

Singular Value Decomposition, Generalised Inverse, Principal Components

Consider a real $m \times n$ matrix \mathbf{A} , with $m \geq n$. Its **singular value decomposition** (SVD) takes the form

$$\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{V}^\top,$$

where \mathbf{U} is $m \times n$, \mathbf{D} is $n \times n$, and \mathbf{V} is $n \times n$. Further, \mathbf{U} has orthonormal columns, which means that $\mathbf{U}^\top\mathbf{U} = \mathbf{I}_n$, and $\mathbf{U}\mathbf{U}^\top$ is symmetric and idempotent, that is, an orthogonal projection matrix; \mathbf{D} is diagonal, with non-negative diagonal elements arranged in decreasing order down the main diagonal; \mathbf{V} is an orthogonal matrix, so that $\mathbf{V}^\top\mathbf{V} = \mathbf{V}\mathbf{V}^\top = \mathbf{I}_n$, and $\mathbf{V}^\top = \mathbf{V}^{-1}$. The SVD always exists if \mathbf{A} is real. The columns of \mathbf{U} are the **left singular vectors**, those of \mathbf{V} are the **right singular vectors**, and the elements of \mathbf{D} are the **singular values**.

The elements of \mathbf{D} are the square roots of the eigenvalues of $\mathbf{A}^\top\mathbf{A}$, and the right singular vectors are the eigenvectors:

$$\mathbf{A}^\top\mathbf{A} = \mathbf{V}\mathbf{D}\mathbf{U}^\top\mathbf{U}\mathbf{D}\mathbf{V}^\top = \mathbf{V}\mathbf{D}^2\mathbf{V}^\top,$$

whence

$$(\mathbf{A}^\top\mathbf{A})\mathbf{V} = \mathbf{V}\mathbf{D}^2.$$

Similarly, the squares of the elements of \mathbf{D} are some of the eigenvalues of $\mathbf{A}\mathbf{A}^\top$, with the left singular vectors as the corresponding eigenvectors:

$$(\mathbf{A}\mathbf{A}^\top)\mathbf{U} = \mathbf{U}\mathbf{D}^2.$$

The other eigenvalues of $\mathbf{A}\mathbf{A}^\top$ are all zero. Some of the bottom right diagonal elements of \mathbf{D} may also be zero: this happens when \mathbf{A} does not have full column rank. In this case, we have this setup:

$$\mathbf{D} = \begin{bmatrix} \mathbf{D}_1 & \mathbf{O} \\ \mathbf{O} & \mathbf{O} \end{bmatrix}. \tag{1}$$

If \mathbf{A} has full column rank of n , then \mathbf{D} is non-singular, and \mathbf{D}^{-1} exists. Otherwise, make the definition

$$\mathbf{D}^+ = \begin{bmatrix} \mathbf{D}_1^{-1} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} \end{bmatrix}.$$

We will see that \mathbf{D}^+ is a **generalised inverse**.

Any matrix \mathbf{A}^+ that satisfies the two requirements

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

is a generalised inverse of \mathbf{A} . If we define $\mathbf{A}^+ = \mathbf{V}\mathbf{D}^+\mathbf{U}^\top$, then we can see at once that this \mathbf{A}^+ is indeed a generalised inverse. If in addition the matrices

$\mathbf{A}\mathbf{A}^+$ and $\mathbf{A}^+\mathbf{A}$ are symmetric, we obtain the **Moore-Penrose** generalised inverse. This property is also satisfied with our definition of \mathbf{A}^+ .

Note that $\mathbf{A}\mathbf{A}^+$ and $\mathbf{A}^+\mathbf{A}$ are idempotent, where \mathbf{A}^+ is *any* generalised inverse of \mathbf{A} :

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

This means that these two matrices are projection matrices; with the Moore-Penrose inverse, they are orthogonal projection matrices.

It is easy to check that \mathbf{D}^+ also satisfies the requirements for it to be the Moore-Penrose generalised inverse of \mathbf{D} . Of course, if \mathbf{A} has full column rank, $\mathbf{D}^+ = \mathbf{D}^{-1}$.

The aim of **Principal Components Analysis** (PCA) is dimension reduction, with as little loss of information as possible. Think first of \mathbf{A} as a matrix of data, where each row is an instance, or example, and each column corresponds to a feature. There are too many features, and the aim is to replace them by a smaller number of linear combinations, while still retaining as much of the information as possible.

The idea is to set equal to zero some of the smallest singular values. Even if \mathbf{D} had no zero diagonal elements, this would replace the original \mathbf{D} by a matrix like (1), where we suppose that \mathbf{D}_1 is a $k \times k$ matrix, with $k < n$. Thus causes the last $n - k$ columns of \mathbf{U} to have no importance, as also the last $n - k$ rows of \mathbf{V} . Let

