

# Principal components

## 1. Introduction

The idea behind the principal components is roughly the following. Suppose one has a random vector  $\mathbf{x} \in \mathbb{R}^d$  (or just a finite number of vectors  $\mathbf{x}_j \in \mathbb{R}^d$ ,  $j = 1, \dots, n$  each chosen with the same probability  $\frac{1}{n}$ ) and that the mean of these vectors is  $\mathbf{0}$ . Suppose that one wants to choose project these vectors on a subspace with smaller dimension e.g.  $m < d$ ) without losing too much information. This is here taken to mean that if  $P$  is the projection, then the variance of  $P\mathbf{x}$  is as large as possible, which turns out to be the same as that the variance of  $\mathbf{x} - P\mathbf{x}$  is as small as possible. If we take  $\mathbf{x}$  to be a column vector we can write the matrix  $P$  as  $P = QQ^T$  where the columns in  $Q$  are an orthonormal basis for the range of  $P$  so that  $Q^TQ = I$ . Thus the (euclidean) length of  $P\mathbf{x}$  is equal to the length of  $Q^T\mathbf{x}$ . The variance of  $Q^T\mathbf{x}$  can be written as  $E(\mathbf{x}^TQQ^T\mathbf{x})$  or  $E(\text{trace}(Q^T\mathbf{x}\mathbf{x}^TQ))$ , where  $E$  denotes the expectation (that is the integral or sum over the probability space). If we let  $R = E(\mathbf{x}\mathbf{x}^T)$  we see that we should choose  $Q$  in such a way that  $\text{trace}(Q^TRQ)$  is as large as possible. We have the following result.

**Lemma 26.** *Let  $R$  be a nonnegative definite symmetric  $d \times d$  matrix with eigenvalues  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$ . Then the maximum value of  $\text{trace}(Q^TRQ)$  where  $Q$  is a  $d \times m$  matrix with  $Q^TQ = I$  is  $\sum_{j=1}^m \lambda_j$  and it is achieved when the columns in  $Q$  form an orthonormal basis of the space spanned by the eigenvectors associated with  $\lambda_1, \dots, \lambda_m$ .*

Another way of looking at this result is to consider the case where one is given  $n$  vectors  $\mathbf{x}_j$ ,  $j = 1, \dots, n$  such that  $\sum_{j=1}^n \mathbf{x}_j = \mathbf{0}$ . Then one can form a  $d \times n$ -matrix  $A$  with the vectors  $\mathbf{x}_j$  as columns. If now we form the

singular value decomposition  $A = USV^T$ . In this case  $E(\mathbf{x}\mathbf{x}^T) = \frac{1}{n}AA^T = \frac{1}{n}USS^T U^T$ . From this it follows that  $Q$  consists of the first  $m$  columns of  $U$ .

Next we consider the case where the mean value of  $\mathbf{x}_j$  is not zero (or  $E(\mathbf{x}) \neq \mathbf{0}$ ). In that case we calculate  $\bar{\mathbf{x}} = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j$  and consider instead the vectors  $\mathbf{x}_j - \bar{\mathbf{x}}$ . If we formulate the problem in term of a minimization problem we have to calculate

$$\min_Q \sum_{j=1}^n \left| \mathbf{x}_j - Q(Q^T \mathbf{x}_j - Q\bar{\mathbf{x}}) - \bar{\mathbf{x}} \right|^2.$$

It turns out that we get the same result by calculating

$$\min_{Q, \mathbf{y}_j, \mathbf{c}} \sum_{j=1}^n \left| \mathbf{x}_j - Q\mathbf{y}_j - \mathbf{c} \right|^2,$$

that is, we get  $\mathbf{c} = \bar{\mathbf{x}}$  and  $\mathbf{y}_j = Q^T \mathbf{x}_j - Q\bar{\mathbf{x}}$ . (Clearly  $\mathbf{c}$  and  $\mathbf{y}_j$  are not uniquely determined since one can always add a constant to all  $\mathbf{y}_j$  and subtract it from  $\mathbf{c}$ .) One can even go further and formulate the following minimization problems:

$$\begin{aligned} \min_{f \in A(m,d), \mathbf{y}_j} \sum_{j=1}^2 \left| \mathbf{x}_j - f(\mathbf{y}_j) \right|^2 \\ \min_{g \in A(d,m), f \in A(d,m)} \sum_{j=1}^2 \left| \mathbf{x}_j - f(g(\mathbf{x}_j)) \right|^2, \end{aligned}$$

where  $A(k, n)$  is the set of all functions of the form  $\mathbf{z} \in \mathbb{R}^k \mapsto A\mathbf{z} + \mathbf{c} \in \mathbb{R}^n$  where  $A$  is an  $n \times m$  matrix.

## 2. Nonlinear principal components and neural networks

The theory of principal components briefly presented in the previous section does not really have anything to do neural networks, but every function in the set  $A(k, n)$  can, of course be realized as a neural network with linear activation function. This gives immediately the possibility of defining several kinds of "nonlinear principal components" if one gives up the requirement that the functions  $f$  and  $g$  are affine. Clearly, one cannot take arbitrary functions because then one could, using e.g. a space filling curve make the expression to be minimized equal to 0 in infinitely quite different ways so no information could be extracted. Neural networks, however, provide an easily parametrized class of functions that can be used for this purpose.

Here we just present an example of what one can get in a very simple case. In the picture below the dots are given data points. We use a neural

network with dimensions  $[2, 5, 1, 5, 2]$  and the line drawn in the picture is the graph of the mapping from level 2 to level 4, that is points in the plane are "projected" onto the real line and reconstructed from the line with this mapping.

