

Module 6: Principal components regression

6.1	The SVD factorization	1
6.2	The NIPALS algorithm for PCA	2
6.2.1	Iterations	2
6.2.2	Properties	2
6.3	How many components?	3
6.4	Principal components regression	4
6.4.1	Prediction	5
6.4.2	PCR and MLR	5

6.1 The SVD factorization

You're saying this only to make me go. [Ilsa Lund Laszlo, *Casablanca*, 1942]

In this module we turn to the Principal Components Regression (PCR) method, in which the PCA (Principal Components Analysis) method from the previous module is put to work in regression. To this end we consider the principal components of $\mathbf{X}^\top \mathbf{X}$, where \mathbf{X} is a centered $n \times k$ data matrix.

There are several ways of finding the principal components of the $\mathbf{X}^\top \mathbf{X}$ matrix. One possibility is to apply the SVD method to \mathbf{X} , writing the reduced form of SVD as follows:

$$\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{P}^\top,$$

where \mathbf{U} ($n \times r$) and \mathbf{P} ($k \times r$) are orthogonal matrices corresponding to r singular values, in the notation of Module 5.

Let the scores matrix be defined by

$$\mathbf{T} = \mathbf{U} \mathbf{D},$$

a matrix with orthogonal, but not necessarily orthonormal columns. In fact

$$\begin{aligned} \mathbf{T}^\top \mathbf{T} &= \mathbf{D} \mathbf{U}^\top \mathbf{U} \mathbf{D} \\ &= \mathbf{D}^2 \\ &= \mathbf{\Lambda}_r, \end{aligned}$$

where $\mathbf{\Lambda}_r = \text{diag}\{\lambda_1, \dots, \lambda_r\}$ contains the non-zero eigenvalues of $\mathbf{X}^\top \mathbf{X}$ in its diagonal. We assume that the eigenvalues are in decreasing order, $\lambda_1 \geq \dots \geq \lambda_r > 0$.

Since

$$\mathbf{X} = \mathbf{T}\mathbf{P}^\top \quad (6.1)$$

we find that

$$\begin{aligned} \mathbf{X}^\top \mathbf{X} &= \mathbf{P}\mathbf{T}^\top \mathbf{T}\mathbf{P}^\top \\ &= \mathbf{P}\mathbf{\Lambda}_r \mathbf{P}^\top, \end{aligned}$$

which is the spectral decomposition for $\mathbf{X}^\top \mathbf{X}$, except that columns of \mathbf{P} corresponding to zero eigenvalues have been left out. By using that \mathbf{P} is orthogonal, we may also write (6.1) as follows:

$$\mathbf{T} = \mathbf{X}\mathbf{P}, \quad (6.2)$$

which follows by noting that $\mathbf{X}\mathbf{P} = \mathbf{T}\mathbf{P}^\top \mathbf{P} = \mathbf{T}$. Recall from Module 5 that the columns of \mathbf{T} are known as *scores*, and those of \mathbf{P} as *loadings*.

6.2 The NIPALS algorithm for PCA

Now we consider the NIPALS (Nonlinear Iterative Partial Least Squares) algorithm for finding the principal components of $\mathbf{X}^\top \mathbf{X}$. We want to find the first g principal components of $\mathbf{X}^\top \mathbf{X}$, starting with the largest eigenvalue λ_1 and down. g must be less than or equal to r .

6.2.1 Iterations

The NIPALS algorithm starts with the initialization $j = 1$ and $\mathbf{X}_1 = \mathbf{X}$. The algorithm then iterates through the following steps:

1. Choose \mathbf{t}_j as any column of \mathbf{X}_j .
2. Let $\mathbf{p}_j = \mathbf{X}_j^\top \mathbf{t}_j / \|\mathbf{X}_j^\top \mathbf{t}_j\|$.
3. Let $\mathbf{t}_j = \mathbf{X}_j \mathbf{p}_j$.
4. If \mathbf{t}_j is unchanged continue; otherwise return to Step 2.
5. Let $\mathbf{X}_{j+1} = \mathbf{X}_j - \mathbf{t}_j \mathbf{p}_j^\top$.
6. Stop if $j = g$; otherwise let $j = j + 1$ and return to Step 1.

Assume first that $g = r$, so we have found all the principal components. Now form the matrices \mathbf{T} and \mathbf{P} with columns \mathbf{t}_j and \mathbf{p}_j , respectively; these matrices now satisfy (6.1).

It is possible to modify the NIPALS algorithm to take missing data into account, see Bro (1996), pp. 43–44.

6.2.2 Properties

Let us consider some properties of the NIPALS algorithm, which also help understand the PCA method.

