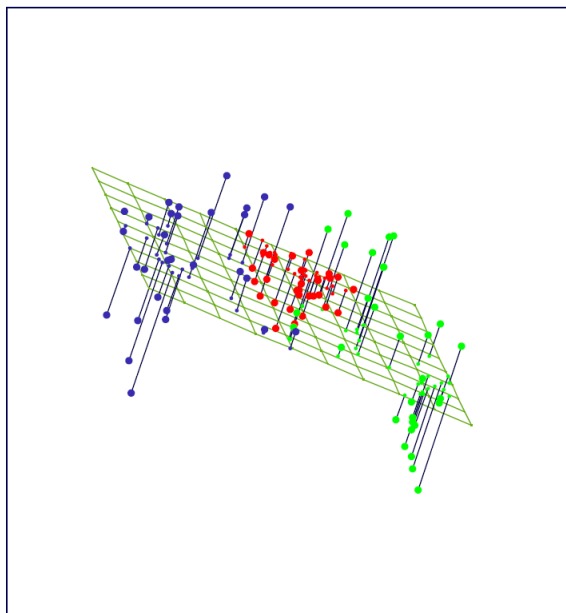


Principal Components Analysis

Principal Components Analysis (PCA) is a standard technique for **dimensionality reduction** of data sets. It is a way to automatically find a **subspace** which approximates the data. It is used everywhere in signal processing, machine learning, and statistics.

There are actually two equivalent ways to think about PCA. The first is statistical: we are trying to find a transform that is carefully tuned to the (second-order) statistics of the data. The second perspective, which is what we will adopt in this course, is more geometrical: given a set of vectors, we are trying to find a subspace of a certain dimension that comes closest to containing this set.

Specifically, suppose that we have data points $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^D$, and want to find the K -dimensional affine space (subspace plus offset) that comes closest to containing them. Here is a picture¹



¹From Ch. 14 of Tibshirani and Hastie's *Elements of Statistical Learning*.

Our goal is to find an offset $\boldsymbol{\mu} \in \mathbb{R}^D$ and a matrix \mathbf{Q} with orthonormal columns such that

$$\mathbf{x}_n \approx \boldsymbol{\mu} + \mathbf{Q}\boldsymbol{\theta}_n \quad \text{for all } n = 1, \dots, N,$$

for some $\boldsymbol{\theta}_n \in \mathbb{R}^K$. We cast this as the following optimization problem. Given $\mathbf{x}_1, \dots, \mathbf{x}_N$, solve

$$\underset{\boldsymbol{\mu}, \mathbf{Q}, \{\boldsymbol{\theta}_n\}}{\text{minimize}} \quad \sum_{n=1}^N \|\mathbf{x}_n - \boldsymbol{\mu} - \mathbf{Q}\boldsymbol{\theta}_n\|_2^2 \quad \text{subject to} \quad \mathbf{Q}^T \mathbf{Q} = \mathbf{I}.$$

Note that if we fix $\boldsymbol{\mu}$ and define $\tilde{\mathbf{x}}_n = \mathbf{x}_n - \boldsymbol{\mu}$, then we can recast the optimization with respect to \mathbf{Q} and the $\boldsymbol{\theta}_n$ as

$$\underset{\mathbf{Q}, \{\boldsymbol{\theta}_n\}}{\text{minimize}} \quad \sum_{n=1}^N \|\tilde{\mathbf{x}}_n - \mathbf{Q}\boldsymbol{\theta}_n\|_2^2 \quad \text{subject to} \quad \mathbf{Q}^T \mathbf{Q} = \mathbf{I}.$$

If $\tilde{\mathbf{X}}$ and $\boldsymbol{\Theta}$ denote the matrices whose columns are given by $\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_N$ and $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N$ respectively, then we can also write this as

$$\underset{\substack{\mathbf{Q}: D \times K \\ \boldsymbol{\Theta}: K \times N}}{\text{minimize}} \quad \|\tilde{\mathbf{X}} - \mathbf{Q}\boldsymbol{\Theta}\|_F^2 \quad \text{subject to} \quad \mathbf{Q}^T \mathbf{Q} = \mathbf{I}.$$

This is exactly the optimization problem that we looked at previously in our Subspace Approximation Lemma! Thus the solution is given by computing the SVD of $\tilde{\mathbf{X}} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}$ and then taking as our solution

$$\begin{aligned} \hat{\mathbf{Q}} &= \mathbf{U}_K, \\ \hat{\boldsymbol{\Theta}} &= \mathbf{U}_K^T \tilde{\mathbf{X}}, \end{aligned}$$

where $\mathbf{U}_K = [\mathbf{u}_1 \ \mathbf{u}_2 \ \cdots \ \mathbf{u}_K]$ contains the first K columns of \mathbf{U} .

