

Probabilistic dimensionality reduction

Course of Machine Learning
Master Degree in Computer Science
University of Rome “Tor Vergata”
a.a. 2024-2025

Giorgio Gambosi

1 Factor Analysis

Factor analysis is one of the simplest and most fundamental generative latent models, the first one we consider here where both the observed variable \mathbf{x} and the latent variable \mathbf{z} are real. At the same time, the model is also simple enough to make it possible to make it feasible to compute the conditional probability $p(\mathbf{z}|\mathbf{x})$, and this Hypothesis 3 not holding.

In particular, the model assume that each element $\mathbf{x}_i \in \mathbb{R}^d$ in the observable dataset is related to the value of a latent variable (also called a **factor** here) $\mathbf{z}_i \in \mathbb{R}^p$ through:

- a linear projection from the p -dimensional space \mathbb{R}^p of \mathbf{z} to the d -dimensional space \mathbb{R}^d of \mathbf{x}
- a translation of the result within \mathbb{R}^d
- an additional (smaller) random translation within \mathbb{R}^d

This is specified by the equation

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

where (see Figure 1)

- $\mathbf{z} \in \mathbb{R}^p$ is a latent variable whose distribution is assumed gaussian with 0 mean and unitary covariance matrix: hence $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$
- $\mathbf{W} \in \mathbb{R}^{d \times p}$ is a linear projection of any point in \mathbb{R}^p to a point in \mathbb{R}^d
- $\boldsymbol{\mu} \in \mathbb{R}^d$ is a translation of points in \mathbb{R}^d
- $\boldsymbol{\epsilon} \in \mathbb{R}^d$ is a gaussian noise for the final random translation: noise covariance on different dimensions is assumed to be 0. That is, its distribution is $\mathcal{N}(\boldsymbol{\epsilon}; \mathbf{0}, \boldsymbol{\Psi})$, where $\boldsymbol{\Psi} \in \mathbb{R}^{d \times d}$ is a diagonal matrix, with Ψ_{ii} the noise variance along the i -th dimension.

Background on Multivariate Gaussian Distribution

Consider two random variables $\mathbf{x} \in \mathbb{R}^d$ and $\mathbf{z} \in \mathbb{R}^p$ and let

$$\mathbf{y} = \begin{bmatrix} \mathbf{x} \\ \mathbf{z} \end{bmatrix} \in \mathbb{R}^{d+p}$$

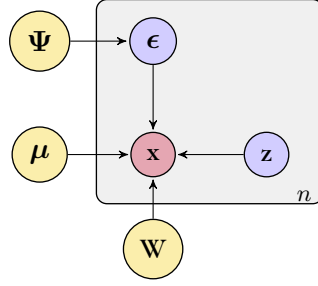


Figure 1: The latent variables ϵ and \mathbf{z} are normally distributed on the observed and the latent space, respectively: they can be both seen as random noise $p(\epsilon) = \mathcal{N}(\epsilon; \mathbf{0}, \Psi)$ and $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$. The observed variable \mathbf{x} is deterministically dependent from them as $\mathbf{x} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \epsilon$. However, a probabilistic dependence from \mathbf{z} alone can be expressed through the conditional distribution $p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{z}; \mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \Psi)$.

Assume that \mathbf{x} and \mathbf{z} are jointly multivariate Gaussian; hence, the variable \mathbf{y} has a multivariate Gaussian distribution, i.e., $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}_y, \Sigma_y)$. The mean and covariance of such distribution can be decomposed as:

$$\boldsymbol{\mu}_y = \begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_z \end{bmatrix} \in \mathbb{R}^{d+p} \quad \Sigma_y = \begin{bmatrix} \Sigma_x & \Sigma_{xz} \\ \Sigma_{zx} & \Sigma_z \end{bmatrix} \in \mathbb{R}^{(d+p) \times (d+p)},$$

where $\boldsymbol{\mu}_x \in \mathbb{R}^d$, $\boldsymbol{\mu}_y \in \mathbb{R}^p$, $\Sigma_x \in \mathbb{R}^{d \times d}$, $\Sigma_z \in \mathbb{R}^{p \times p}$, $\Sigma_{xz} \in \mathbb{R}^{d \times p}$, and $\Sigma_{zx} = \Sigma_{xz}^T \in \mathbb{R}^{p \times d}$.

