

Total Least-Squares

Our main approach thus far to “solving” $\mathbf{y} \approx \mathbf{A}\mathbf{x}$ is to optimize

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2.$$

Thought of another way, if we can't find a \mathbf{x} such that $\mathbf{y} = \mathbf{A}\mathbf{x}$ exactly, we are looking for the smallest possible perturbation we could add to \mathbf{y} so that there is an exact solution. Mathematically, the standard least-squares program above is equivalent to solving

$$\underset{\Delta\mathbf{y}, \mathbf{x}}{\text{minimize}} \quad \|\Delta\mathbf{y}\|_2^2 \quad \text{subject to} \quad (\mathbf{y} + \Delta\mathbf{y}) = \mathbf{A}\mathbf{x}.$$

This reformulation makes it clear that least-squares implicitly assumes that all of the error (i.e. all of the reasons we can't find an exact solution) lies in the measured data \mathbf{y} .

But what if the entries of \mathbf{A} are also subject to error? That is, how can we account for *modeling error* as well as measurement error? *Total least-squares* (TLS) is a framework for doing exactly this in a principled manner. TLS finds the smallest perturbations $\Delta\mathbf{y}, \Delta\mathbf{A}$ such that

$$(\mathbf{y} + \Delta\mathbf{y}) = (\mathbf{A} + \Delta\mathbf{A})\mathbf{x}$$

has an exact solution. It does this by solving

$$\underset{\Delta\mathbf{A}, \Delta\mathbf{y}, \mathbf{x}}{\text{minimize}} \quad \|\Delta\mathbf{A}\|_F^2 + \|\Delta\mathbf{y}\|_2^2 \quad \text{subject to} \quad (\mathbf{y} + \Delta\mathbf{y}) = (\mathbf{A} + \Delta\mathbf{A})\mathbf{x},$$

where $\|\Delta\mathbf{A}\|_F^2$ is the **Frobenius norm** of $\Delta\mathbf{A}$ and is given by

$$\|\Delta\mathbf{A}\|_F^2 = \sum_{m=1}^M \sum_{n=1}^N |\Delta A[m, n]|^2.$$

Example: 1D linear regression

Say we are given a set of points

$$(a_1, y_1), (a_2, y_2), \dots, (a_M, y_M)$$

Suppose that the goal is to find the “best” line that fits these points. (For simplicity, we will only consider lines that pass through the origin.) That is, we are looking for the slope x such that the $a_m x$ are as close to the y_m as possible.

The standard least-squares framework models this problem as follows. We observe

$$y_m = a_m x + \text{noise},$$

or in matrix form,

$$\mathbf{y} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_M \end{bmatrix} x + \text{noise}.$$

