

# ToolBox Documentation

Joseph Raso

January 2021

## Contents

<b>1</b>	<b>Principal Component Analysis</b>	<b>2</b>
1.1	Theory . . . . .	2
1.2	Tools . . . . .	3
1.2.1	( <i>class</i> ) PrincipalComponents.PCA(X) . . . . .	3

# 1 Principal Component Analysis

## 1.1 Theory

PCA is a "statistical interpretation of the singular value decomposition". It is a bedrock tool for discovering the axes along which a data set varies the most. Starting from a data matrix  $X$ , with each row corresponding to one of  $n$  samples and each column corresponding to one of  $m$  dimensions<sup>1</sup>,

$$X = \begin{bmatrix} - & x_1 & - \\ - & x_2 & - \\ & \dots & \\ - & x_n & - \end{bmatrix} \quad (1)$$

we first find the mean sample

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad (2)$$

and subtract it off the data matrix:

$$B = X - \bar{X} = \begin{bmatrix} - & x_1 & - \\ - & x_2 & - \\ & \dots & \\ - & x_n & - \end{bmatrix} - \begin{bmatrix} - & \bar{x} & - \\ - & \bar{x} & - \\ & \dots & \\ - & \bar{x} & - \end{bmatrix} \quad (3)$$

Now  $B$  the same data matrix (i.e. same distribution), but shifted to be mean-zero.

Next, now that the data is mean-zero, it's pretty trivial to compute the covariance between the dimensions of  $B$  as the covariance matrix  $C$ .

$$C = \frac{B^T B}{n} \quad (4)$$

$C$  is called the covariance matrix because it's elements are the variance and covariance between the various dimensions in  $B$ ,<sup>2</sup> i.e. for element  $C_{ij}$ ,

$$C_{ij} = \frac{b_i \cdot b_j}{n} = \frac{1}{n} \sum_{k=1}^n (b_i)_k (b_j)_k = \mathbb{E}[b_i \cdot b_j] = \text{Cov}(b_i, b_j) \quad (5)$$

Eigenvalue decomposition of  $C$  therefor yields:

1. A dimensional basis  $V$  in which the dimensional distributions are independent (have covariance 0); i.e. the basis in which  $C$  is diagonal.
2. The diagonalized covariance  $\Lambda$ , with diagonal elements  $\lambda_1, \lambda_2, \dots, \lambda_m$  which are equal to the variance of the data along the corresponding eigen-axis.

It follows that the eigenvector corresponding to the largest eigenvalue is the independent axis along which the data has the highest variance. The matrix  $V$  of eigenvectors is referred to as the *loadings*. Finally, to find the amount of each principal component in the original samples, we find

$$T = BV. \quad (6)$$

---

<sup>1</sup>This is the transpose of the setup for singular value decomposition.

<sup>2</sup>Note the dimensions:  $B$  is  $n$  samples by  $m$  dimensions just like  $X$ , and consequently  $C$  is  $m \times m$ .

**Connection to SVD** If we look at the singular value decomposition of  $B$  into left eigenvectors  $U$  and right eigenvectors  $V$ ,

$$B = U\Sigma V^T \quad (7)$$

it follows that

$$C = B^T B = (V\Sigma^T U^T)(U\Sigma V^T) = V(\Sigma^T \Sigma)V^T \quad (8)$$

$$\implies CV = V(\Sigma^T \Sigma). \quad (9)$$

Since the elements of the square matrix  $\Sigma^T \Sigma$  are equal to the variances along the corresponding right eigenvectors  $V$ ,  $\Sigma^T \Sigma = \Lambda$ , giving the right eigenvectors  $V$  the same values as in the PCA case. We can immediately then see that the loadings  $V$ , variances  $\Lambda$  and principal-component transformed data  $T$  is immediately available from an SVD breakdown of  $B$ :

$$V = V \quad (10)$$

$$\Lambda = \Sigma^T \Sigma \quad (11)$$

$$T = BV = U\Sigma \quad (12)$$

**Dimensionality Reduction** PCA can be used to sift out components dimensions in the data that are unimportant, allowing the size of the data set to be reduced to a set of axis that are most meaningful.

Deciding which components to keep is done by looking at the fraction  $f_k$ , where for the first  $k$  eigenvalues (for  $n$  total, sorted by magnitude):

$$f_k = \frac{\sum_{i=1}^k \lambda_i}{\sum_{j=1}^n \lambda_i} \quad (13)$$

## 1.2 Tools

### 1.2.1 (*class*) PrincipalComponents.PCA(X)

Object representing the principal component analysis of input X.

#### Parameters:

---

**X** (*ndarray*) The data matrix, with rows as samples, and columns as dimensions. (Is not required to be mean-zero.)

---

#### Attributes:

---

**X** (*ndarray*) The original input data.

**xbar** (*ndarray*) the mean sample (row) of X.

**B** (*ndarray*) The data X, shifted to mean-zero.

**C** (*ndarray*) The covariance matrix of B, normalized to the number of samples (rows of X).

**s** (*ndarray*) Variance along each principal component, sorted by magnitude, normalized to the number of samples.

- V** (*ndarray*) Loadings of each principal component as column vectors, sorted by magnitude of the variance (same order as *v*), oriented to put the first dimension of the first principal component in the positive direction.
- T** (*ndarray*) The data *X*, rotated into the PCA basis.
-