It can be shown that the marginal distributions $p(\mathbf{x})$ and $p(\mathbf{z})$ are Gaussian distributions with $E[\mathbf{x}] = \boldsymbol{\mu}_x$ and $E[\mathbf{z}] = \boldsymbol{\mu}_z$. The covariance matrix of the joint distribution can be simplified as:

$$\Sigma_y = \mathbb{E}_{\mathbf{y}} \left[(\mathbf{y} - \boldsymbol{\mu}_y)(\mathbf{y} - \boldsymbol{\mu}_y)^T \right] = \begin{bmatrix} \mathbb{E}_{\mathbf{x}} \left[(\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{x} - \boldsymbol{\mu}_x)^T \right], \mathbb{E}_{\mathbf{x}, \mathbf{z}} \left[(\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{z} - \boldsymbol{\mu}_z)^T \right] \\ \mathbb{E}_{\mathbf{x}, \mathbf{z}} \left[(\mathbf{z} - \boldsymbol{\mu}_z)(\mathbf{x} - \boldsymbol{\mu}_x)^T \right], \mathbb{E}_{\mathbf{z}} \left[(\mathbf{z} - \boldsymbol{\mu}_z)(\mathbf{z} - \boldsymbol{\mu}_z)^T \right] \end{bmatrix}$$

This shows that:

$$p(\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_x, \Sigma_x),$$

$$p(\mathbf{z}) = \mathcal{N}(\boldsymbol{\mu}_z, \Sigma_z).$$

According to the definition of the multivariate Gaussian distribution, the conditional distribution is also Gaussian, i.e., $p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\boldsymbol{\mu}_{x|\mathbf{z}}, \Sigma_{x|\mathbf{z}})$ where:

$$\boldsymbol{\mu}_{x|\mathbf{z}} = \boldsymbol{\mu}_x + \Sigma_{xz}\Sigma_z^{-1}(\mathbf{z} - \boldsymbol{\mu}_z) \in \mathbb{R}^d$$

$$\Sigma_{x|\mathbf{z}} = \Sigma_x - \Sigma_{xz}\Sigma_z^{-1}\Sigma_{zx} = \Sigma_x - \Sigma_{xz}\Sigma_z^{-1}\Sigma_{zx}^T \in \mathbb{R}^{d \times d}$$

and likewise for $p(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_{z|\mathbf{x}}, \Sigma_{z|\mathbf{x}})$:

$$\boldsymbol{\mu}_{z|\mathbf{x}} = \boldsymbol{\mu}_z + \Sigma_{zx}\Sigma_x^{-1}(\mathbf{x} - \boldsymbol{\mu}_x) \in \mathbb{R}^p,$$

$$\Sigma_{z|\mathbf{x}} = \Sigma_z - \Sigma_{zx}\Sigma_x^{-1}\Sigma_{xz} = \Sigma_z - \Sigma_{zx}^T\Sigma_x^{-1}\Sigma_{xz} \in \mathbb{R}^{p \times p}.$$

All marginal and conditional distributions turn out to be Gaussian also under the following different hypotheses:

1. \mathbf{z} is normally distributed $\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}_z, \Sigma_z)$
2. there exist $\mathbf{A} \in \mathbb{R}^{d \times p}$, $\mathbf{b} \in \mathbb{R}^d$ such that the conditional distribution of \mathbf{x} given \mathbf{z} is a gaussian $p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{A}\mathbf{z} + \mathbf{b}, \Sigma_{xz})$

this is denoted as **linear gaussian model** and, in this framework, both the marginal distribution of \mathbf{x} and the inverse conditional distribution of $\mathbf{z}|\mathbf{x}$ are also Gaussian. In particular

- For the marginal distribution, $p(\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_x, \Sigma_x)$, with

$$\begin{aligned}\boldsymbol{\mu}_x &= \mathbf{A}\boldsymbol{\mu}_z + \mathbf{b} \\ \Sigma_x &= \Sigma_{xz} + \mathbf{A}\Sigma_z\mathbf{A}^T\end{aligned}$$

- For the conditional distribution, $\mathbf{z}|\mathbf{x} = \mathcal{N}(\boldsymbol{\mu}_{z|x}, \Sigma_{z|x})$, with

$$\begin{aligned}\boldsymbol{\mu}_{z|x} &= (\Sigma_z^{-1} + \mathbf{A}^T\Sigma_{xz}^{-1}\mathbf{A})^{-1}(\mathbf{A}^T\Sigma_{xz}^{-1}(\mathbf{x} - \mathbf{b}) + \Sigma_z^{-1}\boldsymbol{\mu}_x) \\ \Sigma_{z|x} &= (\Sigma_z^{-1} + \mathbf{A}^T\Sigma_{xz}^{-1}\mathbf{A})^{-1}\end{aligned}$$

The Factor Analysis Model

As already stated, the prior distribution of the latent variable is assumed to be a multivariate Gaussian distribution.