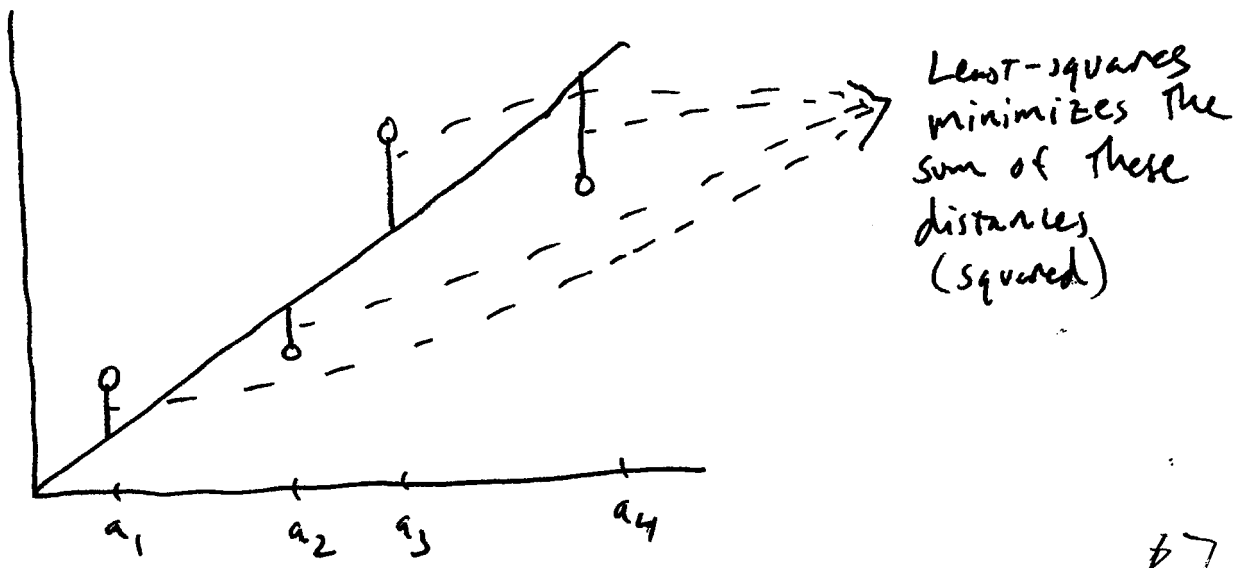
The solution is of course

$$\hat{x} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y} = \frac{\sum_{m=1}^M a_m y_m}{\sum_{m=1}^M a_m^2}.$$

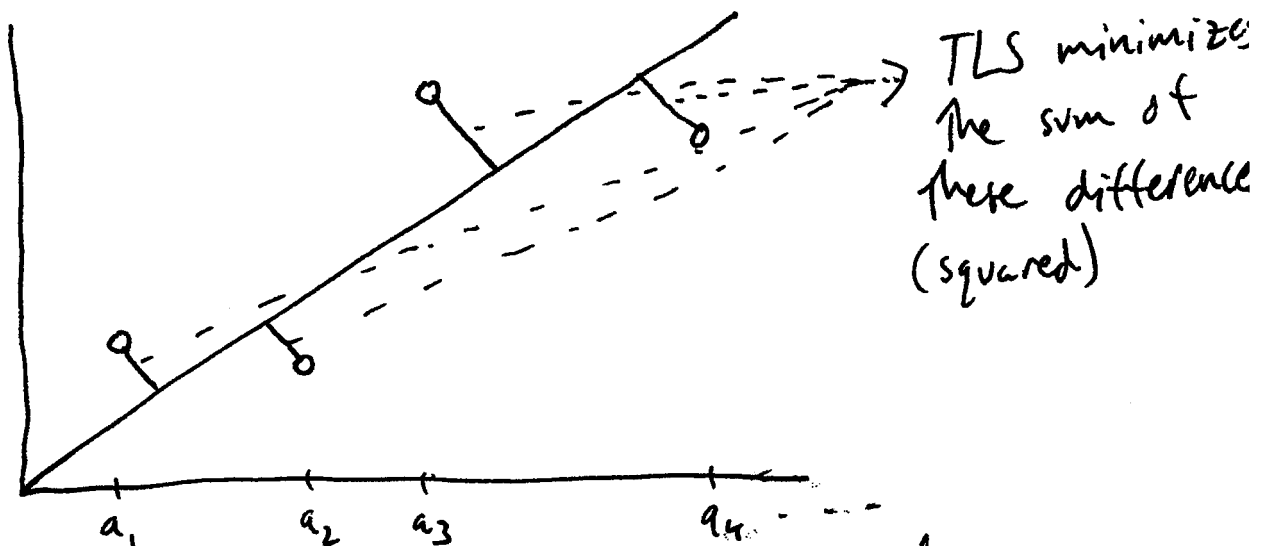
This solution minimizes the size of the residual

$$\|\mathbf{r}\|_2^2 = \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 = \sum_{m=1}^M |y_m - a_m x|^2.$$

Geometrically, we are choosing the slope that minimizes the sum of the squares of the *vertical distances* of the points to the line we choose to approximate them:



In contrast, the TLS estimate (which we will see how to compute below) minimizes the *distance in the plane* of the points to the line we choose:



This distance includes changes in both the a_m and y_m .

Solving TLS – Part I

We will assume that \mathbf{A} is an $M \times N$ matrix, with $M > N$, and $\text{rank}(\mathbf{A}) = N$ (i.e. \mathbf{A} is overdetermined with full column rank). The problem only really makes sense if $\text{rank}(\mathbf{A}) < M$, otherwise there is always an exact solution. By being careful with the details, the method we present here can also be extended to the case where $\text{rank}(\mathbf{A}) < N < M$, but we will leave it to you to fill in those gaps.

