

A DATA VISUALIZATION FOR PCA ANALYSIS / MAI718 3/2017 FINAL PROJECT

Roberto Stelling *

PPGI - 117.335.792

Universidade Federal do Rio de Janeiro

Rio de Janeiro, RJ 21941-916, Brazil

roberto@stelling.cc

ABSTRACT

Traditionally PCA Analysis for features reduction in images is performed without a systematic visual aid approach. This paper describes the implementation of a PCA Analysis tool in JavaScript that uses data visualization techniques aiming to improve the analysis process.

1 PROBLEM DESCRIPTION

There are many reasons to reduce a data set, some examples are: noise reduction, outlier removal, lossy image compression or even as a preliminary step in various types of data exploration and data analysis.

Principal component analysis (PCA) is well suited as a lossy image compression solution, as you can apply PCA on a set of points of a data set and reduce its dimensionality with a certain, desirably controllable, loss of precision. Of course, one want to loose as little precision as possible while compressing as much as possible. There is a clear trade off between compression and precision. With images, the main difficulty resides in this very trade off: how many dimensions can the algorithm throw away and still retain the desired level of image quality or sharpness? Is there a general rule where you can certainly decide how many dimensions will be cut off the original data set? Of course, the type and amount of data available for compression, the data set, has a big influence on the final point where the cut will be, but can the analyst decide simply on the number of dimensions or amount of variance that will be thrown away and be sure that the results will be satisfactory? We propose that using a visual helping tool during the decision process can have a positive impact on the cut off selection.

2 BRIEF INTRODUCTION TO PCA

2.1 OBJECTIVE

According to Jolliffe (1986), the central idea of principal component analysis (PCA) is to reduce the dimensionality of a data set in which there are a large number of interrelated variables, while retaining as much as possible of the variation present in the data set.

This reduction is achieved by transforming the data set to a new set of variables, the principal components, which are not correlated, and which are ordered so that the first few retain most of the variation present in all the original variables.

Principal component analysis was first described by Pearson (1901) and later developed independently by Hotelling (1933) (Jolliffe, 1986).

Craw & Cameron (1992) describe a method for face recognition using principal component analysis, showing that PCA can be used as an effective tool in image analysis.

*PPGI/UFRJ Graduate Student, <http://stelling.cc>.

2.2 INTUITION

PCA can be thought of as the problem of fitting an n -dimensional ellipsoid to the m -dimensional data, where $n \leq m$ and each axis of the ellipsoid represents a principal component. The larger the axis of a component, the larger the variance for that component. So, the objective of PCA is to build a transformation of m -dimensional space to n -dimensional space while preserving most of the m -dimensional space variance. To find that transformation and the components, we compute the singular value decomposition of the data. The singular value decomposition will provide a computationally efficient method of finding the principal components and the scaled versions of the principal component scores.

2.3 SINGULAR VALUE DECOMPOSITION

Given an arbitrary $D_{m \times n}$ matrix, then D can be written as

$$D = USV^T \quad (1)$$

where

- (i) $U_{m \times r}, V_{n \times r}$ each of which with orthonormal¹ columns so that $U^T U = I_r, V^T V = I_r$;
- (ii) $S_{r \times r}$ is a diagonal matrix;
- (iii) r is the rank² of D .

S is a diagonal matrix such as:

$$S = \begin{bmatrix} s_1 & & \\ & \ddots & \\ & & s_r \end{bmatrix}$$

s_1 to s_r are the principal components scores and $s_1 \geq s_2 \geq \dots \geq s_{r-1} \geq s_r$.

U is the eigenvector matrix, with eigenvectors ordered by the component scores, their eigenvalues.

2.4 REDUCING DIMENSIONS

Given the original data set, $D_{m \times n}$ and $U_{m \times r}$ obtained from equation (1), we can build a new reduced data set $P_{m \times k}$ with the first $k < r$ eigenvectors. This new data set is computed as:

$$P_{m \times k} = D_{m \times n} U_{n \times k} \quad (2)$$

where $U_{n \times k}$, or U_k , is the eigenvector matrix truncated to the first k eigenvectors.

2.5 RESTORING DATA

To restore the original data from P we compute

$$PU_k^T = DU_k U_k^T = DI_n = D$$

3 PCA FOR IMAGE COMPRESSION

A computer image is usually thought of as a two dimensional matrix, with m lines and n columns representing the horizontal e vertical pixels of the image. A simpler, albeit equally meaningful, representation is a single vector with $m \times n$ cells for the whole image. The content of each cell, in either representation, depend on the selected image mode. For example: RGB, RGB grayscale, CMYK, etc. For the purposes of the following argument and the presented implementation, we assume that each cell is an integer between 0 and 255, representing the grayscale RGB value of the pixel.

¹both orthogonal and normalized

²corresponds to the maximal number of linearly independent columns of D

Lets assume that D is a data set with m images where each image has n pixels. The method will work even if $m < n$ but a larger m will result in better compression gains and finer eigenvector tuning.

