

# ANU Machine Learning Summer School 2010

## Handout for Lab 4: Principal Component Analysis

Daniel Visentin

School of Computer Science  
College of Engineering and Computer Science  
The Australian National University  
27<sup>th</sup> September 2010

**Abstract.** In this lab we learn about principal component analysis and apply it to compress a set of images.

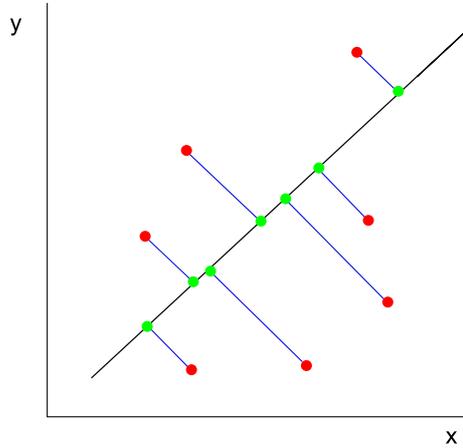
### 1 Introduction

Principal component analysis (PCA) [1–3] is a technique for projecting a data set onto a lower-dimensional subspace (see Figure 1 for an example). Intuitively, it works by projecting the data such that most of the “information” is contained in the first few dimensions of the projection. This makes it suitable for a number of applications:

- Dimensionality reduction: If most of the important information is contained in the first few dimensions of the projected data, the remaining dimensions can often be safely discarded without significant detriment to subsequent analysis. Furthermore, this reduction in dimension will often result in accelerated performance.
- Lossy data compression: By performing dimensionality reduction on a data set, we can represent each data point by a smaller number of components. Applying the inverse projection to the projected data gives an approximation of the original data. Since most of the information is conserved while the dimensionality is reduced, this can be used as a compression scheme.
- Feature extraction: The projection may elucidate some structures/features of the data set which are more useful or natural than its original representation.
- Data visualisation: It is often difficult to visualise high-dimensional data sets. Since most of the information is contained in the initial dimensions of the projected data, visualising the projected data might still offer some insights into the original data set.

### 2 The Algorithm

Principal component analysis involves computing the orthogonal projection of a data set onto a space of dimensionality  $M$  such that the variance of the projected



**Fig. 1.** Example of PCA – here we have projected the red points onto the central line, resulting in the green points. So we have reduced the two dimensional data (red) set to a one dimensional (the line) data set (green). The line was chosen by minimizing the sum of squared projection errors, which is the length of the blue lines.

data is maximised. This is equivalent to computing the orthogonal projection that minimises the average projection cost (squared distance between the data points and their projections). By solving either of these formulations we find that this projection is defined by the  $M$  eigenvectors of the data covariance matrix  $\mathbf{S}$  corresponding to the  $M$  largest eigenvalues of  $\mathbf{S}$ . We refer to these eigenvectors as the *principal components* and the space they span as the *principal subspace*.

## 2.1 Deriving the Principal Components

Let  $(\mathbf{x}_1, \dots, \mathbf{x}_N)$  be a data set of  $D$ -dimensional column vectors. Let  $\mathbf{X}$  denote the  $D \times N$  data matrix whose  $i^{\text{th}}$  column is given by  $\mathbf{x}_i$ . Our goal is to project the data onto an  $M$ -dimensional space in a way that maximises the projected variance.

We proceed iteratively, beginning with  $M = 1$ . In this case, the principal subspace is defined by a single  $D$ -dimensional unit vector  $\mathbf{u}_1$ . A data point  $\mathbf{x}_i$  is projected onto the principal subspace by taking the dot product  $\mathbf{u}_1^T \mathbf{x}_i$ . We choose  $\mathbf{u}_1$  which maximises the variance of the projected data

$$\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 = \frac{1}{N} \sum_{i=1}^N (\mathbf{u}_1^T \mathbf{x}_i - \mathbf{u}_1^T \bar{\mathbf{x}})^2 \quad (1)$$

where  $\bar{\mathbf{x}}$  is the mean vector

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i,$$

$\mathbf{S}$  is the  $D \times D$  data covariance matrix

$$\mathbf{S} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T = \frac{1}{N} (\mathbf{X} - \bar{\mathbf{X}}) (\mathbf{X} - \bar{\mathbf{X}})^T,$$