We want to find $\Delta\mathbf{A}, \Delta\mathbf{y}, \mathbf{x}$ such that

$$(\mathbf{y} + \Delta\mathbf{y}) = (\mathbf{A} + \Delta\mathbf{A})\mathbf{x},$$

for $\Delta\mathbf{y}, \Delta\mathbf{A}$ of minimal size. Rewrite this as

$$\begin{aligned}(\mathbf{A} + \Delta\mathbf{A})\mathbf{x} - (\mathbf{y} + \Delta\mathbf{y}) &= \mathbf{0} \\ \Rightarrow [\mathbf{A} + \Delta\mathbf{A} \quad \mathbf{y} + \Delta\mathbf{y}] \begin{bmatrix} \mathbf{x} \\ -1 \end{bmatrix} &= \mathbf{0} \\ \Rightarrow (\mathbf{C} + \Delta) \begin{bmatrix} \mathbf{x} \\ -1 \end{bmatrix} &= \mathbf{0}\end{aligned}$$

where

$$\mathbf{C} = [\mathbf{A} \quad \mathbf{y}], \quad \Delta = [\Delta\mathbf{A} \quad \Delta\mathbf{y}].$$

Note that both \mathbf{C} and Δ are $M \times (N + 1)$ matrices.

The result of the progression of equations above says that we are looking for a Δ (of minimal size) such that there is a vector $\begin{bmatrix} \mathbf{x} \\ -1 \end{bmatrix}$ in the nullspace of $\mathbf{C} + \Delta$. Since

$$\mathbf{v} \in \text{Null}(\mathbf{C} + \Delta) \Leftrightarrow \alpha\mathbf{v} \in \text{Null}(\mathbf{C} + \Delta) \quad \text{for all } \alpha \in \mathbb{R},$$

and \mathbf{x} in arbitrary, we are really just asking that $\mathbf{C} + \Delta$ has a nullspace; as long as there is at least one vector in the nullspace

whose last entry is nonzero, we can find a vector of the required form just by normalizing. In short, this means that our task is to find $\mathbf{\Delta}$ such that the $M \times (N + 1)$ matrix $\mathbf{C} + \mathbf{\Delta}$ is **rank deficient**, that is $\text{rank}(\mathbf{C} + \mathbf{\Delta}) < N + 1$.

Put another way, we want to solve the optimization program

$$\underset{\mathbf{\Delta}}{\text{minimize}} \quad \|\mathbf{\Delta}\|_F^2 \quad \text{subject to} \quad \text{rank}(\mathbf{C} + \mathbf{\Delta}) = N.$$

Making the substitution $\mathbf{X} = \mathbf{C} + \mathbf{\Delta}$, this is equivalent to solving

$$\underset{\mathbf{X}}{\text{minimize}} \quad \|\mathbf{C} - \mathbf{X}\|_F^2 \quad \text{subject to} \quad \text{rank}(\mathbf{X}) = N,$$

and then taking $\widehat{\mathbf{\Delta}} = \widehat{\mathbf{X}} - \mathbf{C}$.

This is a low-rank approximation problem.¹ We now consider how to solve problems of this form (which arise in many other important contexts).

The SVD and Matrix Approximation

Let \mathbf{A} be an $M \times N$ matrix with rank R . We are often interested in determining the closest matrix to \mathbf{A} that has rank r ²? More precisely, we would like to solve

$$\underset{\mathbf{X}}{\text{minimize}} \quad \|\mathbf{A} - \mathbf{X}\|_F^2 \quad \text{subject to} \quad \text{rank}(\mathbf{X}) = r. \quad (1)$$

The functional above is standard least-squares, but the constraint set (the set of rank- r matrices) is a complicated entity. Nevertheless, as with many things in this class, the SVD reveals the solution immediately.

¹Or at least a “lower rank” approximation problem.

²We will assume that $r < R$, as for $r = R$ the answer is easy, and for $R < r \leq \min(M, N)$ the question is not well-posed.

Low-rank approximation.

Let \mathbf{A} be a matrix with SVD

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \sum_{p=1}^R \sigma_p \mathbf{u}_p \mathbf{v}_p^T.$$

Then (1) is solved simply by truncating the SVD:

$$\widehat{\mathbf{X}} = \sum_{p=1}^r \sigma_p \mathbf{u}_p \mathbf{v}_p^T = \mathbf{U}_r \mathbf{\Sigma}_r \mathbf{V}_r^T,$$

where \mathbf{U}_r contains the first r columns of \mathbf{U} , \mathbf{V}_r contains the first r columns of \mathbf{V} , and $\mathbf{\Sigma}_r$ is the first r columns and r rows of $\mathbf{\Sigma}$.