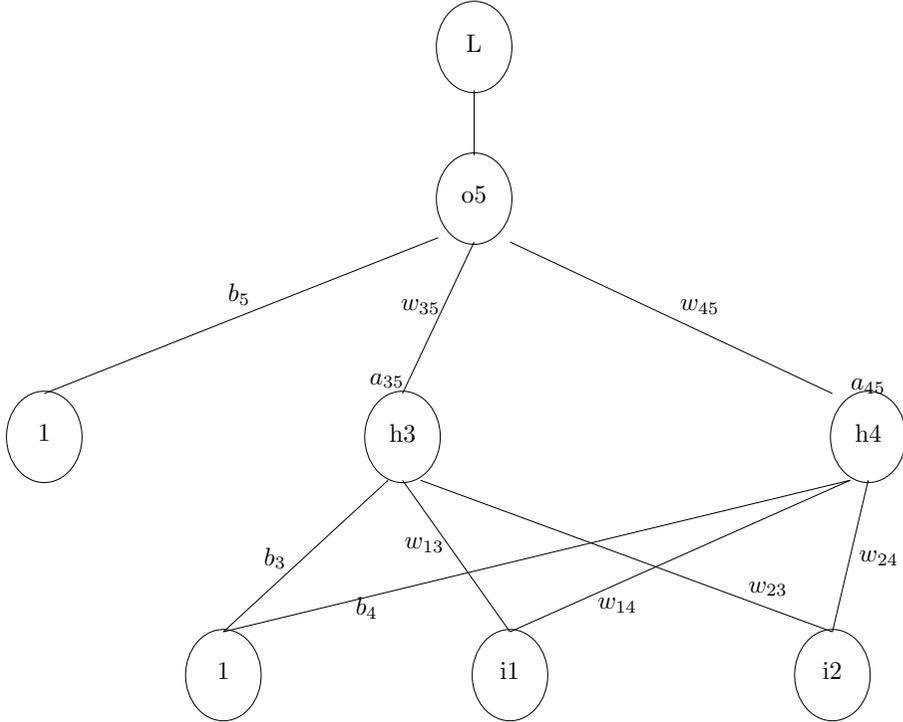
$$\mathbf{U} = [\mathbf{U}_1 \quad \mathbf{O}] \quad \text{and} \quad \mathbf{V}^\top = \begin{bmatrix} \mathbf{V}_1^\top \\ \mathbf{O} \end{bmatrix}.$$

Then $\mathbf{B} \equiv \mathbf{U}_1\mathbf{D}_1\mathbf{V}_1^\top$ is an $m \times k$ matrix, with $k < n$.

The name of principal components analysis could in principle be applied to the full matrix \mathbf{U} , which is of the same dimensions as \mathbf{A} itself, with the columns of \mathbf{U} being the principal components, or orthonormalised linear combinations of those of \mathbf{A} , ordered by how much of the total variation of \mathbf{A} they account for, but normally the analysis is used to trim away the directions of only small variation.

Back-Propagation

Below is a representation of a feedforward network with an input layer containing two units and a constant (bias), a hidden layer with two units and another constant, an output layer with just one unit, and a loss function that depends on the content of that unit .



The inputs $i1$ and $i2$ are fed in, and then, for one iteration, the contents of the other units are computed as follows:

$$\begin{aligned}
 h3 &= b_3 + w_{13}i1 + w_{23}i2 \\
 h4 &= b_4 + w_{14}i1 + w_{24}i2 \\
 o5 &= b_5 + w_{35}a_{35}(h3) + w_{45}a_{45}(h4) \\
 \text{loss} &= L(o5)
 \end{aligned} \tag{2}$$

The **parameters** are the biases b_3 , b_4 , and b_5 , and the weights w_{13} , w_{14} , w_{23} , w_{24} , w_{35} , and w_{45} ; a_{35} and a_{45} are the **activation functions**.

The feedforward step involves the computations (2), the results of which are stored, as also the values of the activation functions $a_{35}(h3)$ and $a_{45}(h4)$. Also stored are the derivatives $a'_{35}(h3)$, $a'_{45}(h4)$, and $L'(o5)$, which are also stored in the corresponding units.

Next comes back-propagation, which involves computing the partial derivatives of the loss function with respects to all of the parameters. This starts with the parameters that take us from the hidden layer to the output layer, namely b_5 , w_{35} , and w_{45} . We have

$$\begin{aligned}
 \frac{\partial L}{\partial b_5} &= L'(o5) \frac{\partial o5}{\partial b_5} = L'(o5); & \frac{\partial L}{\partial w_{35}} &= L'(o5) \frac{\partial o5}{\partial w_{35}} = L'(o5) a_{35}(h3); & \text{and} \\
 \frac{\partial L}{\partial w_{45}} &= L'(o5) \frac{\partial o5}{\partial w_{45}} = L'(o5) a_{45}(h4)
 \end{aligned}$$

Notice that everything needed for computing these derivatives has been stored. To proceed further, we need the derivatives with respect to $h3$ and $h4$. These are

$$\frac{\partial L}{\partial h3} = L'(o5) w_{35} a'_{35}(h3) \quad \text{and} \quad \frac{\partial L}{\partial h4} = L'(o5) w_{45} a'_{45}(h4)$$

But these are just intermediary calculations for the derivatives with respect to the parameters that take us from the input layer to the hidden layer. For these, we have

$$\begin{aligned}
 \frac{\partial L}{\partial b_3} &= \frac{\partial L}{\partial h3} \frac{\partial h3}{\partial b_3} = L'(o5) w_{35} a'_{35}(h3); \\
 \frac{\partial L}{\partial w_{13}} &= \frac{\partial L}{\partial h3} \frac{\partial h3}{\partial w_{13}} = L'(o5) w_{35} a'_{35}(h3) i1 \\
 \frac{\partial L}{\partial w_{14}} &= \frac{\partial L}{\partial h3} \frac{\partial h3}{\partial w_{14}} = L'(o5) w_{35} a'_{35}(h3) i2
 \end{aligned}$$

Again, everything needed was stored during the feedforward step. The derivatives with respect to b_4 , w_{14} , and w_{24} are computed in exactly the same way.

We see that what we have done is indeed *back*-propagation. We start with the derivative of the loss function with respect to the outputs of the output layer (here there is only one, $o5$), and go back down, systematically applying the chain rule, for which everything is available when needed.

What we have computed is the **gradient** of the loss function, and so we can move to the next iteration by updating all the parameters by, for instance, stochastic gradient descent.