

Dimensionality reduction-principal component analysis¹

Ningxin Yang

March 5, 2024

This document describes a specific type of dimensionality reduction technique-principal component analysis

Principal component analysis (PCA) is a widely covered machine learning method on the web. And while there are some great articles about it, many go into too much detail. Below we cover how principal component analysis works in a simple step-by-step way.

Dimensionality reduction (DR)

Let us consider a computational model \mathcal{M} with a specific realization of input sets $\mathbf{x} = (x_1, \dots, x_M)$. This model enables the analyst to predict certain *QoI*, represented in a vector $\mathbf{y} \in \mathbb{R}^N$ as a function of input parameters \mathbf{x} :

$$\mathcal{M} : \mathbf{x} \in D_{\mathbf{X}} \subset \mathbb{R}^M \rightarrow \mathbf{y} = \mathcal{M}(\mathbf{x}) \in \mathbb{R}^N \quad (1)$$

One possible set of high-dimensional outputs can be represented as $\mathbf{y} = (y_1, \dots, y_N)^T$. Considering different realizations in the experimental design $\mathcal{X} = \{\mathbf{x}^1, \dots, \mathbf{x}^K\}$, where $\mathbf{x}^k \in \mathbb{R}^M$ for $k = 1, \dots, K$, the computed responses can be collected into a data matrix as:

$$\mathbf{y} = \begin{bmatrix} (\mathbf{y}^1)^T \\ (\mathbf{y}^2)^T \\ \vdots \\ (\mathbf{y}^K)^T \end{bmatrix} = \begin{bmatrix} y_1^1 & y_2^1 & \cdots & y_N^1 \\ y_1^2 & y_2^2 & \cdots & y_N^2 \\ \vdots & \vdots & \ddots & \vdots \\ y_1^K & y_2^K & \cdots & y_N^K \end{bmatrix} \quad (2)$$

The superscript denotes the realizations of *experimental design* (or *observed samples*) and the subscript denotes the number of outputs. In an abstract form, the transformation from the *original space* $\mathcal{D}_{\mathbf{Y}} \subseteq \mathbb{R}^N$ to a *reduced space* $\mathcal{D}_{\mathbf{Z}} \subseteq \mathbb{R}^{N'}$ ($N' \ll N$) is the general form of a DR mapping:

$$\mathcal{T}_{DR} : \mathcal{D}_{\mathbf{Y}} \rightarrow \mathcal{D}_{\mathbf{Z}} \quad (3)$$

where the underlying assumption is that $\mathcal{D}_{\mathbf{Z}}$ is embedded inside $\mathcal{D}_{\mathbf{Y}}$. And the nature and number of the parameters N' is depending on the specific DR technique chosen. So far, there exists a large set of DR techniques ranging from linear to nonlinear approaches. To clarify the principles of DR, this study will adopts a simple but effective linear technique called *princepal component analysis* (PCA)². PCA is selected

¹ Motivated by Dr Truong Le to write this document

² P-R Wagner, Reto Fahrni, Michael Klippel, Andrea Frangi, and Bruno Sudret. Bayesian calibration and sensitivity analysis of heat transfer models for fire insulation panels. *Engineering structures*, 205:110063, 2020; Paul-Remo Wagner. *Stochastic Spectral Embedding in Forward and Inverse Uncertainty Quantification*. PhD thesis, ETH Zurich, 2021; Joseph B Nagel, Jörg Rieckermann, and Bruno Sudret. Principal component analysis and sparse polynomial chaos expansions for global sensitivity analysis and model calibration: Application to urban drainage simulation. *Reliability Engineering & System Safety*, 195:106737, 2020; and Manav Vohra, Paromita Nath, Sankaran Mahadevan, and Yung-Tsun Tina Lee. Fast surrogate modeling using dimensionality reduction in model inputs and field output: Application to additive manufacturing. *Reliability engineering & system safety*, 201:106986, 2020

not only for its simplicity but also for its enduring favor across diverse disciplines, underscoring its continued widespread use. It is closely related to the *Karhunen-Loève expansion*, *Hotelling transform*, and *proper orthogonal decomposition*.