The framed result above, known as the Eckart-Young theorem, is an immediate consequence of the following lemma.

Subspace Approximation Lemma. For fixed \mathbf{A} with SVD $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$, the optimization program

$$\underset{\substack{\mathbf{Q}: M \times r \\ \mathbf{\Theta}: r \times N}}{\text{minimize}} \|\mathbf{A} - \mathbf{Q}\mathbf{\Theta}\|_F^2 \quad \text{subject to} \quad \mathbf{Q}^T \mathbf{Q} = \mathbf{I}, \quad (2)$$

has solution

$$\begin{aligned} \widehat{\mathbf{Q}} &= \mathbf{U}_r, \\ \widehat{\mathbf{\Theta}} &= \mathbf{U}_r^T \mathbf{A}, \end{aligned}$$

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

We prove this lemma in the technical details section at the end of the notes. To see how it implies the Eckart-Young theorem, we can

interpret the search over $M \times r$ matrices \mathbf{Q} with orthonormal columns as a search over all possible *column spaces* of dimension r . Then the search over $\mathbf{\Theta}$ finds the best linear combinations in a column spaces to approximate the columns of \mathbf{A} . Since any rank- r matrix can be represented this way, the optimization program (2) is equivalent to (1); if $\hat{\mathbf{Q}}, \hat{\mathbf{\Theta}}$ solve (2), then $\hat{\mathbf{A}} = \hat{\mathbf{Q}}\hat{\mathbf{\Theta}}$ solves (1).

Also note that

$$\hat{\mathbf{\Theta}} = \mathbf{U}_r^T \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T = [\mathbf{I} \ \mathbf{0}] \mathbf{\Sigma} \mathbf{V}^T,$$

where \mathbf{I} is the $r \times r$ identity matrix, and $\mathbf{0}$ is a $r \times (R - r)$ matrix of zeros. This matrix of all zeros has the same effect as removing all but the first r terms along the diagonal of $\mathbf{\Sigma}$ and all but the first r rows of \mathbf{V}^T . Thus

$$\hat{\mathbf{Q}}\hat{\mathbf{\Theta}} = \mathbf{U}_r [\mathbf{I} \ \mathbf{0}] \mathbf{\Sigma} \mathbf{V}^T = \mathbf{U}_r \mathbf{\Sigma}_r \mathbf{V}_r^T.$$

What is the error between \mathbf{A} and its best rank- r approximation $\hat{\mathbf{A}}$? Well,

$$\mathbf{A} - \hat{\mathbf{A}} = \sum_{p=r+1}^R \sigma_p \mathbf{u}_p \mathbf{v}_p^T,$$

and so the error matrix has singular values $\sigma_{r+1}, \dots, \sigma_R$. Since the Frobenius norm (squared) can be calculated by summing the squares of the singular values,

$$\|\mathbf{A} - \hat{\mathbf{A}}\|_F^2 = \sum_{p=r+1}^R \sigma_p^2.$$

Solving TLS – Part II

Recall that the TLS problem reduced to solving a problem of the form

$$\underset{\mathbf{X}}{\text{minimize}} \|\mathbf{C} - \mathbf{X}\|_F^2 \quad \text{subject to} \quad \text{rank}(\mathbf{X}) = N,$$

and then taking $\widehat{\Delta} = \widehat{\mathbf{X}} - \mathbf{C}$. In light of our discussion above, we take the SVD of \mathbf{C} ,

$$\mathbf{C} = \mathbf{W}\mathbf{\Gamma}\mathbf{Z}^T = \sum_{n=1}^{N+1} \gamma_n \mathbf{w}_n \mathbf{z}_n^T,$$

and create $\widehat{\mathbf{X}}$ by leaving out the last term in the sum above³:

$$\widehat{\mathbf{X}} = \sum_{n=1}^N \gamma_n \mathbf{w}_n \mathbf{z}_n^T.$$

Then

$$\widehat{\Delta} = \widehat{\mathbf{X}} - \mathbf{C} = -\gamma_{N+1} \mathbf{w}_{N+1} \mathbf{z}_{N+1}^T.$$

