

Supporting information for Data driven estimation of imputation error -A strategy for imputation with a reject option

S1 Appendix

Imputation with Probabilistic PCA

We focus on latent variable models for dimensional reduction; in case so-called “probabilistic principal component analysis” (for a general discussion and references to original work see [2]). Under the probabilistic PCA model data are generated by the process

$$\bar{x} = \bar{A}\bar{z} + \bar{\epsilon} \quad (1)$$

where the latent variables (principal components) are multivariate normal with unit covariance matrix $\bar{z} \sim \mathcal{N}(\bar{0}, I)$, and the additive noise is normal with variance σ^2 , hence, $\bar{\epsilon} \sim \mathcal{N}(\bar{0}, \sigma^2 I)$. The observations are multivariate normal with covariance matrix

$$\bar{\Sigma} = \bar{A}\bar{A}^T + \sigma^2 I. \quad (2)$$

Using

$$p(\bar{x}, \bar{z} | \bar{A}, \sigma_0^2) = p(\bar{x} | \bar{z}) p(\bar{z}) \propto e^{-\frac{1}{2\sigma^2} \|\bar{A}\bar{z} - \bar{x}\|^2} e^{-\frac{1}{2} \bar{z}^T \bar{z}}, \quad (3)$$

we get the distribution of the principal components conditioned on observations

$$\log p(\bar{z} | \bar{x}) = -\frac{1}{2\sigma^2} \|\bar{A}\bar{z} - \bar{x}\|^2 - \frac{1}{2} \bar{z}^T \bar{z} + \text{const.}, \quad (4)$$

$$= -\frac{1}{2\sigma^2} \bar{z}^T \bar{A}\bar{A}^T \bar{z} + \frac{1}{\sigma^2} \bar{x}^T \bar{A}\bar{z} - \frac{1}{2} \bar{z}^T \bar{z} + \text{const.} \quad (5)$$

Hence, the conditional distribution of the principal components is the normal distribution $\mathcal{N}(\bar{\mu}_{z|x}, \bar{\Sigma}_{z|x})$, with

$$\begin{aligned} \bar{\mu}_{z|x} &= \frac{1}{\sigma^2} \bar{\Sigma}_{z|x} \bar{A}^T \bar{x}, \\ \bar{\Sigma}_{z|x}^{-1} &= \frac{\bar{A}\bar{A}^T}{\sigma^2} + I \end{aligned} \quad (6)$$

Inference based on a sample of N complete data points $\bar{x} \in \mathbb{R}^d$ forming the data matrix

$$\bar{X} = [\bar{x}_1, \bar{x}_2, \dots, \bar{x}_N] (\text{centered}) \quad (7)$$

is based on the singular value decomposition,

$$\bar{X} = \bar{U} \bar{S} \bar{V}^T. \quad (8)$$

Let k be the selected subspace dimension, then the subspace of interest is spanned by the columns of the matrix \bar{A} estimated by

$$\hat{\bar{A}} = \bar{U}_{1:k} \bar{S}_{1:k}. \quad (9)$$

The noise variance is estimated from the variance outside the subspace of interest

$$\widehat{\sigma^2} = \frac{\text{Tr}(\bar{\bar{X}}\bar{\bar{X}}^T) - \text{Tr}(\bar{\bar{X}}_r\bar{\bar{X}}_r^T)}{N(d-k)}. \quad (10)$$

where the subspace reconstruction of data is given by

$$\bar{\bar{X}}_r = \bar{\bar{U}}_{1:k}\bar{\bar{U}}_{1:k}^T \cdot \bar{x}. \quad (11)$$

Now consider inference based on missing data, i.e., the remaining features indexed by the set m $\bar{x} \rightarrow \bar{x}_m$. The relevant distribution of principal components conditioned on the features present \bar{x}_m is given simply by $\mathcal{N}(\bar{\mu}_{z|\bar{x}_m}, \bar{\Sigma}_{z|\bar{x}_m})$, with

$$\widehat{\bar{\Sigma}_{z|\bar{x}_m}^{-1}} = \frac{\hat{\hat{A}}_m \hat{\hat{A}}_m^T}{\widehat{\sigma^2}} + I \quad (12)$$

$$\widehat{\bar{\mu}_{z|\bar{x}_m}} = \frac{1}{\widehat{\sigma^2}} \widehat{\bar{\Sigma}_{z|\bar{x}_m}} \hat{\hat{A}}_m^T \bar{x}_m \quad (13)$$

where $\hat{\hat{A}}_m$ are the rows of subspace spanning vectors corresponding to the features present \bar{x}_m .

References

1. Bishop CM. Pattern Recognition and Machine Learning (Information Science and Statistics). Secaucus, NJ, USA: Springer-Verlag New York, Inc.; 2006.