

1 Principal Component Analysis

1.1 Motivation

PCA wants to identify a *meaningful* basis to re-express the dataset [1]. PCA assumes that a *meaningful* data representation is one which

1. the features with large variance have meaningful structure and should be preserved
2. the features with small variance are noise and should be discarded
3. correlated features indicate redundancy and should be made uncorrelated

Suppose we have observations $\{x_i\}_{i=1}^N$ where $x_i \in \mathbb{R}^p$ for some random variable \mathbf{x} . We want to find linear transformation of \mathbf{x} to obtain \mathbf{y} . In particular, let $\mathbf{X} \in \mathbb{R}^{N \times p}$ be stacked observations, we want to find a linear map $\mathbf{P} \in \mathbb{R}^{p \times q}$, where columns of \mathbf{P} are orthonormal basis for feature space, i.e. $\text{row}(\mathbf{X})$, to re-express data \mathbf{X} to $\mathbf{Y} \in \mathbb{R}^{N \times q}$.

$$\mathbf{Y} = \mathbf{XP}$$

\mathbf{Y} has a meaningful representation if $\text{cov}(\mathbf{Y})$ is a diagonal matrix (decorrelated), and that successive dimension in \mathbf{Y} are rank-ordered according to variance (preserve struture, discard noise).

1.2 Empirical Covariance Matrix

Note that for a random variable \mathbf{x} with stacked observations $\mathbf{X} \in \mathbb{R}^{N \times p}$, the empirical covariance for $\mathbf{x}_i, \mathbf{x}_j$ is given by

$$\hat{\sigma}^2(\mathbf{x}_i, \mathbf{x}_j) = \frac{1}{N} \sum_i (x_i - \bar{x}_i)(x_j - \bar{x}_j) = \frac{1}{N} (\mathbf{X}_i - \bar{\mathbf{X}}_i \mathbf{1}_N)^T (\mathbf{X}_j - \bar{\mathbf{X}}_j \mathbf{1}_N)$$

where $\mathbf{X}_i, \mathbf{X}_j$ are i and j -th column of \mathbf{X} and $\bar{\mathbf{X}}_i = \frac{1}{N} \sum_j \mathbf{X}_{ji}$. So then,

$$\widehat{\text{Cov}}(\mathbf{x}) = [\hat{\sigma}(\mathbf{x}_i, \mathbf{x}_j)]_{i,j=1}^p = \frac{1}{N} (\mathbf{X} - \bar{\mathbf{X}} \mathbf{1}_N)^T (\mathbf{X} - \bar{\mathbf{X}} \mathbf{1}_N)$$

where $\bar{\mathbf{X}}$ is column wise feature average of \mathbf{X} . For zero mean observation matrix, the empirical covariance matrix is simply $\frac{1}{N} \mathbf{X}^T \mathbf{X}$

1.3 Solving PCA using Eigenvector Decomposition

We first write covariance matrix for \mathbf{Y} ,

$$\widehat{\text{Cov}}(\mathbf{y}) = \frac{1}{N} \mathbf{Y}^T \mathbf{Y} = \frac{1}{N} (\mathbf{XP})^T (\mathbf{XP}) = \mathbf{P}^T \left(\frac{1}{N} \mathbf{X}^T \mathbf{X} \right) \mathbf{P} = \mathbf{P}^T \widehat{\text{Cov}}(\mathbf{x}) \mathbf{P}$$

We know that $\widehat{\text{Cov}}(\mathbf{x})$ is a symmetric matrix and therefore can be written as $\widehat{\text{Cov}}(\mathbf{x}) = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^T$ where $\mathbf{Q} \in \mathbb{R}^{p \times p}$ are eigenvectors of $\widehat{\text{Cov}}(\mathbf{x})$ with corresponding eigenvalues along diagonal entries in $\mathbf{\Lambda}$. Setting projection to be eigenvectors of $\widehat{\text{Cov}}(\mathbf{x})$ diagonalizes $\widehat{\text{Cov}}(\mathbf{y})$,

$$\mathbf{P} \leftarrow \mathbf{Q} \quad \Rightarrow \quad \widehat{\text{Cov}}(\mathbf{y}) = \mathbf{P}^T \widehat{\text{Cov}}(\mathbf{x}) \mathbf{P} = \mathbf{Q}^T \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^T \mathbf{Q} = \mathbf{\Lambda}$$

where $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}$. The *principal components* of \mathbf{X} are column vectors of \mathbf{P} , i.e. eigenvectors for $\widehat{\text{Cov}}(\mathbf{x})$. \mathbf{y} is decorrelated and $\hat{\sigma}^2(\mathbf{y}_i)$ is the variance of \mathbf{x} along i -th principal component.

