

The Singular Value Decomposition

The singular value decomposition (henceforth SVD) of an $m \times n$ matrix \mathbf{X} , with $m \geq n$, can be expressed in the form $\mathbf{X} = \mathbf{U}\mathbf{W}\mathbf{V}^\top$. Here \mathbf{U} is $m \times n$, like \mathbf{X} , \mathbf{W} is a diagonal $n \times n$ matrix, and \mathbf{V} is a $n \times n$ orthogonal matrix, that is, such that $\mathbf{V}^\top = \mathbf{V}^{-1}$ or, equivalently, $\mathbf{V}\mathbf{V}^\top = \mathbf{V}^\top\mathbf{V} = \mathbf{I}_n$, the $n \times n$ identity matrix. In addition, the columns of \mathbf{U} are orthonormal, which means that $\mathbf{U}^\top\mathbf{U} = \mathbf{I}_n$.

The first question to be asked is whether the SVD exists for an arbitrary $m \times n$ matrix \mathbf{X} with $m \geq n$. Suppose that it does. Then it follows that

$$\mathbf{X}^\top\mathbf{X} = \mathbf{V}\mathbf{W}\mathbf{U}^\top\mathbf{U}\mathbf{W}\mathbf{V}^\top = \mathbf{V}\mathbf{W}^2\mathbf{V}^\top.$$

On postmultiplying by \mathbf{V} , we get $\mathbf{X}^\top\mathbf{X}\mathbf{V} = \mathbf{V}\mathbf{W}^2$. This shows that the elements of the diagonal matrix \mathbf{W}^2 are the eigenvalues of $\mathbf{X}^\top\mathbf{X}$, which must be non-negative since $\mathbf{X}^\top\mathbf{X}$ is by construction positive semi-definite, and that the columns of \mathbf{V} are the corresponding eigenvectors. This establishes the existence of \mathbf{W} and \mathbf{V} , although they are by no means unique. But given some choice of \mathbf{W} and \mathbf{V} , then, if \mathbf{U} exists, postmultiplying the defining equation by \mathbf{V} gives $\mathbf{X}\mathbf{V} = \mathbf{U}\mathbf{W}$. Then, if there are no zero elements on the diagonal of \mathbf{W} , we see that $\mathbf{U} = \mathbf{X}\mathbf{V}\mathbf{W}^{-1}$, and so in this case \mathbf{U} exists, and does indeed satisfy the requirement that $\mathbf{U}^\top\mathbf{U} = \mathbf{I}_n$, because

$$\mathbf{U}^\top\mathbf{U} = \mathbf{W}^{-1}\mathbf{V}^\top\mathbf{X}^\top\mathbf{X}\mathbf{V}\mathbf{W}^{-1} = \mathbf{W}^{-1}\mathbf{V}^\top\mathbf{V}\mathbf{W}^2\mathbf{W}^{-1},$$

since $\mathbf{X}^\top\mathbf{X}\mathbf{V} = \mathbf{V}\mathbf{W}^2$. Then, since $\mathbf{V}^\top\mathbf{V} = \mathbf{I}_n$, we see that $\mathbf{U}^\top\mathbf{U} = \mathbf{W}^{-1}\mathbf{W}^2\mathbf{W}^{-1} = \mathbf{I}_n$.

In fact, the eigenvalues of $\mathbf{X}^\top\mathbf{X}$ are strictly positive if (and only if) \mathbf{X} is of full rank. If not, this means that there exists a nonzero n -vector \mathbf{v} such that $\mathbf{X}\mathbf{v} = \mathbf{0}$. Then $\mathbf{v}^\top\mathbf{X}^\top\mathbf{X}\mathbf{v} = 0$, which implies that the positive semi-definite matrix $\mathbf{X}^\top\mathbf{X}$ is not positive definite, and so it has at least one zero eigenvalue. Conversely, if there is a zero eigenvalue of $\mathbf{X}^\top\mathbf{X}$, with corresponding (nonzero) eigenvector \mathbf{v} say, then $\mathbf{v}^\top\mathbf{X}^\top\mathbf{X}\mathbf{v} = \|\mathbf{X}\mathbf{v}\|^2 = 0$, which implies that $\mathbf{X}\mathbf{v} = \mathbf{0}$, so that \mathbf{X} does not have full rank.