So $D_{m \times n}$ is a data set with m data points $\in \mathbb{R}^n$ where $m \geq n$. Then we define $D_{m \times n}^*$ as the normalized data set,

$$D^* = \frac{D - \bar{D}}{s}$$

where \bar{D} is the mean of D and s is the sample standard deviation. Let $\Sigma_{n \times n} = \text{cov}(D^*)$ be the covariance matrix of D^*

$$\Sigma = \frac{1}{m} D^{*T} D^*$$

Then, according to equation (1), the singular value decomposition of Σ is:

$$\Sigma = U S V^T$$

where:

- U is an $n \times n$ unitary matrix³
- S is an $n \times n$ diagonal matrix with non-negative numbers on the diagonal
- V is an $n \times n$ unitary matrix and V^T is V transposed.

4 HOW TO SELECT THE NUMBER OF COMPONENTS TO RETAIN

The problem of selecting how many components to retain is not new, Zwick & Velicer (1986), present the results of a Monte Carlo evaluation of five methods that have been proposed for determining how many factors or components to retain: Horn’s parallel analysis, Velicer’s minimum average partial [MAP], Cattell’s scree test, Bartlett’s chi-square test, and Kaiser’s eigenvalue greater than 1.0 rule. The determination of the number of components or factors to retain is likely to be the most important decision a researcher will make Zwick & Velicer (1986).

5 IMPLEMENTATION

The code for the implementation and supporting documents can be found at GitHub.

6 CITATIONS, FIGURES, TABLES, REFERENCES

These instructions apply to everyone, regardless of the formatter being used.

6.1 CITATIONS WITHIN THE TEXT

Citations within the text should be based on the `natbib` package and include the authors’ last names and year (with the “et al.” construct for more than two authors). When the authors or the publication are included in the sentence, the citation should not be in parenthesis (as in “See Goodfellow et al. (2016) for more information.”). Otherwise, the citation should be in parenthesis (as in “Deep learning shows promise to make progress towards AI (Goodfellow et al., 2016).”).

The corresponding references are to be listed in alphabetical order of authors, in the REFERENCES section. As to the format of the references themselves, any style is acceptable as long as it is used consistently.

6.2 FOOTNOTES

Indicate footnotes with a number⁴ in the text. Place the footnotes at the bottom of the page on which they appear. Precede the footnote with a horizontal rule of 2 inches (12 picas).⁵

³If U is unitary then $U^* U = U U^* = I$

⁴Sample of the first footnote

⁵Sample of the second footnote

Table 1: Sample table title

PART	DESCRIPTION
Dendrite	Input terminal
Axon	Output terminal
Soma	Cell body (contains cell nucleus)

6.3 FIGURES

All artwork must be neat, clean, and legible. Lines should be dark enough for purposes of reproduction; art work should not be hand-drawn. The figure number and caption always appear after the figure. Place one line space before the figure caption, and one line space after the figure. The figure caption is lower case (except for first word and proper nouns); figures are numbered consecutively.

Make sure the figure caption does not get separated from the figure. Leave sufficient space to avoid splitting the figure and figure caption.

You may use color figures. However, it is best for the figure captions and the paper body to make sense if the paper is printed either in black/white or in color.

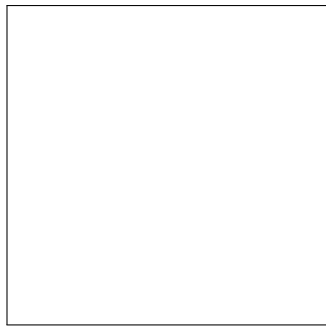


Figure 1: Sample figure caption.

6.4 TABLES

All tables must be centered, neat, clean and legible. Do not use hand-drawn tables. The table number and title always appear before the table. See Table 1.

Place one line space before the table title, one line space after the table title, and one line space after the table. The table title must be lower case (except for first word and proper nouns); tables are numbered consecutively.

ACKNOWLEDGMENTS

Use unnumbered third level headings for the acknowledgments. All acknowledgments, including those to funding agencies, go at the end of the paper.

REFERENCES

Ian Craw and Peter Cameron. Face recognition by computer. In *BMVC92*, pp. 498–507. Springer, 1992.

GitHub. <https://github.com/robstelling/imgpca>.

Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*, volume 1. MIT press Cambridge, 2016.

- Harold Hotelling. Analysis of a complex of statistical variables into principal components. *Journal of educational psychology*, 24(6):417, 1933.
- Ian T Jolliffe. Principal component analysis and factor analysis. In *Principal component analysis*, pp. 115–128. Springer, 1986.
- Karl Pearson. Principal components analysis. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, 6(2):559, 1901.
- William R Zwick and Wayne F Velicer. Comparison of five rules for determining the number of components to retain. *Psychological bulletin*, 99(3):432, 1986.