

1 Principal Component Analysis

Given a collection of measured data points in a high dimensional space, it is often desirable to reduce the dimensionality of such data to a lower dimensional latent space, this could be to aid to processing the data or enable visualisation. In such instances, it is desirable for the mapping to the latent space to maintain as much information from the original data as possible.

Principal Component Analysis (PCA) aims to reduce the dimensionality of a collection of data points while maximising the variance in the latent space. Alternatively, PCA aims to find a mapping from the original data to a new latent space which allow for the original data to be reconstructed with minimal error. These two aims are in fact equivalent.

1.1 Title?

Let $\mathbf{x}_i \in \mathbb{R}^f$ represent a data point in f dimensional space and $\mathbf{y}_i \in \mathbb{R}^d$ represent the point which \mathbf{x}_i is transformed to by a linear transformation \mathbf{W} equation (1), where $d \ll f$.

$$\mathbf{y}_i = \mathbf{W}^\top \mathbf{x}_i \quad (1)$$

where,

$$\mathbf{y}_i = \begin{bmatrix} y_{i1} \\ \vdots \\ y_{id} \end{bmatrix}, \quad \mathbf{x}_i = \begin{bmatrix} x_{i1} \\ \vdots \\ x_{if} \end{bmatrix}, \quad \mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_d], \quad \mathbf{w}_i = \begin{bmatrix} w_{i1} \\ \vdots \\ w_{if} \end{bmatrix}$$

The optimal transformation \mathbf{W} will maximise the variance in the data, where variance is expressed by equation (2). This then follows that the optimal solution is found by maximising equation (4). To prevent the trivial solution where $\mathbf{w}_k = \infty$, constrain a fixed magnitude $\mathbf{W}^\top \mathbf{W} = \mathbf{I}$.

$$\sigma_y^2 = \frac{1}{N} \sum_{i=1}^N (y_{ik} - \mu_k)^2 \quad (2)$$

where,

$$\begin{aligned} \mu_k &= \frac{1}{N} \sum_{i=1}^N y_{ik} \\ \mathbf{W} &= \arg \max_{\mathbf{W}} \frac{1}{N} \sum_{k=1}^d \sum_{i=1}^N (y_{ik} - \mu_k)^2 \\ &= \arg \max_{\mathbf{W}} \frac{1}{N} \sum_{k=1}^d \sum_{i=1}^N \mathbf{w}_k^\top (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top \mathbf{w}_k \\ &= \arg \max_{\mathbf{W}} \sum_{k=1}^d \mathbf{w}_k^\top \mathbf{S}_t \mathbf{w}_k \end{aligned}$$

where \mathbf{S}_t is the covariance matrix,

$$\mathbf{S}_t = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top \quad (3)$$

and,

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

$$\mathbf{W} = \arg \max_{\mathbf{W}} \text{tr}[\mathbf{W}^\top \mathbf{S}_t \mathbf{W}], \quad \text{subject to} \quad \mathbf{W}^\top \mathbf{W} = \mathbf{I} \quad (4)$$

By constructing the Lagangian from equation (4) and solving the solution in equation (5) can be obtained. From eigendecomposition, \mathbf{W} has columns of columns of eigenvectors which correspond to the d largest non-zero eigenvalues of the covariance matrix, \mathbf{S}_t .

$$L(\mathbf{W}, \boldsymbol{\Lambda}) = \text{tr}[\mathbf{W}^\top \mathbf{S}_t \mathbf{W}] - \text{tr}[\boldsymbol{\Lambda}(\mathbf{W}^\top \mathbf{W} - \mathbf{I})]$$

$$\frac{\partial \text{tr}[\mathbf{W}^\top \mathbf{S}_t \mathbf{W}]}{\partial \mathbf{W}} = 2\mathbf{S}_t \mathbf{W}, \quad \frac{\partial \text{tr}[\boldsymbol{\Lambda}(\mathbf{W}^\top \mathbf{W} - \mathbf{I})]}{\partial \mathbf{W}} = 2\mathbf{W} \boldsymbol{\Lambda}$$

$$L(\mathbf{W}, \boldsymbol{\Lambda}) = 0$$

$$\mathbf{S}_t \mathbf{W} = \mathbf{W} \boldsymbol{\Lambda} \quad (5)$$

References