In practical applications, a PCA process is mainly calculating the estimation of the expectation μ_Y and the covariance matrix Σ_Y :

$$\mu_Y = \mathbb{E}[Y] = \begin{bmatrix} \mu_{y_1} \\ \mu_{y_2} \\ \vdots \\ \mu_{y_N} \end{bmatrix}^T \quad \text{with } \mu_{y_i} = \frac{1}{K} \sum_{k=1}^K y_i^k, \quad i = 1, \dots, N \quad (4)$$

$$\Sigma_Y \approx \text{Cov}[Y] = \mathbb{E}[(Y - \mu_Y)(Y - \mu_Y)^T] \quad (5)$$

Since Σ_Y is symmetric and positive definite, one can find linearly independent eigenvectors ϕ_i with positive eigenvalues λ_i for $i = 1, \dots, N$. The characteristic vectors and values should satisfy:

$$\Sigma_Y \phi_i = \lambda_i \phi_i \quad (6)$$

The N eigenvectors of this covariance matrix (new coordinations to be projected on) are collected into a matrix $\Phi_N = \{\phi_1, \dots, \phi_N\}$ for $i = 1, \dots, N$. The corresponding eigenvalue λ_i signifies the variance of Y in direction of the i -th principle component showing the descending order $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_N$. Since original data Y now has been centered and decorrelated in Equations (4) to (6), e.g., the linearly transformed random vector has mean zero $\mathbb{E}[Z] = 0$ and the diagonal covariance matrix $\text{Cov}[Z] = \mathbb{E}[ZZ^T] = \Lambda$. Orthogonal components can be then represented as:

$$Z_{\text{full}} = \Phi_N^T (Y - \mu_Y) \quad (7)$$

Figure 1 shows an example of PCA process centering and decorrelating on a two-dimensional data. In which, linear non-correlated coordinations (e.g., PC1 and PC2) can be found based on the descending orders of total variance.

By retaining only those N' principal components with the highest variance, the model output Y can then be compressed to a lower dimensional subspace Z .

$$Z = \Phi_{N'}^T (Y - \mu_Y) \approx Z_{\text{full}} \quad (8)$$

Each realization (*observation sample*) can reduce the number of outputs though Equation (8). A specific realization x^k can be visualized in Figure 2 based on the selected eigenvectors $\Phi_{N'}^T$.

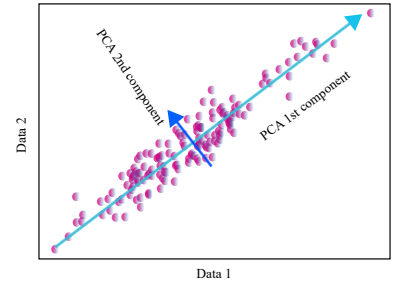


Figure 1: An example of PCA process on a two-dimensional data

$$z^k = N' \left\{ \begin{matrix} \Phi_{N'}^T \\ \text{ } \end{matrix} \right\} \times \begin{matrix} (y^k - \mu_{y^k}) \\ \text{ } \end{matrix}$$

Figure 2: A visualized realization of data compressing with selected eigenvectors

All dataset can be then compressed while retaining most of the total variation by:

$$\mathbf{z} = \begin{bmatrix} (\mathbf{z}^1)^T \\ (\mathbf{z}^2)^T \\ \vdots \\ (\mathbf{z}^K)^T \end{bmatrix} = \begin{bmatrix} z_1^1 & z_2^1 & \cdots & z_{N'}^1 \\ z_1^2 & z_2^2 & \cdots & z_{N'}^2 \\ \vdots & \vdots & \ddots & \vdots \\ z_1^K & z_2^K & \cdots & z_{N'}^K \end{bmatrix} \quad (9)$$

The selection of the number N' is determined such that $\sum_{i=1}^{N'} \lambda_i = (1 - \varepsilon_{\text{DR}}^{\text{threshold}}) \sum_{i=1}^N \lambda_i$, where $\varepsilon_{\text{DR}}^{\text{threshold}}$ is a prespecified threshold. An adequate number of principal components is to represent the system in an optimal way. If few PCs are selected than required, a poor model will be obtained. On the contrary, if more PCs than necessary are selected, the model will still face the pressure from high dimensionality. In real high-dimensional data outputs, it may also encounter the problems of over-parameterized and will include noise. To avoid these problems, several criteria for selecting the optimum number of PCs were proposed such as scree plot, explain-variance, permutation test, cross-validation and variance of reconstruction error³. With a suitable threshold and choosing criteria, the original space can be reconstructed as \mathbf{Y}^{Re} through its optimum N' principal components using Equation (10):

$$\mathbf{Y} = \boldsymbol{\mu}_Y + \sum_{i=1}^N \mathbf{z}_i \boldsymbol{\phi}_i \approx \mathbf{Y}^{Re} = \boldsymbol{\mu}_Y + \sum_{i=1}^{N'} \mathbf{z}_i \boldsymbol{\phi}_i \quad (10)$$

Similar to PCA, most DR techniques share the same ingredients as illustrated in Figure 3: (1) calculate the reduced latent spaces \mathbf{Z} (e.g., principal components in PCA) and (2) based on a predefined threshold $\varepsilon_{\text{DR}}^{\text{threshold}}$ and transform process $\mathcal{T}_{\text{DR}}^{-1}$, the outputs \mathbf{Y}^{Re} can be reconstructed. A suitable number of latent spaces N' can be retained. While PCA is a highly effective and powerful compression tool, the main idea of PCA is still based on linear decomposition. As data complexity increases, the limitations of PCA become apparent. In such cases, more advanced DR techniques, like *MDS*, *kPCA*, and *autoencoder*, come into play.