1.4 Singular Value Decomposition

The singular value decomposition of an arbitrary matrix $\mathbf{X} \in \mathbb{R}^{N \times p}$ is

$$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$$

where left singular vectors $\mathbf{U} \in \mathbb{R}^{N \times N}$ is orthogonal, singular values $\mathbf{\Sigma} \in \mathbb{R}^{N \times p}$ is diagonal, right singular vectors $\mathbf{V} \in \mathbb{R}^{p \times p}$ is orthogonal. If $\mathbf{X}^T\mathbf{X}$ has rank r , then column vectors of \mathbf{V} are eigenvectors with eigenvalues $\{\lambda_i\}_{i=1}^r$ (assuming descending ordering) for symmetric matrix $\mathbf{X}^T\mathbf{X}$, i.e. $(\mathbf{X}^T\mathbf{X})\mathbf{v}_i = \lambda_i\mathbf{v}_i$. Entries along the diagonals of $\mathbf{\Sigma}$ are singular values $\sigma_i = \sqrt{\lambda_i}$. The column vectors of \mathbf{U} are given by $\mathbf{u}_i = \frac{1}{\sigma_i}\mathbf{X}\mathbf{v}_i$. We can show that column vectors of \mathbf{U} are unit orthonormal vectors. Grouping linear relationships $\mathbf{X}\mathbf{v}_i = \sigma_i\mathbf{u}_i$ yield $\mathbf{X}\mathbf{V} = \mathbf{\Sigma}\mathbf{U}$. Note, \mathbf{V} acts similarly to the projection matrix \mathbf{P} .

1.5 Solving PCA using SVD

From previous, the principal components of \mathbf{X} are the eigenvectors of $\widehat{\text{Cov}}(\mathbf{x})$. Let $\mathbf{W} = \frac{1}{\sqrt{N}}\mathbf{X}$, then right singular vectors of \mathbf{W} are the principal components desired,

$$\mathbf{W}^T\mathbf{W} = \left(\frac{1}{\sqrt{N}}\mathbf{X}\right)^T \left(\frac{1}{\sqrt{N}}\mathbf{X}\right) = \frac{1}{N}\mathbf{X}^T\mathbf{X} = \widehat{\text{Cov}}(\mathbf{x})$$

and that $\hat{\sigma}^2(y_i) = \mathbf{\Sigma}_{ii}^2$

1.6 Limitations

PCA works well with Gaussian observations, in particular the transformed data is guaranteed to be independent. If \mathbf{x} is jointly Gaussian, then any linear function of \mathbf{x} is also jointly Gaussian. Suppose $\mathbf{b} \sim \mathcal{N}(\boldsymbol{\mu}_b, \mathbf{\Sigma}_b)$

$$\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_x, \mathbf{\Sigma}_x) \quad \Rightarrow \quad A\mathbf{x} + \mathbf{b} \sim \mathcal{N}(A\boldsymbol{\mu}_x + \boldsymbol{\mu}_b, A\mathbf{\Sigma}_x A^T + \mathbf{\Sigma}_b)$$

Therefore, the transformed data \mathbf{y} are jointly Gaussian. Any two variables y_i, y_j are uncorrelated (by diagonal $\widehat{\text{Cov}}(\mathbf{y})$) and therefore independent (by \mathbf{y} jointly Gaussian). For non jointly Gaussian data, we can not assume independence in the transformed data. In other words, PCA is not able to reveal non-linear relationships between features.

2 Eigenfaces for Recognition

Eigenfaces project a set of faces to the *face space*, spanned by a set of orthonormal *eigenfaces*, which best encode variation amongst faces [2, 3, 4]. In practice this means doing SVD on the set of zero mean faces \mathbf{X} , pick first M right singular vectors $\mathbf{V} \in \mathbb{R}^{p \times M}$ associated with largest singular values. columns of \mathbf{V} are called *eigenfaces* and $\text{col}(\mathbf{V})$ is the *face space*. We can project a new image $\mathbf{x} \in \mathbb{R}^{1 \times p}$ to the face space, $\mathbf{y} = (\mathbf{x} - \bar{\mathbf{X}})\mathbf{V}$ and classify faces to class $k = \arg \min_{k \in 1:N} \|\mathbf{y} - (\mathbf{X}\mathbf{V})_k\|$.

References

- [1] Jonathon Shlens. “A Tutorial on Principal Component Analysis”. In: *arXiv:1404.1100 [cs, stat]* (Apr. 3, 2014). arXiv: [1404.1100](https://arxiv.org/abs/1404.1100). URL: <http://arxiv.org/abs/1404.1100> (visited on 01/31/2020).
- [2] Alex P Pentland. “Face Recognition Using Eigenfaces”. In: (1991), p. 6.
- [3] Matthew Turk and Alex Pentland. “Eigenfaces for Recognition”. In: *Journal of Cognitive Neuroscience* 3.1 (Jan. 1991), pp. 71–86. ISSN: 0898-929X, 1530-8898. DOI: [10.1162/jocn.1991.3.1.71](https://doi.org/10.1162/jocn.1991.3.1.71). URL: <http://www.mitpressjournals.org/doi/10.1162/jocn.1991.3.1.71> (visited on 02/04/2020).
- [4] Jun Zhang, Yong Yan, and M. Lades. “Face recognition: eigenface, elastic matching, and neural nets”. In: *Proceedings of the IEEE* 85.9 (Sept. 1997), pp. 1423–1435. ISSN: 1558-2256. DOI: [10.1109/5.628712](https://doi.org/10.1109/5.628712).