Let us suppose that \mathbf{X} has full rank. We see immediately that the three matrices of the SVD exist. If one requires the elements of \mathbf{W} to be the positive square roots of the eigenvalues of $\mathbf{X}^\top\mathbf{X}$, and to be sorted in decreasing order from top left to bottom right, then the SVD is essentially unique (up to the signs of the columns of \mathbf{U} and \mathbf{V}), except in the non-generic case in which there is degeneracy of the eigenvalues.

The matrix $\mathbf{X}\mathbf{X}^\top$ is of dimension $m \times m$, and it given by

$$\mathbf{X}\mathbf{X}^\top = \mathbf{U}\mathbf{W}\mathbf{V}^\top\mathbf{V}\mathbf{W}\mathbf{U}^\top = \mathbf{U}\mathbf{W}^2\mathbf{U}^\top$$

Thus $\mathbf{X}\mathbf{X}^\top\mathbf{U} = \mathbf{U}\mathbf{W}^2$, which shows that the n diagonal elements of \mathbf{W}^2 are eigenvalues of $\mathbf{X}\mathbf{X}^\top$, as well as of $\mathbf{X}^\top\mathbf{X}$, while the columns of \mathbf{U} are the

corresponding eigenvectors. If $n < m$, there are also $m - n$ zero eigenvalues, the corresponding eigenspace being the complement in E^m of the span of the columns of \mathbf{U} , or equivalently those of \mathbf{X} . Note that $\mathbf{U}\mathbf{U}^\top$ is an orthogonal projection matrix:

$$(\mathbf{U}\mathbf{U}^\top)^2 = \mathbf{U}\mathbf{U}^\top\mathbf{U}\mathbf{U}^\top = \mathbf{U}\mathbf{I}_k\mathbf{U}^\top = \mathbf{U}\mathbf{U}^\top,$$

which establishes idempotency, symmetry being obvious. In fact, $\mathbf{U}\mathbf{U}^\top = \mathbf{P}_\mathbf{X}$, the orthogonal projection on to the span of the columns of \mathbf{X} , since the columns of \mathbf{U} and \mathbf{X} span the same space. The columns of \mathbf{U} therefore provide an orthonormal basis for this span. This suggests that the SVD can provide a way to implement ordinary least squares, and indeed it is one of the most numerically stable ways to do so.

The columns of \mathbf{U} are called the left singular vectors of \mathbf{X} , and those of \mathbf{V} the right singular vectors of \mathbf{X} .

Generalised Inverse of a Matrix

A generalised inverse of an $m \times n$ matrix \mathbf{X} is an $n \times m$ matrix \mathbf{X}^+ that satisfies the two relations

$$\mathbf{X}\mathbf{X}^+\mathbf{X} = \mathbf{X}, \quad \text{and} \quad \mathbf{X}^+\mathbf{X}\mathbf{X}^+ = \mathbf{X}^+.$$

If \mathbf{X} is a non-singular square matrix, then we can pre- and post-multiply the first of these requirements by \mathbf{X}^{-1} to get $\mathbf{X}^+ = \mathbf{X}^{-1}$. This shows that the notion of a generalised inverse does indeed generalise that of the inverse, when the latter exists.

If $m > n$ and \mathbf{X} has full column rank n , a possible generalised inverse is given by $(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}$. However, it is not unique. If the nonzero $m \times n$ matrix \mathbf{B} has columns that are orthogonal to those of \mathbf{X} , it is easy to check that $(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top + \mathbf{B}^\top$ is also a generalised inverse of \mathbf{X} . Provided that $m > n$, it is always possible to find such a matrix \mathbf{B} .

