

2.12 Example: Principal Components Analysis

One simple machine learning algorithm, **principal components analysis** or PCA can be derived using only knowledge of basic linear algebra.

Suppose we have a collection of m points $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ in \mathbb{R}^n . Suppose we would like to apply lossy compression to these points. Lossy compression means storing the points in a way that requires less memory but may lose some precision. We would like to lose as little precision as possible.

One way we can encode these points is to represent a lower-dimensional version of them. For each point $\mathbf{x}^{(i)} \in \mathbb{R}^n$ we will find a corresponding code vector $\mathbf{c}^{(i)} \in \mathbb{R}^l$. If l is smaller than n , it will take less memory to store the code points than the original data. We will want to find some encoding function that produces the code for an input, $f(\mathbf{x}) = \mathbf{c}$, and a decoding function that produces the reconstructed input given its code, $\mathbf{x} \approx g(f(\mathbf{x}))$.

PCA is defined by our choice of the decoding function. Specifically, to make the decoder very simple, we choose to use matrix multiplication to map the code back into \mathbb{R}^n . Let $g(\mathbf{c}) = \mathbf{D}\mathbf{c}$, where $\mathbf{D} \in \mathbb{R}^{n \times l}$ is the matrix defining the decoding.

Computing the optimal code for this decoder could be a difficult problem. To keep the encoding problem easy, PCA constrains the columns of \mathbf{D} to be orthogonal to each other. (Note that \mathbf{D} is still not technically “an orthogonal matrix” unless $l = n$)

With the problem as described so far, many solutions are possible, because we can increase the scale of $\mathbf{D}_{:,i}$ if we decrease c_i proportionally for all points. To give the problem a unique solution, we constrain all of the columns of \mathbf{D} to have unit norm.

In order to turn this basic idea into an algorithm we can implement, the first thing we need to do is figure out how to generate the optimal code point \mathbf{c}^* for each input point \mathbf{x} . One way to do this is to minimize the distance between the input point \mathbf{x} and its reconstruction, $g(\mathbf{c}^*)$. We can measure this distance using a norm. In the principal components algorithm, we use the L^2 norm:

$$\mathbf{c}^* = \arg \min_{\mathbf{c}} \|\mathbf{x} - g(\mathbf{c})\|_2. \quad (2.54)$$

We can switch to the squared L^2 norm instead of the L^2 norm itself, because both are minimized by the same value of \mathbf{c} . Both are minimized by the same value of \mathbf{c} because the L^2 norm is non-negative and the squaring operation is