That the NIPALS algorithm gives PCA may be seen as follows. Let $\lambda_j = \|\mathbf{X}^\top \mathbf{t}_j\|$ and write Step 2 as follows:

$$\mathbf{X}^\top \mathbf{t}_j = \lambda_j \mathbf{p}_j.$$

Now insert $\mathbf{t}_j = \mathbf{X} \mathbf{p}_j$ from Step 3, giving

$$\mathbf{X}^\top \mathbf{X} \mathbf{p}_j = \lambda_j \mathbf{p}_j. \quad (6.3)$$

This equation is satisfied upon convergence of the loop 2–4. This shows that λ_j and \mathbf{p}_j are an eigenvalue and eigenvector of $\mathbf{X}^\top \mathbf{X}$, respectively. Also note that using $\mathbf{t}_j = \mathbf{X} \mathbf{p}_j$ and (6.3) we obtain

$$\begin{aligned} \mathbf{t}_j^\top \mathbf{t}_j &= \mathbf{p}_j^\top \mathbf{X}^\top \mathbf{X} \mathbf{p}_j \\ &= \mathbf{p}_j^\top (\mathbf{X}^\top \mathbf{X} \mathbf{p}_j) \\ &= \lambda_j \mathbf{p}_j^\top \mathbf{p}_j \\ &= \lambda_j, \end{aligned} \quad (6.4)$$

where in the last step we have used the fact that \mathbf{p}_j is a unit vector (see Step 2).

After the first run through the loop 1–5, Step 5 with $j = 1$ gives that

$$\mathbf{X} = \mathbf{t}_1 \mathbf{p}_1^\top + \mathbf{X}_2. \quad (6.5)$$

Let us show that \mathbf{t}_1 and $\mathbf{X}_2 = \mathbf{X} - \mathbf{t}_1 \mathbf{p}_1^\top$ are orthogonal. In fact

$$\begin{aligned} (\mathbf{X} - \mathbf{t}_1 \mathbf{p}_1^\top)^\top \mathbf{t}_1 &= \mathbf{X}^\top \mathbf{t}_1 - \mathbf{p}_1^\top \mathbf{t}_1^\top \mathbf{t}_1 \\ &= \mathbf{X}^\top \mathbf{X} \mathbf{p}_1 - \mathbf{p}_1^\top \lambda_1 \\ &= \mathbf{0}, \end{aligned}$$

as seen from (6.3) with $j = 1$. Since \mathbf{t}_2 was initially picked as a column of \mathbf{X}_2 , it is hence orthogonal to \mathbf{t}_1 and remains so to the end of the loop.

After the second run through the loop 1–5, we obtain

$$\mathbf{X} = \mathbf{t}_1 \mathbf{p}_1^\top + \mathbf{t}_2 \mathbf{p}_2^\top + \mathbf{X}_3.$$

After g runs through the loop, we similarly have

$$\mathbf{X} = \mathbf{t}_1 \mathbf{p}_1^\top + \mathbf{t}_2 \mathbf{p}_2^\top + \cdots + \mathbf{t}_g \mathbf{p}_g^\top + \mathbf{X}_{g+1}, \quad (6.6)$$

where $\mathbf{X}_{g+1} = \mathbf{0}$ in the case $g = r$ (compare with (6.1)).

6.3 How many components?

As the example suggests, the essence of the PCA method is to decompose the \mathbf{X} matrix as in (6.6),

$$\begin{aligned}\mathbf{X} &= \mathbf{t}_1 \mathbf{p}_1^\top + \mathbf{t}_2 \mathbf{p}_2^\top + \cdots + \mathbf{t}_g \mathbf{p}_g^\top + \mathbf{X}_{g+1} \\ &= \mathbf{T}_g \mathbf{P}_g^\top + \mathbf{X}_{g+1},\end{aligned}\tag{6.7}$$

say, where \mathbf{T}_g and \mathbf{P}_g contain the first g columns of \mathbf{T} and \mathbf{P} , respectively. We want to choose g in such a way that \mathbf{X}_{g+1} is small and represents only noise, while the term $\mathbf{T}_g \mathbf{P}_g^\top$ represents the salient features of \mathbf{X} . In order to accomplish this, g must be chosen in such a way that the $r - g$ terms that are ignored correspond to zero or negligible eigenvalues.