Note that, for any generalised inverse \mathbf{X}^+ , both $\mathbf{X}\mathbf{X}^+$ and $\mathbf{X}^+\mathbf{X}$ are idempotent, so that, if either of them is symmetric, it is an orthogonal projection matrix. We have the following result:

Moore-Penrose Generalised Inverse

For every $m \times n$ matrix \mathbf{X} , $m \geq n$, there exists a unique matrix \mathbf{X}^+ , called the Moore-Penrose generalised inverse, or the pseudo-inverse, satisfying the four conditions

$$\begin{aligned} \mathbf{X}\mathbf{X}^+\mathbf{X} &= \mathbf{X}, & \mathbf{X}^+\mathbf{X}\mathbf{X}^+ &= \mathbf{X}^+, \\ (\mathbf{X}\mathbf{X}^+)^\top &= \mathbf{X}\mathbf{X}^+, & (\mathbf{X}^+\mathbf{X})^\top &= \mathbf{X}^+\mathbf{X}. \end{aligned}$$

I omit the full proof.

However, the existence part can be shown constructively by using the SVD. The Moore-Penrose generalised inverse of \mathbf{X} is $\mathbf{X}^+ = \mathbf{V}\mathbf{W}^+\mathbf{U}^\top$, where \mathbf{W}^+ is a specific generalised inverse of the diagonal matrix \mathbf{W} . If \mathbf{W} has no zero singular values, then \mathbf{W}^+ is the usual inverse, the diagonal matrix the elements of which are the inverses of the singular values in \mathbf{W} . If there are singular values equal to zero, they occur in the bottom right corner of \mathbf{W} . Let the rank of \mathbf{W} be $p < n$, then we can write

$$\mathbf{W} = \begin{bmatrix} \mathbf{D}_{p \times p} & \mathbf{O}_{p \times (n-p)} \\ \mathbf{O}_{(n-p) \times p} & \mathbf{O}_{(n-p) \times (n-p)} \end{bmatrix},$$

where the matrix \mathbf{D} contains all the nonzero singular values, so that \mathbf{D}^{-1} is well defined. Then define

$$\mathbf{W}^+ = \begin{bmatrix} \mathbf{D}_{p \times p}^{-1} & \mathbf{O}_{p \times (n-p)} \\ \mathbf{O}_{(n-p) \times p} & \mathbf{O}_{(n-p) \times (n-p)} \end{bmatrix}.$$

We see immediately that \mathbf{W}^+ thus defined is indeed a generalised inverse of \mathbf{W} .

Calculate as follows:

$$\mathbf{X}^+\mathbf{X} = \mathbf{V}\mathbf{W}^+\mathbf{U}^\top\mathbf{U}\mathbf{W}\mathbf{V}^\top = \mathbf{V}\mathbf{W}^+\mathbf{W}\mathbf{V}^\top.$$

This is clearly symmetric, and it is easy to show that $\mathbf{X}\mathbf{X}^+$ is so as well. Then

$$\begin{aligned} \mathbf{X}\mathbf{X}^+\mathbf{X} &= \mathbf{U}\mathbf{W}\mathbf{V}^\top\mathbf{V}\mathbf{W}^+\mathbf{W}\mathbf{V}^\top = \mathbf{U}\mathbf{W}\mathbf{V}^\top = \mathbf{X}, \text{ and} \\ \mathbf{X}^+\mathbf{X}\mathbf{X}^+ &= \mathbf{V}\mathbf{W}^+\mathbf{W}\mathbf{V}^\top\mathbf{V}\mathbf{W}^+\mathbf{U}^\top = \mathbf{V}\mathbf{W}^+\mathbf{U}^\top = \mathbf{X}^+, \end{aligned}$$

where I have used the fact that \mathbf{W}^+ is a generalised inverse of \mathbf{W} . This shows that \mathbf{X}^+ satisfies the defining properties of the Moore-Penrose inverse. Only uniqueness remains to be proved.

The *Deep Learning* book claims that \mathbf{X}^+ can also be defined by the formula

$$\mathbf{X}^+ = \lim_{\alpha \downarrow 0} (\mathbf{X}^\top\mathbf{X} + \alpha\mathbf{I})^{-1}\mathbf{X}^\top$$

This formula is harder to implement than the one based on the SVD, as very efficient algorithms exist for the computation of the SVD. The right-hand side is easy enough to compute, but the limit is not.