and  $\bar{\mathbf{X}}$  is the  $D \times N$  matrix whose columns are  $\bar{\mathbf{x}}$ . Since  $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 \rightarrow \infty$  as  $\|\mathbf{u}_1\| \rightarrow \infty$ , we constrain  $\mathbf{u}_1$  to be a unit vector  $\|\mathbf{u}_1\| = 1$ . Introducing the Lagrange multiplier  $\lambda_1$  and maximising gives the result [3]

$$\mathbf{S} \mathbf{u}_1 = \lambda_1 \mathbf{u}_1.$$

So  $\mathbf{u}_1$  must be an eigenvector of the data covariance  $\mathbf{S}$  and the projected covariance  $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$  is maximised when  $\mathbf{u}_1$  is the eigenvector with the largest eigenvalue  $\lambda_1$ .

For  $M = 2$ , we can proceed by trying to find a  $\mathbf{u}_2$  that minimizes the projected variance in the direction orthogonal to  $\mathbf{u}_1$ . Generalizing this for any  $M \leq D$ , we find that the required  $\mathbf{u}_i$ 's are the  $M$  eigenvectors of  $\mathbf{S}$  corresponding to the  $M$  largest eigenvalues  $\lambda_i$  (sorted in decreasing order). These eigenvectors are called the *principle components*. For future convenience, we let  $\mathbf{U}$  be the  $D \times M$  matrix whose  $i^{\text{th}}$  column is given by the  $i^{\text{th}}$  principal component  $\mathbf{u}_i$ .

## 2.2 Projection into the principal subspace

The projection of the data point  $\mathbf{x}_i$  into the principal subspace is the  $M$ -dimensional vector  $\mathbf{y}_i$  defined by

$$\mathbf{y}_i := (\mathbf{u}_1^T \mathbf{x}_i, \mathbf{u}_2^T \mathbf{x}_i, \dots, \mathbf{u}_M^T \mathbf{x}_i)^T = \mathbf{U}^T \mathbf{x}_i.$$

In words, the  $k^{\text{th}}$  component of the projection is given by the dot product of the  $k^{\text{th}}$  principal component  $\mathbf{u}_k$  with the data point  $\mathbf{x}_i$ . We can extend this to calculate the projection  $\mathbf{Y}$  of the entire data set in one matrix operation

$$\mathbf{Y} = \mathbf{U}^T \mathbf{X}$$

where the  $i^{\text{th}}$  column of  $\mathbf{Y}$  corresponds to  $\mathbf{y}_i$ .

## 2.3 PCA approximation

We can use PCA as a compression technique by projecting a data set into the principal subspace and storing the mean, the projections, and the principal components. We can then “decompress” these stored values by performing the inverse transformation back into the original space. A suitable choice of the principal subspace dimension  $M$  will result in a satisfactory reproduction of the original data and potentially save on space. In particular, we will need to store

$$\underbrace{D}_{\text{mean}} + \underbrace{NM}_{\text{projections}} + \underbrace{MD}_{\text{principal components}} = M(N + D) + D$$

numbers as opposed to  $ND$  numbers for the original data set. For  $M$  much less than  $N$  and large  $D$  this can result in a substantial saving.

The PCA approximation  $\tilde{\mathbf{x}}_i$  of the data point  $\mathbf{x}_i$  is given by

$$\begin{aligned}\tilde{\mathbf{x}}_i &:= \sum_{j=1}^M \mathbf{u}_j (\mathbf{u}_j^T \mathbf{x}) + \sum_{j=M+1}^D \mathbf{u}_j (\mathbf{u}_j^T \bar{\mathbf{x}}) \\ &= \bar{\mathbf{x}} + \sum_{j=1}^M \mathbf{u}_j (\mathbf{u}_j^T (\mathbf{x}_i - \bar{\mathbf{x}})) \\ &= \bar{\mathbf{x}} + \left( \sum_{j=1}^M \mathbf{u}_j \mathbf{u}_j^T \right) (\mathbf{x}_i - \bar{\mathbf{x}}) \\ &= \bar{\mathbf{x}} + \mathbf{U} \mathbf{U}^T (\mathbf{x}_i - \bar{\mathbf{x}})\end{aligned}$$