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

and the observed value \mathbf{x} is obtained from \mathbf{z} through

1. the linear projection of \mathbf{z} by $\mathbf{W} \in \mathbb{R}^{d \times p}$,
2. applying some linear translation $\boldsymbol{\mu} \in \mathbb{R}^d$, and
3. adding a Gaussian noise $\boldsymbol{\epsilon} \in \mathbb{R}^d$ with mean $\mathbf{0}$ and covariance $\boldsymbol{\Psi} \in \mathbb{R}^{d \times d}$.

As a consequence, the conditional distribution of \mathbf{x} given \mathbf{z} is

$$p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}; \mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \boldsymbol{\Psi})$$

Factor Analysis is then a linear Gaussian model with $\boldsymbol{\mu}_z = \mathbf{0}$, $\Sigma_z = \mathbf{I}$, $\mathbf{A} = \mathbf{W}$, $\mathbf{b} = \boldsymbol{\mu}$, $\Sigma_{x|z} = \boldsymbol{\Psi}$. By applying its properties, we get:

- the marginal distribution, $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^T + \boldsymbol{\Psi})$
- the inverse conditional distribution, $p(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}|\mathbf{x}; \boldsymbol{\mu}_{z|x}, \Sigma_{z|x})$, with

$$\begin{aligned}\Sigma_{z|x} &= (\mathbf{I} + \mathbf{W}^T\boldsymbol{\Psi}^{-1}\mathbf{W})^{-1} \\ \boldsymbol{\mu}_{z|x} &= \Sigma_{z|x}(\mathbf{W}^T\boldsymbol{\Psi}^{-1}(\mathbf{x} - \boldsymbol{\mu}))\end{aligned}$$

This distribution can be exploited to map points onto the latent space. In particular, the conditional expectation

$$\mathbb{E}_{p(\mathbf{z}|\mathbf{x})}[\mathbf{z}] = (\mathbf{I} + \mathbf{W}^T\boldsymbol{\Psi}^{-1}\mathbf{W})^{-1}\mathbf{W}^T\boldsymbol{\Psi}^{-1}(\mathbf{x} - \boldsymbol{\mu})$$

can be assumed as the point in latent space corresponding to \mathbf{x} .

Maximization of likelihood in FA

The log-likelihood of the observed dataset in the model is

$$\begin{aligned}\log p(\mathbf{X}|\mathbf{W}, \boldsymbol{\mu}, \boldsymbol{\Psi}) &= \sum_{i=1}^n \log p(\mathbf{x}_i|\mathbf{W}, \boldsymbol{\mu}, \boldsymbol{\Psi}) = \sum_{i=1}^n \log \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}, \boldsymbol{\Psi} + \mathbf{W}\mathbf{W}^T) \\ &= -\frac{nd}{2} \log(2\pi) - \frac{n}{2} \log |\boldsymbol{\Psi} + \mathbf{W}\mathbf{W}^T| - \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})(\boldsymbol{\Psi} + \mathbf{W}\mathbf{W}^T)^{-1} (\mathbf{x}_i - \boldsymbol{\mu})^T\end{aligned}$$

Setting the gradient wrt $\boldsymbol{\mu}$ to 0 results into

$$\boldsymbol{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$$

However, no closed form solution for \mathbf{W} and $\boldsymbol{\Psi}$ can be obtained by setting the corresponding gradients to $\mathbf{0}$. Iterative techniques such as EM can then be applied to maximize the log-likelihood with respect to these parameters.