Now we are ready to construct the actual estimate $\widehat{\mathbf{x}}$. Recall that we want a vector such that

$$(\mathbf{C} + \widehat{\Delta}) \begin{bmatrix} \mathbf{x} \\ -1 \end{bmatrix} = \mathbf{0}, \quad \text{meaning} \quad \widehat{\mathbf{X}} \begin{bmatrix} \mathbf{x} \\ -1 \end{bmatrix} = \mathbf{0}.$$

The null space of $\widehat{\mathbf{X}}$ is (by construction) simply the span of \mathbf{z}_{N+1} , meaning we need to find a scalar α such that

$$\begin{bmatrix} \mathbf{x} \\ -1 \end{bmatrix} = \alpha \mathbf{z}_{N+1}.$$

³If \mathbf{C} has fewer than $N + 1$ non-zero singular values, then it is already rank deficient, and we can take $\widehat{\mathbf{X}} = \mathbf{C} \Rightarrow \widehat{\Delta} = \mathbf{0}$.

Thus we can take

$$\hat{\mathbf{x}}_{\text{TLS}} = \frac{-1}{z_{N+1}[N+1]} \cdot \begin{bmatrix} z_{N+1}[1] \\ z_{N+1}[2] \\ \vdots \\ z_{N+1}[N] \end{bmatrix}.$$

If it happens that $z_{N+1}(N+1) = 0$, this means $\widehat{\Delta\mathbf{y}} = \mathbf{0}$, and we would need an \mathbf{x} such that

$$(\mathbf{A} + \widehat{\Delta\mathbf{A}})\mathbf{x} = \mathbf{y}.$$

Such an \mathbf{x} may or may not exist (and probably doesn't), so in this case there is no TLS solution.

In the special case where the smallest singular value of $\mathbf{C} = [\mathbf{A} \ \mathbf{y}]$ is not unique, i.e.

$$\gamma_1 \geq \gamma_2 \geq \gamma_q > \gamma_{q+1} = \gamma_{q+2} = \cdots = \gamma_{N+1}, \quad \text{for some } q < N,$$

then the TLS solution may not be unique. We take

$$\mathbf{Z}' = [\mathbf{z}_{q+1} \ \mathbf{z}_{q+2} \ \cdots \ \mathbf{z}_{N+1}],$$

and try to find a vector in the span that has the right form; any vector \mathbf{x} such that

$$\begin{bmatrix} \mathbf{x} \\ -1 \end{bmatrix} \in \text{Span}(\{\mathbf{z}_{q+1}, \dots, \mathbf{z}_{N+1}\})$$

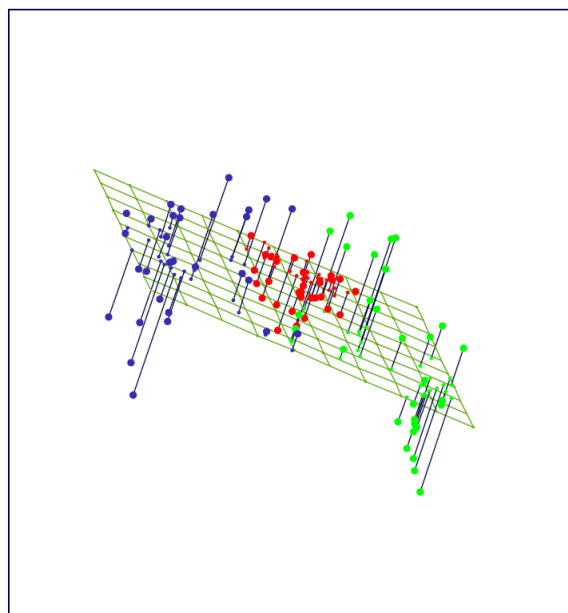
is equally good. All we need is a $\boldsymbol{\beta}$ such that the last entry of $\mathbf{Z}'\boldsymbol{\beta}$ is equal to -1 .

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 picture⁴



⁴From 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 as a K -dimensional

subspace, with $\widehat{\mathbf{Q}}$ 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

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

taking an eigenvalue decomposition $\mathbf{S} = \mathbf{V}\boldsymbol{\Lambda}\mathbf{V}^T$, 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.⁵ 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.

⁵Recall the relationship between the SVD of $\widetilde{\mathbf{X}}$ and the eigendecomposition of $\widetilde{\mathbf{X}}\widetilde{\mathbf{X}}^T$.