Finally, let us return to the question of how to set $\boldsymbol{\mu}$. For any given $\boldsymbol{\mu}$, the solution for \mathbf{Q} and $\boldsymbol{\Theta}$ is given by the Subspace Approximation Lemma. This results in setting

$$\boldsymbol{\theta}_n = \mathbf{Q}^T(\mathbf{x}_n - \boldsymbol{\mu}).$$

Plugging this in for $\boldsymbol{\theta}_n$ in our objective function, we have that

$$\begin{aligned}\mathbf{x}_n - \boldsymbol{\mu} - \mathbf{Q}\boldsymbol{\theta}_n &= \mathbf{x}_n - \boldsymbol{\mu} - \mathbf{Q}\mathbf{Q}^T(\mathbf{x}_n - \boldsymbol{\mu}) \\ &= (\mathbf{I} - \mathbf{Q}\mathbf{Q}^T)(\mathbf{x}_n - \boldsymbol{\mu}).\end{aligned}$$

Hence, the problem of selecting $\boldsymbol{\mu}$ reduces to the optimization problem

$$\underset{\boldsymbol{\mu}}{\text{minimize}} \quad \sum_{n=1}^N \|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^T)(\mathbf{x}_n - \boldsymbol{\mu})\|_2^2$$

The vector $\boldsymbol{\mu}$ is unconstrained; we can solve for the optimal $\boldsymbol{\mu}$ by taking a gradient and setting it equal to zero. To make this easier, note that

$$\|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^T)(\mathbf{x}_n - \boldsymbol{\mu})\|_2^2 = (\mathbf{x}_n - \boldsymbol{\mu})^T(\mathbf{I} - \mathbf{Q}\mathbf{Q}^T)(\mathbf{x}_n - \boldsymbol{\mu})$$

by simply expanding out the norm squared as an inner product and then using the fact that $\mathbf{I} - \mathbf{Q}\mathbf{Q}^T$ is a projector, i.e., it is symmetric and $(\mathbf{I} - \mathbf{Q}\mathbf{Q}^T)^2 = \mathbf{I} - \mathbf{Q}\mathbf{Q}^T$. Thus, by taking a gradient and setting it equal to zero we have

$$\begin{aligned}\mathbf{0} &= -2 \sum_{n=1}^N (\mathbf{I} - \mathbf{Q}\mathbf{Q}^T)(\mathbf{x}_n - \boldsymbol{\mu}) \\ &= -2(\mathbf{I} - \mathbf{Q}\mathbf{Q}^T) \left(\left(\sum_{n=1}^N \mathbf{x}_n \right) - N\boldsymbol{\mu} \right).\end{aligned}$$

We can satisfy this condition by taking the offset $\boldsymbol{\mu}$ to be the sample mean (average of all the observed vectors):

$$\hat{\boldsymbol{\mu}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n.$$

Note that this choice is not unique – any choice of $\boldsymbol{\mu}$ that results in $\sum(\mathbf{x}_n - \boldsymbol{\mu})$ living in the nullspace of $\mathbf{I} - \mathbf{Q}\mathbf{Q}^T$ would also suffice – but $\boldsymbol{\mu}$ is the easy and obvious choice, and also what is usually done in practice, because it makes computing the solution to the PCA problem straightforward.

Computing the PCA solution

Specifically, in practice you would typically proceed by first computing the mean $\hat{\boldsymbol{\mu}}$ of your data as described above. Given $\hat{\boldsymbol{\mu}}$, you can then form the matrix $\widetilde{\mathbf{X}}$ whose columns are given by

$$\widetilde{\mathbf{x}}_n = \mathbf{x}_n - \hat{\boldsymbol{\mu}}.$$

Alternatively, if you know a priori that your columns of zero mean (or should have zero mean) based on the underlying process generating the data, then you can skip this step, setting $\widetilde{\mathbf{X}} = \mathbf{X}$.

In either case, once you have formed $\widetilde{\mathbf{X}}$, you simply compute the SVD of $\widetilde{\mathbf{X}} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}$ and then set

$$\begin{aligned}\widehat{\mathbf{Q}} &= \mathbf{U}_K, \\ \widehat{\boldsymbol{\theta}}_n &= \mathbf{U}_K^T \widetilde{\mathbf{x}}_n,\end{aligned}$$

where $\mathbf{U}_K = [\mathbf{u}_1 \ \mathbf{u}_2 \ \cdots \ \mathbf{u}_K]$ contains the first K columns of \mathbf{U} . We can think of $\widehat{\boldsymbol{\theta}}_n$ as a representation of \mathbf{x}_n is a K -dimensional

subspace, with $\hat{\mathbf{Q}}$ giving us a basis for that subspace (which is useful for projecting vectors $\mathbf{x} \in \mathbb{R}^N$ into the subspace).

Note that if you look up a discussion of PCA in most textbooks or online, you will typically see a slightly different presentation. Specifically, most texts describe an approach to the problem that involves forming the matrix

$$\mathbf{S} = \sum_{n=1}^N (\mathbf{x}_n - \hat{\boldsymbol{\mu}})(\mathbf{x}_n - \hat{\boldsymbol{\mu}})^T,$$

taking an eigenvalue decomposition $\mathbf{S} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T$, and then taking

$$\mathbf{Q} = [\mathbf{v}_1 \ \mathbf{v}_2 \ \cdots \ \mathbf{v}_K],$$

where $\mathbf{v}_1, \dots, \mathbf{v}_K$ are the eigenvectors of \mathbf{S} corresponding to the K largest eigenvalues.

This approach is **completely equivalent** to our approach above.² The reason that PCA is typically presented in this way is that \mathbf{S} can be interpreted as a scaled version of an empirical estimate of the covariance matrix for the underlying distribution generating the data. While this provides a nice connection with the other (statistical) interpretation of PCA, I personally find the SVD approach more intuitive. In PCA, we are simply trying to find a low-rank approximation to our dataset, which is directly and optimally handled by computing a truncated SVD.

²Recall the relationship between the SVD of $\widetilde{\mathbf{X}}$ and the eigendecomposition of $\widetilde{\mathbf{X}}\widetilde{\mathbf{X}}^T$.