In order to help rationalize the choice of g , the relative size of the eigenvalues are expressed as a percentage of the sum of all eigenvalues,

$$\frac{\lambda_1}{\lambda_1 + \cdots + \lambda_r} \times 100$$

and this percentage is interpreted as the *percent variation explained* by the corresponding principal component. Often, the cumulated percentages are used, so that the percent variation explained by the first g components is

$$\frac{\lambda_1 + \cdots + \lambda_g}{\lambda_1 + \cdots + \lambda_r} \times 100$$

As a rule, g should be chosen so that at least about 80–90 percent of the variation is explained.

The justification for the above terminology is that the variance of the score vector \mathbf{t}_j is

$$\begin{aligned}s^2(\mathbf{t}_j) &= \frac{1}{n-1} \|\mathbf{t}_j\|^2 \\ &= \frac{1}{n-1} \mathbf{t}_j^\top \mathbf{t}_j \\ &= \frac{1}{n-1} \lambda_j,\end{aligned}$$

so that λ_j is proportional to the variance of the corresponding score. In particular, all components with $\lambda_j = 0$ should be left out. Also, since the covariance matrix of \mathbf{X} is

$$\begin{aligned}\mathbf{V}_X &= \frac{1}{n-1} \mathbf{X}^\top \mathbf{X} \\ &= \frac{1}{n-1} \sum_{j=1}^k \lambda_j \mathbf{p}_j \mathbf{p}_j^\top,\end{aligned}$$

we may interpret λ_j as the contribution of \mathbf{t}_j to the total (co-)variance for \mathbf{X} . Note also that the sum of the eigenvalues is equal to

$$\text{tr}\{\mathbf{X}^\top \mathbf{X}\} = (n-1) \{s^2(\mathbf{x}_1) + \cdots + s^2(\mathbf{x}_k)\},$$

which is interpreted as the total variance in \mathbf{X} .

6.4 Principal components regression

The basic idea in Principal Components Regression (PCR) is that after choosing a suitable value for g in (6.7), the important features of \mathbf{X} have been retained by \mathbf{T}_g . We then perform the MLR with \mathbf{T}_g in place of \mathbf{X} for an $n \times m$ calibration data matrix \mathbf{Y} ,

$$\mathbf{Y} = \mathbf{T}_g \mathbf{C} + \mathbf{F}. \quad (6.8)$$

The least squares method then gives

$$\hat{\mathbf{C}} = \left(\mathbf{T}_g^\top \mathbf{T}_g \right)^{-1} \mathbf{T}_g^\top \mathbf{Y},$$

where $\mathbf{T}_g^\top \mathbf{T}_g$, being diagonal, is easy to invert. The fact that we have left out the loadings matrix \mathbf{P}_g in (6.8) is of no consequence for prediction, because the scores \mathbf{t}_j are linear combinations of the columns of \mathbf{X} , and the PCR method amounts to singling out those linear combinations that are best for predicting \mathbf{Y} .

6.4.1 Prediction

For prediction with PCR, it is necessary to turn to \mathbf{X} again, and using (6.2) we may write the regression equation as follows:

$$\begin{aligned} \mathbf{Y} &= \mathbf{T}_g \mathbf{C} + \mathbf{F}. \\ &= \mathbf{X} \mathbf{P}_g \mathbf{C} + \mathbf{F}. \end{aligned} \quad (6.9)$$

Consider a new sample spectrum \mathbf{z} and predicted value $\hat{\mathbf{y}}$ (both uncentered), and let $\bar{\mathbf{x}}$ and $\bar{\mathbf{y}}$ be the calibration sample averages. Then the prediction takes the form

$$\hat{\mathbf{y}} = \bar{\mathbf{y}} + (\mathbf{z} - \bar{\mathbf{x}}) \mathbf{P}_g \hat{\mathbf{C}}.$$

The matrix $\mathbf{P}_g \hat{\mathbf{C}}$ is called the *regression matrix*, and may be compared with the $\hat{\mathbf{B}}$ matrix of MLR.

6.4.2 PCR and MLR

Just as in MLR, the same prediction would be obtained if the columns of \mathbf{Y} were considered separately. The fact that \mathbf{P}_g appears in the PCR Equation (6.9) may be seen as compensating for the fact that we left out \mathbf{P}_g in (6.8). When comparing with MLR, the role of the \mathbf{B} matrix is now played by $\mathbf{P}_g \mathbf{C}$.

In the case where \mathbf{X} has rank k and $g = k$, the two methods will give identical results. When $g < k$, and still \mathbf{X} has rank k , the results of PCR may differ somewhat from those of MLR, depending on the number and sizes of the components left out.

However, PCR has some major advantages over MLR, in that \mathbf{X} may be singular, and the case $k \geq n$ may be dealt with. Note, however, that in the latter case $\lambda_n = \dots = \lambda_k = 0$, and so at most $g = n - 1$ components may be included.