which is equivalent to

$$\tilde{\mathbf{x}}_i = \bar{\mathbf{x}} + \mathbf{U} \mathbf{y}_i$$

where we define  $\mathbf{y}_i := \mathbf{U}^T (\mathbf{x}_i - \bar{\mathbf{x}})$  to be the projection into the principal subspace of the mean-centered data point  $\mathbf{x}_i - \bar{\mathbf{x}}$ . The approximation of the entire data set can be computed in one step as

$$\widetilde{\mathbf{X}} := \bar{\mathbf{X}} + \mathbf{U} \mathbf{U}^T (\mathbf{X} - \bar{\mathbf{X}})$$

where the  $i^{\text{th}}$  column of  $\widetilde{\mathbf{X}}$  corresponds to  $\tilde{\mathbf{x}}_i$ .

## 2.4 PCA for high-dimensional data

In order to use PCA we need to compute the eigenvectors of the  $D \times D$  data covariance matrix  $\mathbf{S}$ . For high-dimensional data sets this is often unfeasibly slow. As an example, consider a data set of  $100 \times 100$  images. Then we need to find the eigenvectors of a  $10,000 \times 10,000$  matrix! However it is often the case in such circumstances that we have a much smaller number of data points  $N$  than the data dimension  $D$ . We can exploit this by letting  $\widehat{\mathbf{X}} := \mathbf{X} - \bar{\mathbf{X}}$  and noting that since the eigenvectors  $\mathbf{u}_i$  satisfy

$$\frac{1}{N} \mathbf{S} \mathbf{u}_i = \frac{1}{N} \widehat{\mathbf{X}} \widehat{\mathbf{X}}^T \mathbf{u}_i = \lambda_i \mathbf{u}_i$$

they must also satisfy

$$\frac{1}{N} \widehat{\mathbf{X}}^T \widehat{\mathbf{X}} \widehat{\mathbf{X}}^T \mathbf{u}_i = \lambda_i \widehat{\mathbf{X}}^T \mathbf{u}_i.$$

So, by finding the eigenvectors  $\mathbf{v}_i = \widehat{\mathbf{X}}^T \mathbf{u}_i$  of the  $N \times N$  matrix

$$\frac{1}{N} \widehat{\mathbf{X}}^T \widehat{\mathbf{X}}$$

we can retrieve the principal components using

$$\mathbf{u}_i = \frac{1}{(N \lambda_i)^{1/2}} \widehat{\mathbf{X}} \mathbf{v}_i.$$

### 3 Examples

#### 3.1 Eigenfaces

In this example we have access to a data set composed of 400 images of faces. Each image is normalised to  $70 \times 100$  pixels, then flattened and stored as a 7000-dimensional vector. By applying PCA to this data set we can derive a set of “eigenfaces” by interpreting the principal components as images (figure 2). We can also use PCA approximation to compress and reconstruct the faces as seen in figure 3. Because of the high-dimensionality of the data, it is necessary to use the formulation in section 2.4.



Fig. 2. First 10 eigenfaces.



Fig. 3. Compression of a face for  $M = 25, 50, 100, 200, 300, 400$ .

### 4 Navigating the Code

The `load_all_faces` function loads each image from the faces directory into a 2-dimensional numpy array where each column corresponds to an image. The code then calls the `PCA` function to compute the principal components of the data set. In order to compute the eigenvectors and eigenvalues we can use either of

the numpy functions `numpy.linalg.eig` or `numpy.linalg.svd`. The documentation for these functions can be found at:

<http://docs.scipy.org/doc/numpy/reference/routines.linalg.html>

It is recommended you use the singular value decomposition (SVD) function for this purpose since it returns as its first value a matrix of eigenvectors sorted by decreasing eigenvalue (Note: the first return value corresponds to the eigenvectors because  $\mathbf{S}$  and  $\widehat{\mathbf{X}}^T \widehat{\mathbf{X}}$  are real symmetric matrices—this is not true in general). Whereas if you do use the `numpy.linalg.eig` function you will have to sort the eigenvectors yourself. The code then saves the eigenvectors as images into the `eigenfaces` directory. Finally, the code computes the PCA approximations for varying subspace dimensions and stores the results in the `compressed` directory.

## References

1. Smith, L.: A tutorial on principal components analysis. Cornell University, USA **51** (2002) 52
2. Jolliffe, I.: Principal component analysis. Springer verlag (2002)
3. Bishop, C., et al.: Pattern recognition and machine learning. Springer New York (2006)