Technical Details: Subspace Approx. Lemma

We prove the subspace approximation lemma from above. First, with \mathbf{Q} fixed, we can break the optimization over Θ into a series of least-squares problems. Let $\mathbf{a}_1, \dots, \mathbf{a}_N$ be the columns of \mathbf{A} , and $\theta_1, \dots, \theta_N$ be the columns of Θ . Then

$$\underset{\Theta}{\text{minimize}} \quad \|\mathbf{A} - \mathbf{Q}\Theta\|_F^2$$

is exactly the same as

$$\underset{\theta_1, \dots, \theta_N}{\text{minimize}} \quad \sum_{n=1}^N \|\mathbf{a}_n - \mathbf{Q}\theta_n\|_2^2.$$

The above is our classic closest point problem, and is optimized by taking $\theta_n = \mathbf{Q}^T \mathbf{a}_n$ (since the columns of \mathbf{Q} are orthonormal). Thus we can write the original problem (2) as

$$\underset{\mathbf{Q}:M \times r}{\text{minimize}} \quad \sum_{n=1}^N \|\mathbf{a}_n - \mathbf{Q}\mathbf{Q}^T \mathbf{a}_n\|_2^2 \quad \text{subject to} \quad \mathbf{Q}^T \mathbf{Q} = \mathbf{I},$$

and then take $\hat{\Theta} = \hat{\mathbf{Q}}^T \mathbf{A}$.

Expanding the functional and using the fact that $(\mathbf{I} - \mathbf{Q}\mathbf{Q}^T)^2 = (\mathbf{I} - \mathbf{Q}\mathbf{Q}^T)$, we have

$$\begin{aligned} \sum_{n=1}^N \|\mathbf{a}_n - \mathbf{Q}\mathbf{Q}^T \mathbf{a}_n\|_2^2 &= \sum_{n=1}^N \mathbf{a}_n^T (\mathbf{I} - \mathbf{Q}\mathbf{Q}^T) \mathbf{a}_n \\ &= \sum_{n=1}^N \|\mathbf{a}_n\|_2^2 - \mathbf{a}_n^T \mathbf{Q}\mathbf{Q}^T \mathbf{a}_n. \end{aligned}$$

Since the first term does not depend on \mathbf{Q} , our optimization program is equivalent to

$$\underset{\mathbf{Q}:M \times r}{\text{maximize}} \sum_{n=1}^N \mathbf{a}_n^T \mathbf{Q} \mathbf{Q}^T \mathbf{a}_n \quad \text{subject to} \quad \mathbf{Q}^T \mathbf{Q} = \mathbf{I}.$$

Now recall that for any vector \mathbf{v} , $\langle \mathbf{v}, \mathbf{v} \rangle = \text{trace}(\mathbf{v} \mathbf{v}^T)$. Thus

$$\begin{aligned} \sum_{n=1}^N \mathbf{a}_n \mathbf{Q} \mathbf{Q}^T \mathbf{a}_n &= \sum_{n=1}^N \text{trace}(\mathbf{Q}^T \mathbf{a}_n \mathbf{a}_n^T \mathbf{Q}) \\ &= \text{trace} \left(\mathbf{Q}^T \left(\sum_{n=1}^N \mathbf{a}_n \mathbf{a}_n^T \right) \mathbf{Q} \right) \\ &= \text{trace} \left(\mathbf{Q}^T (\mathbf{A} \mathbf{A}^T) \mathbf{Q} \right). \end{aligned}$$

The matrix $\mathbf{A} \mathbf{A}^T$ has eigenvalue decomposition

$$\mathbf{A} \mathbf{A}^T = \mathbf{U} \mathbf{\Sigma}^2 \mathbf{U}^T,$$

where \mathbf{U} and $\mathbf{\Sigma}$ come from the SVD of \mathbf{A} (we will take \mathbf{U} to be $M \times M$, possibly adding zeros down the diagonal of $\mathbf{\Sigma}^2$). Now

$$\begin{aligned} \text{trace} \left(\mathbf{Q}^T (\mathbf{A} \mathbf{A}^T) \mathbf{Q} \right) &= \text{trace} \left(\mathbf{Q}^T \mathbf{U} \mathbf{\Sigma}^2 \mathbf{U}^T \mathbf{Q} \right) \\ &= \text{trace} \left(\mathbf{W}^T \mathbf{\Sigma}^2 \mathbf{W} \right), \end{aligned}$$