Expectation-Maximization for FA

By definition, the algorithm operates by alternatively computing (in the E-step)

$$p(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta}^{(k)}) = \sum_{i=1}^n p(\mathbf{z}_i|\mathbf{x}_i; \boldsymbol{\theta}^{(k)})$$

given the parameter value $\boldsymbol{\theta}^{(k)}$ and then (in the M-step) maximizing

$$\mathbb{E}_{p(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta}^{(k)})} [\log p(\mathbf{X}, \mathbf{Z}; \boldsymbol{\theta})] = \sum_{i=1}^n \mathbb{E}_{p(\mathbf{z}_i|\mathbf{x}_i; \boldsymbol{\theta}^{(k)})} [\log p(\mathbf{x}_i, \mathbf{z}_i; \boldsymbol{\theta})]$$

with respect to the parameter $\boldsymbol{\theta}$, obtaining the new value $\boldsymbol{\theta}^{(k+1)}$.

M-step Let us first observe that in the case of FA, we have $\boldsymbol{\theta} = (\mathbf{W}, \boldsymbol{\mu}, \boldsymbol{\Psi})$.

For what regards maximization wrt $\boldsymbol{\mu}$, we already observed that the optimum value for such parameter is

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$$

regarding maximization wrt \mathbf{W} and $\boldsymbol{\Psi}$, we skip some technical details, stating, without proof, that

$$\begin{aligned}\bar{\mathbf{W}} &= \left(\sum_{i=1}^n \bar{\mathbf{x}}_i \boldsymbol{\mu}_{\mathbf{z}_i|\mathbf{x}_i} \right) \left(\sum_{i=1}^n \boldsymbol{\mu}_{\mathbf{z}_i|\mathbf{x}_i} \boldsymbol{\mu}_{\mathbf{z}_i|\mathbf{x}_i}^T \right)^{-1} \\ \bar{\boldsymbol{\Psi}} &= \frac{1}{n} \text{diag} \left(\mathbf{S} - \mathbf{W} \sum_{i=1}^n \boldsymbol{\mu}_{\mathbf{z}_i|\mathbf{x}_i} \bar{\mathbf{x}}_i^T \right)\end{aligned}$$

where

1. $\boldsymbol{\mu}_{\mathbf{z}_i|\mathbf{x}_i}$ and $\boldsymbol{\mu}_{\mathbf{z}_i\mathbf{z}_i^T|\mathbf{x}_i}$ are the expectations wrt distribution $p(\mathbf{z}_i|\mathbf{x}_i; \boldsymbol{\theta}^{(k)})$ of \mathbf{z}_i and $\mathbf{z}_i\mathbf{z}_i^T$, respectively

$$\begin{aligned}\boldsymbol{\mu}_{\mathbf{z}_i|\mathbf{x}_i} &\triangleq \mathbb{E}_{p(\mathbf{z}_i|\mathbf{x}_i; \boldsymbol{\theta}^{(k)})} [\mathbf{z}_i] \\ \boldsymbol{\mu}_{\mathbf{z}_i\mathbf{z}_i^T|\mathbf{x}_i} &\triangleq \mathbb{E}_{p(\mathbf{z}_i|\mathbf{x}_i; \boldsymbol{\theta}^{(k)})} [\mathbf{z}_i\mathbf{z}_i^T]\end{aligned}$$

2. $\bar{\mathbf{x}}_i$ is the difference between \mathbf{x}_i and the centroid $\bar{\mathbf{x}}$

$$\bar{\mathbf{x}}_i \triangleq \mathbf{x}_i - \bar{\mathbf{x}}$$

3. the **diag** operator sets to 0 all non diagonal elements
4. \mathbf{S} is the scatter matrix of \mathbf{X}

$$\mathbf{S} \triangleq \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$

E-step For the M-step, the conditional expectations $\boldsymbol{\mu}_{\mathbf{z}_i|\mathbf{x}_i}$ and $\boldsymbol{\mu}_{\mathbf{z}_i\mathbf{z}_i^T|\mathbf{x}_i}$ are computed in the E-step. They can be shown to be

$$\begin{aligned}\boldsymbol{\mu}_{\mathbf{z}|\mathbf{x}} &= \Sigma_{\mathbf{z}|\mathbf{x}} \Sigma_{\mathbf{z}|\mathbf{x}}^{-1} \mathbf{W}^T \boldsymbol{\Psi}^{-1} (\mathbf{x} - \bar{\mathbf{x}}) \\ \boldsymbol{\mu}_{\mathbf{z}\mathbf{z}^T|\mathbf{x}} &= \boldsymbol{\mu}_{\mathbf{z}|\mathbf{x}} \boldsymbol{\mu}_{\mathbf{z}|\mathbf{x}}^T + \Sigma_{\mathbf{z}|\mathbf{x}}\end{aligned}$$

where, as shown above,

$$\Sigma_{\mathbf{z}|\mathbf{x}} = (\mathbf{I} + \mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W})^{-1}$$