Principal Components Analysis

Another use of the SVD is to implement Principal Components Analysis (PCA). PCA has been around for a long time as a tool for the analysis of data. For a given full-rank $m \times n$ matrix \mathbf{X} , an orthogonal linear transformation of the columns of \mathbf{X} is sought, such that the transformed matrix has mutually orthogonal columns, and has the property that the successive

columns have progressively decreasing variance. Before performing the transformation, the elements in each column of \mathbf{X} must be centred by subtracting the column means.

The desired transformation is given by $\mathbf{T} = \mathbf{X}\mathbf{V}$. We saw that $\mathbf{X}\mathbf{V} = \mathbf{U}\mathbf{W}$, and so $\mathbf{T}^\top\mathbf{T} = \mathbf{W}\mathbf{U}^\top\mathbf{U}\mathbf{W} = \mathbf{W}^2$, a diagonal matrix. This implies that the columns of \mathbf{T} are indeed mutually orthogonal. The relation between the columns of \mathbf{T} and those of \mathbf{U} can be written as $\mathbf{t}_i = w_i\mathbf{u}_i$, $i = 1, \dots, n$, where \mathbf{t}_i and \mathbf{u}_i are the i^{th} columns of \mathbf{T} and \mathbf{U} respectively, and w_i is the i^{th} diagonal element of \mathbf{W} . Thus $\mathbf{t}_i^\top\mathbf{t}_i = w_i^2\mathbf{u}_i^\top\mathbf{u}_i = w_i^2$ since the columns of \mathbf{U} have unit norm. Now we chose to sort the w_i in decreasing order, and so the norms of the columns of \mathbf{T} are also in decreasing order, as desired.

One purpose of PCA is dimension reduction. One seeks an $m \times n$ matrix \mathbf{Y} , of rank $p < n$, that approximates the original matrix \mathbf{X} as well as possible. In the formula $\mathbf{X} = \mathbf{U}\mathbf{W}\mathbf{V}^\top$, we can set the $n - p$ smallest diagonal elements of \mathbf{W} equal to zero, and denote the result by \mathbf{W}_p , which has only p nonzero diagonal elements, so that the approximation $\mathbf{Y} \equiv \mathbf{U}\mathbf{W}_p\mathbf{V}^\top$ has rank p . It can be shown that \mathbf{Y} is closer to \mathbf{X} than any other $m \times n$ matrix of rank p , where distance is measured by the sum of the squares of the elements of the difference between \mathbf{X} and a rank- p approximation to \mathbf{X} . This result is known as the Eckart-Young theorem. The name of Mirsky is sometimes added to the other two names, as Mirsky is responsible for an extension of the result.

Proof: The square of the **Frobenius norm** of an $m \times n$ matrix \mathbf{A} is the sum of the squares of the elements of \mathbf{A} . It is easy to check that this is equal to $\text{Tr}(\mathbf{A}^\top\mathbf{A}) = \text{Tr}(\mathbf{A}\mathbf{A}^\top)$.

Consider an $m \times n$ rank- p matrix \mathbf{Y} , and the Frobenius norm of the difference $\mathbf{X} - \mathbf{Y}$. We have $\mathbf{X} - \mathbf{Y} = \mathbf{X} - \mathbf{P}_\mathbf{X}\mathbf{Y} - \mathbf{M}_\mathbf{X}\mathbf{Y}$, and so

$$\text{Tr}(\mathbf{X} - \mathbf{Y})^\top(\mathbf{X} - \mathbf{Y}) = \text{Tr}(\mathbf{X} - \mathbf{P}_\mathbf{X}\mathbf{Y})^\top(\mathbf{X} - \mathbf{P}_\mathbf{X}\mathbf{Y}) + \text{Tr}(\mathbf{Y}^\top\mathbf{M}_\mathbf{X}\mathbf{Y}).$$