where $\mathbf{W} = \mathbf{U}^T \mathbf{Q}$. Notice that \mathbf{W} also has orthonormal columns, as $\mathbf{W}^T \mathbf{W} = \mathbf{Q}^T \mathbf{U} \mathbf{U}^T \mathbf{Q} = \mathbf{Q}^T \mathbf{Q} = \mathbf{I}$. Thus our optimization program has become

$$\underset{\mathbf{W}:M \times r}{\text{maximize}} \text{trace}(\mathbf{W}^T \mathbf{\Sigma}^2 \mathbf{W}) \quad \text{subject to} \quad \mathbf{W}^T \mathbf{W} = \mathbf{I}.$$

After we solve this, we can take any $\widehat{\mathbf{Q}}$ such that $\widehat{\mathbf{W}} = \mathbf{U}^T \widehat{\mathbf{Q}}$.

This last optimization program is equivalent to a simple linear program that is solvable by inspection. Let $\mathbf{w}_1, \dots, \mathbf{w}_r$ be the columns of \mathbf{W} . Then

$$\begin{aligned} \text{trace}(\mathbf{W}^T \boldsymbol{\Sigma}^2 \mathbf{W}) &= \sum_{p=1}^r \mathbf{w}_p^T \boldsymbol{\Sigma}^2 \mathbf{w}_p \\ &= \sum_{p=1}^r \sum_{m=1}^M |w_p[m]|^2 \sigma_m^2 \\ &= \sum_{m=1}^M h[m] \sigma_m^2, \quad \text{where } h[m] = \sum_{p=1}^r |w_p[m]|^2. \end{aligned}$$

Notice that

$$h[m] = \sum_{p=1}^r |W[p, m]|^2$$

is a sum of the squares of a *row* of \mathbf{W} . Since the sum of the squares of every column of \mathbf{W} is one, the sum of the squares of every entry in \mathbf{W} must be r , and so

$$\sum_{m=1}^M h[m] = r.$$

It is clear that $h[m]$ is non-negative, but it also true that $h[m] \leq 1$. Here is why: since the columns of \mathbf{W} are orthonormal, they can be considered as part of an orthonormal basis for \mathbb{R}^M . That is, there is a $M \times (M - r)$ matrix \mathbf{W}_0 such that the $M \times M$ matrix $[\mathbf{W} \ \mathbf{W}_0]$ has both orthonormal columns and orthonormal rows — thus the sum of the squares of each row are equal to one. Thus the sum of the squares of the first r entries cannot be larger than this.

Thus the maximum value $\text{trace}(\mathbf{W}^T \boldsymbol{\Sigma}^2 \mathbf{W})$ can take is given by the linear program

$$\underset{\mathbf{h} \in \mathbb{R}^M}{\text{maximize}} \quad \sum_{m=1}^M h[m] \sigma_m^2 \quad \text{subject to} \quad \sum_{m=1}^M h[m] = r, \quad 0 \leq h[m] \leq 1.$$

We can intuit the answer to this program. Since all of the σ_m^2 and all of the $h[m]$ are positive, we want to have as much weight as possible assigned to the largest singular values. Since the weights are constrained to be less than 1, this simply means we “max out” the first r terms; the solution to the program above is

$$\hat{h}[m] = \begin{cases} 1, & m = 1, \dots, r \\ 0, & m = r + 1, \dots, M. \end{cases}$$

This means that the sum of the squares of the first r rows in $\widehat{\mathbf{W}}$ are equal to one, while the rest are zero. There might be many such matrices that fit this bill, but one of them is

$$\widehat{\mathbf{W}} = \begin{bmatrix} \mathbf{I} \\ \mathbf{0} \end{bmatrix},$$

where above, \mathbf{I} is the $r \times r$ identity matrix, and $\mathbf{0}$ is a $(M - r) \times r$ matrix of all zeros. It is easy to see that choosing

$$\widehat{\mathbf{Q}} = [\mathbf{u}_1 \quad \mathbf{u}_2 \quad \cdots \quad \mathbf{u}_r]$$

satisfies

$$\mathbf{U}^T \widehat{\mathbf{Q}} = \begin{bmatrix} \mathbf{I} \\ \mathbf{0} \end{bmatrix}.$$