The EM algorithm in factor analysis is then summarized as follows. The centroid of data, $\bar{\mathbf{x}}$, is computed and, from it, all $\bar{\mathbf{x}}_i$. Then, at every step k , we iteratively solve as:

for $i = 1, \dots, n$:

$$\begin{aligned}\boldsymbol{\mu}_{\mathbf{z}_i|\mathbf{x}_i}^{(k)} &\leftarrow \Sigma_{\mathbf{z}|\mathbf{x}}^{(k-1)} (\mathbf{W}^{(k-1)})^T (\boldsymbol{\Psi}^{(k-1)})^{-1} \bar{\mathbf{x}}_i \\ \boldsymbol{\mu}_{\mathbf{z}_i\mathbf{z}_i^T|\mathbf{x}_i}^{(k)} &\leftarrow \boldsymbol{\mu}_{\mathbf{z}_i|\mathbf{x}_i}^{(k-1)} \boldsymbol{\mu}_{\mathbf{z}_i|\mathbf{x}_i}^{(k-1)T} + \Sigma_{\mathbf{z}|\mathbf{x}}^{(k-1)} \\ \mathbf{W}^{(k)} &\leftarrow \left(\sum_{i=1}^n \bar{\mathbf{x}}_i (\boldsymbol{\mu}_{\mathbf{z}_i|\mathbf{x}_i}^{(k)})^T \right) \left(\sum_{i=1}^n \boldsymbol{\mu}_{\mathbf{z}_i\mathbf{z}_i^T|\mathbf{x}_i}^{(k)} \right)^{-1} \\ \boldsymbol{\Psi}^{(k)} &\leftarrow \frac{1}{n} \text{diag} \left(\mathbf{S} - \mathbf{W}^{(k)} \sum_{i=1}^n \boldsymbol{\mu}_{\mathbf{z}_i|\mathbf{x}_i}^{(k)} \bar{\mathbf{x}}_i^T \right) \\ \Sigma_{\mathbf{z}|\mathbf{x}}^{(k)} &\leftarrow \left(\mathbf{I} + (\mathbf{W}^{(k)})^T (\boldsymbol{\Psi}^{(k)})^{-1} \mathbf{W}^{(k)} \right)^{-1}\end{aligned}$$

until convergence.

2 Probabilistic PCA

Probabilistic PCA is defined through a simplification of the factor analysis model. In particular, all the rest being equal, the noise covariance matrix is assumed to have equal variance for all dimensions. That is,

$$\boldsymbol{\Psi} = \sigma^2 \mathbf{I}$$

The resulting model is described graphically in Figure 2.

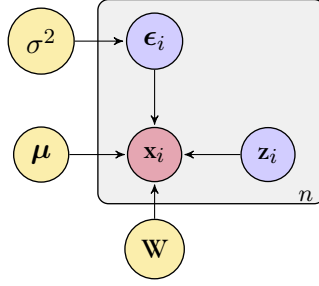


Figure 2: The latent variables ϵ and z are normally distributed on the observed and the latent space, respectively: they can be both seen as random noise $p(\epsilon; \sigma^2) = \mathcal{N}(\epsilon; \mathbf{0}, \sigma^2 \mathbf{I})$ and $p(z) = \mathcal{N}(z; \mathbf{0}, \mathbf{I})$. The observed variable \mathbf{x} is deterministically dependent from them as $\mathbf{x} = \mathbf{W}z + \mu + \epsilon$. However, a probabilistic dependence from z alone can be expressed through the conditional distribution $p(\mathbf{x}|z) = \mathcal{N}(\mathbf{x}; \mathbf{W}z + \mu, \mathbf{I}\sigma^2)$.

Expectation-Maximization for Probabilistic PCA

Clearly, Hypothesis 3 does not hold also in this case, and expectation maximization can still be applied to maximize the log-likelihood of the observed dataset \mathbf{X} wrt the parameters $\mathbf{W}, \mu, \sigma^2$.

Being PPCA a particular case of factor analysis, the E and M steps can be derived from the ones defined for FA, substituting the new noise covariance matrix $\sigma^2 \mathbf{I}$ to the more general Ψ .