The second term on the right-hand side is non-negative, and so, if we wish to minimise the norm of $\mathbf{X} - \mathbf{Y}$, we should choose \mathbf{Y} such that $\mathbf{M}_\mathbf{X}\mathbf{Y} = \mathbf{O}$ and $\mathbf{P}_\mathbf{X}\mathbf{Y} = \mathbf{Y}$. Let us write $\mathbf{Y} = \mathbf{U}\mathbf{Z}$, where the $n \times n$ matrix \mathbf{Z} has rank p . Above, we had $\mathbf{Y} = \mathbf{U}\mathbf{W}_p\mathbf{V}^\top$, which implies that, in this case, $\mathbf{Z} = \mathbf{W}_p\mathbf{V}^\top$. Then $\mathbf{X} - \mathbf{Y} = \mathbf{U}(\mathbf{W} - \mathbf{W}_p)\mathbf{V}^\top$, and so, where $\|\cdot\|_F$ denotes the Frobenius norm,

$$\begin{aligned} \|\mathbf{X} - \mathbf{Y}\|_F^2 &= \text{Tr}(\mathbf{V}(\mathbf{W} - \mathbf{W}_p)\mathbf{U}^\top\mathbf{U}(\mathbf{W} - \mathbf{W}_p)\mathbf{V}^\top) \\ &= \text{Tr}(\mathbf{V}(\mathbf{W} - \mathbf{W}_p)^2\mathbf{V}^\top) = \text{Tr}((\mathbf{W} - \mathbf{W}_p)^2), \end{aligned} \quad (1)$$

where we have used the invariance of the trace of a matrix product under a cyclic permutation of the factors, and the fact that $\mathbf{V}^\top\mathbf{V} = \mathbf{I}$. Now the matrix $\mathbf{W} - \mathbf{W}_p$ is a diagonal matrix with only the last $n - p$ elements nonzero, these being equal to the $n - p$ smallest singular values. Thus

$$\|\mathbf{X} - \mathbf{Y}\|_F^2 = \sum_{i=n-p+1}^n w_i^2. \quad (2)$$

For the theorem, we have to show that, with any other choice of \mathbf{Y} and \mathbf{Z} , the squared norm of $\mathbf{X} - \mathbf{Y}$ is no less than the sum above.

In the general case, then, $\mathbf{X} - \mathbf{Y} = \mathbf{U}(\mathbf{W}\mathbf{V}\mathbf{V}^\top - \mathbf{Z})$. Let $\mathbf{Z} = \mathbf{A}\mathbf{V}^\top$, where the $n \times n$ matrix \mathbf{A} is given by $\mathbf{A} = \mathbf{Z}\mathbf{V}$. With this, $\mathbf{X} - \mathbf{Y} = \mathbf{U}(\mathbf{W} - \mathbf{A})\mathbf{V}^\top$. Similarly to the calculation in (1), we find that

$$\|\mathbf{X} - \mathbf{Y}\|_F^2 = \text{Tr}((\mathbf{W} - \mathbf{A}^\top)(\mathbf{W} - \mathbf{A})) = \text{Tr}((\mathbf{W} - \mathbf{A}^\top)\mathbf{P}_\mathbf{A}(\mathbf{W} - \mathbf{A})) + \text{Tr}(\mathbf{W}\mathbf{M}_\mathbf{A}\mathbf{W}), \blacksquare$$

where $\mathbf{P}_\mathbf{A}$ is the orthogonal projection in the n -dimensional space spanned by the columns of \mathbf{V} , and $\mathbf{M}_\mathbf{A}$ is the complementary projection, so that $\mathbf{M}_\mathbf{A}\mathbf{A} = \mathbf{O}$. Since \mathbf{W} is a diagonal matrix, it follows that

$$\|\mathbf{X} - \mathbf{Y}\|_F^2 \geq \text{Tr}(\mathbf{W}^2\mathbf{M}_\mathbf{A}) = \sum_{i=1}^n w_i^2 (\mathbf{M}_\mathbf{A})_{ii}, \quad (3)$$

where $(\mathbf{A})_{ii}$ is the i^{th} diagonal element of the orthogonal projection matrix $\mathbf{M}_\mathbf{A}$. These diagonal elements all lie between 0 and 1, and they sum to $n - p$, which is the dimension of the image of $\mathbf{M}_\mathbf{A}$. Comparison of the right-hand sides of (2) and (3) shows that the latter cannot be smaller than the former, given that the w_i are in decreasing in i . \blacksquare