References

Baligh Mnassri, Bouchra Ananou, Mustapha Ouladsine, et al. A generalized variance of reconstruction error criterion for determining the optimum number of principal components. In *18th Mediterranean Conference on Control and Automation, MED'10*, pages 868–873. IEEE, 2010.

³ Sergio Valle, Weihua Li, and S Joe Qin. Selection of the number of principal components: the variance of the reconstruction error criterion with a comparison to other methods. *Industrial & Engineering Chemistry Research*, 38(11): 4389–4401, 1999; Edoardo Saccenti and José Camacho. Determining the number of components in principal components analysis: A comparison of statistical, crossvalidation and approximated methods. *Chemometrics and Intelligent Laboratory Systems*, 149:99–116, 2015; S Joe Qin and Ricardo Dunia. Determining the number of principal components for best reconstruction. *Journal of process control*, 10(2-3):245–250, 2000; and Baligh Mnassri, Bouchra Ananou, Mustapha Ouladsine, et al. A generalized variance of reconstruction error criterion for determining the optimum number of principal components. In *18th Mediterranean Conference on Control and Automation, MED'10*, pages 868–873. IEEE, 2010

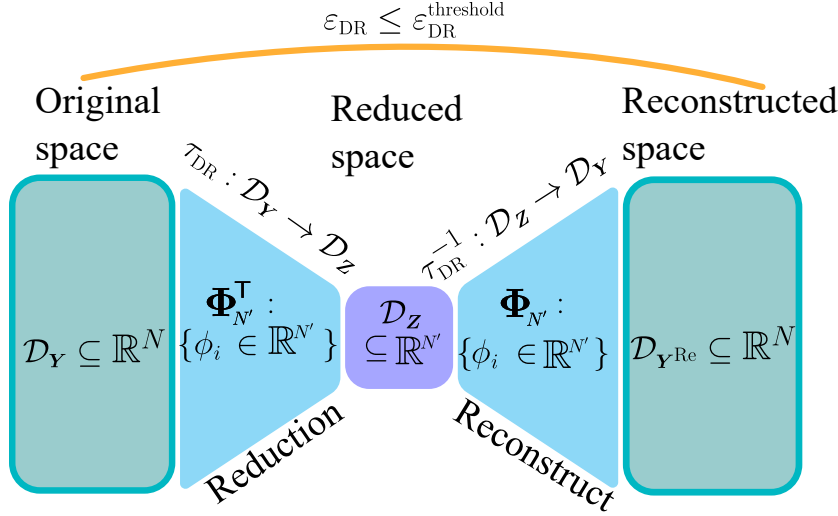


Figure 3: DR-flowchart

Joseph B Nagel, Jörg Rieckermann, and Bruno Sudret. Principal component analysis and sparse polynomial chaos expansions for global sensitivity analysis and model calibration: Application to urban drainage simulation. *Reliability Engineering & System Safety*, 195: 106737, 2020.

S Joe Qin and Ricardo Dunia. Determining the number of principal components for best reconstruction. *Journal of process control*, 10 (2-3):245–250, 2000.

Edoardo Saccenti and José Camacho. Determining the number of components in principal components analysis: A comparison of statistical, crossvalidation and approximated methods. *Chemometrics and Intelligent Laboratory Systems*, 149:99–116, 2015.

Sergio Valle, Weihua Li, and S Joe Qin. Selection of the number of principal components: the variance of the reconstruction error criterion with a comparison to other methods. *Industrial & Engineering Chemistry Research*, 38(11):4389–4401, 1999.

Manav Vohra, Paromita Nath, Sankaran Mahadevan, and Yung-Tsun Tina Lee. Fast surrogate modeling using dimensionality reduction in model inputs and field output: Application to additive manufacturing. *Reliability engineering & system safety*, 201:106986, 2020.

P-R Wagner, Reto Fahrni, Michael Klippel, Andrea Frangi, and Bruno Sudret. Bayesian calibration and sensitivity analysis of heat transfer models for fire insulation panels. *Engineering structures*, 205:110063, 2020.

Paul-Remo Wagner. *Stochastic Spectral Embedding in Forward and Inverse Uncertainty Quantification*. PhD thesis, ETH Zurich, 2021.