This results in the following:

$$\begin{aligned}\mu_{z_i|x_i} &= \beta \Sigma_{z|x} \mathbf{W}^T \bar{\mathbf{x}}_i \\ \mu_{zz^T|x} &= \mu_{z|x} \mu_{z|x}^T + \Sigma_{z|x}\end{aligned}$$

where $\beta = \frac{1}{\sigma^2}$ is the **precision**.

It can be proved that the algorithm behaves, at each step, as follows (d is the dimensionality, that is the size of data items).

for $i = 1, \dots, n$:

$$\begin{aligned}\mu_{z_i|x_i}^{(k)} &\leftarrow \beta^{(k-1)} \Sigma_{z|x}^{(k-1)} (\mathbf{W}^{(k-1)})^T \mathbf{x}_i \\ \mu_{z_i z_i^T|x_i}^{(k)} &\leftarrow \mu_{z_i|x_i}^{(k-1)} \mu_{z_i|x_i}^{(k-1)T} + \Sigma_{z|x}^{(k-1)} \\ \mathbf{W}^{(k)} &\leftarrow \left(\sum_{i=1}^n \mathbf{x}_i (\mu_{z_i|x_i}^{(k)})^T \right) \left(\sum_{i=1}^n \mu_{z_i z_i^T|x_i}^{(k)} \right)^{-1} \\ \beta^{(k)} &\leftarrow nd \left(\sum_{i=1}^n \left(\bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T - 2 \mu_{z_i|x_i}^{(k)T} \mathbf{W}^{(k)} \bar{\mathbf{x}}_i + \text{tr} \left(\mu_{z_i z_i^T|x_i}^{(k)} (\mathbf{W}^{(k)})^T \mathbf{W}^{(k)} \right) \right) \right)^{-1}\end{aligned}$$

Maximization of the observed set log-likelihood

The probabilistic PCA model also makes it possible to analytically maximize its likelihood directly and, as a consequence, to express the linear projection of any D -dimensional point onto the d -dimensional subspace in a closed form.

The log-likelihood of the dataset in the model is

$$\begin{aligned}\log p(\mathbf{X}; \mathbf{W}, \boldsymbol{\mu}, \sigma^2) &= \sum_{i=1}^n \log p(\mathbf{x}_i; \mathbf{W}, \boldsymbol{\mu}, \sigma^2) \\ &= -\frac{nD}{2} \log(2\pi) - \frac{n}{2} \log |\Sigma_{\mathbf{x}}| - \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu}) \Sigma_{\mathbf{x}}^{-1} (\mathbf{x}_i - \boldsymbol{\mu})^T\end{aligned}$$

Maximization wrt $\boldsymbol{\mu}$ can be easily done by setting the corresponding gradient to zero, which results into

$$\boldsymbol{\mu}^* = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$$

Maximization wrt \mathbf{W} is more complex: however, a closed form solution exists:

$$\mathbf{W}^* = \mathbf{U}_d (\mathbf{L}_d - \sigma^2 \mathbf{I})^{1/2}$$

where

- \mathbf{U}_d is the $D \times d$ matrix whose columns $1, \dots, d$ are the eigenvectors corresponding to the d largest eigenvalues of the scatter matrix

$$\mathbf{S} \triangleq \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$

- \mathbf{L}_d is the $d \times d$ diagonal matrix of the largest eigenvalues $\lambda_1, \dots, \lambda_d$

The columns of \mathbf{W}^* are the eigenvectors $1, \dots, d$, each i scaled by the square root of the difference $\lambda_i - \sigma^2$.

Indeed, any rotation of \mathbf{W}^* in latent space is a solution of the likelihood maximization problem. Hence, the general solution is given by

$$\mathbf{W}^* = \mathbf{U}_d (\mathbf{L}_d - \sigma^2 \mathbf{I})^{1/2} \mathbf{R}$$

where \mathbf{R} is an arbitrary $d \times d$ orthogonal matrix, corresponding to a rotation in \mathbb{R}^d .

For what concerns the maximization wrt σ^2 , it results

$$\sigma^2 = \frac{1}{D-d} \sum_{i=d+1}^D \lambda_i$$

Since eigenvalues provide measures of the dataset variance along the corresponding eigenvector direction, this corresponds to the average variance along the discarded directions.