbf{W} & \mathbf{W} \mathbf{W}^T + \sigma^2 \mathbf{I}_d \end{pmatrix} \right)$$

Using the lower (right) partitions, we also obtain the observation distr.:

$$p(\underline{x}) = \mathcal{N}(\underline{x} \mid \underline{\mu}, \mathbf{W} \mathbf{W}^T + \sigma^2 \mathbf{I}_d)$$

# Maximum likelihood

Given centred  $d \times n$  observation matrix

$$X = (\underline{x}_1 \dots \underline{x}_n)$$

$$\text{Objective: } \log p(X|W, \sigma^2) \stackrel{\text{ind}}{=} \sum_{i=1}^n \log p(\underline{x}_i | W, \sigma^2) \\ = \sum_{i=1}^n \log ((2\pi)^d \det \Sigma)^{-\frac{1}{2}} - \frac{1}{2} \underline{x}_i^T \Sigma^{-1} \underline{x}_i$$

$$= -\frac{1}{2} [n d \log 2\pi + n \log (\det \Sigma) + n \frac{1}{n} \sum_i \underline{x}_i^T \Sigma^{-1} \underline{x}_i]$$

$$= -\frac{n}{2} [d \log (2\pi) + \log (\det \Sigma) + \text{trace}(\Sigma^{-1} \hat{\Sigma})]$$

where  $\Sigma = WW^T + \sigma^2 I_d$  and

$$\hat{\Sigma} = \text{cov}(X) = \frac{1}{n} XX^T$$

$$W_{ML} = U_k (\Lambda_k - \sigma^2 I_k)^{1/2} R \quad (\text{Tipping and Bishop 1999})$$

where  $\hat{\Sigma} = U \Lambda U$  the EV decomposition

$R = k \times k$  orth matrix "Solution not unique"

e.g.  $R = I_k$  yielding  $W_{ML} = U_k (\Lambda_k - \sigma^2 I_k)^{1/2}$

$$\sigma_{ML}^2 = \frac{1}{d-k} \sum_{i=k+1}^d \lambda_i$$

$\leadsto \sigma^2$  represents the average lost variance per residual dimension

## Relation to PCA:

Given an observation  $\underline{x}_i$ , what do we know about the latent variable?

A)  $\underline{z}_i = U_k^T \underline{x}_i$  X

B) We get a single  $\underline{z}_i$  from  $\underline{x}_i$  via some other relation X

C) We get a distribution over  $\underline{z}$ 's ✓

D) You don't know X

Relation to PCA:

- PCA: map data  $\underline{x}_i$  to PC scores  $\underline{z}_i$ :

$$\underline{z}_i = U_k^T \underline{x}_i$$

- Probabilistic PCA: map  $\underline{z}_i$  to  $\underline{x}_i$ :

$$\underline{x}_i = W \underline{z}_i + \underline{\epsilon} \text{ stochastic}$$

→ whole posterior distribution  $p(\underline{z} | \underline{x}_i)$   
closest thing to PCA mapping

Posterior  $p(\underline{z} | \underline{x})$ :

Fix  $\underline{x}$  in  $p(\underline{z}, \underline{x})$  and complete  
the square

Completing the square for posterior:

- 2<sup>nd</sup> order terms:

$$\begin{aligned}\text{Cov}[\underline{z}|\underline{x}] &= (I_k + W^T \frac{1}{\sigma^2} W)^{-1} \\ &= \sigma^2 (W^T W + \frac{1}{\sigma^2} I_k)^{-1} = \sigma^2 M^{-1}\end{aligned}$$

where  $M = W^T W + \sigma^2 I_k$

- Linear terms:

$$\begin{aligned}\underline{y} &= \frac{1}{2} W^T \frac{1}{\sigma^2} \underline{x} + \frac{1}{2} (\underline{x}^T \frac{1}{\sigma^2} W)^T - W^T \frac{1}{\sigma^2} \underline{\mu} \\ &= W^T \frac{1}{\sigma^2} (\underline{x} - \underline{\mu})\end{aligned}$$

$$\begin{aligned}\Rightarrow E[\underline{z}|\underline{x}] &= \sigma^2 M^{-1} W^T \frac{1}{\sigma^2} (\underline{x} - \underline{\mu}) \\ &= M^{-1} W^T (\underline{x} - \underline{\mu})\end{aligned}$$

Posterior:  $p(\underline{z}|\underline{x}) = \mathcal{N}(\underline{z} | M^{-1} W^T (\underline{x} - \underline{\mu}), \sigma^2 M^{-1})$

Projections:

For PCA:  $\hat{\underline{x}} = U_k U_k^T \underline{x}$

For Probabilistic PCA:

$$\hat{\underline{x}} = W_{ML} \mathbb{E}[\underline{z} | \underline{x}] = W_{ML} M_{ML}^{-1} W_{ML}^T \underline{x}$$

where  $M_{ML} = W_{ML}^T W_{ML} + \sigma^2 I$

"Expectation of posterior projected back into data space."

For  $\sigma^2 \rightarrow 0$ :

$$W_{ML} M_{ML}^{-1} W_{ML}^T \underline{x}$$

$$= \dots = U_k U_k^T \underline{x}$$

$\leadsto$  PCA is a special case of Probabilistic PCA  
where covariance of noise is  
infinitesimally small