

# IV. Matrix Approximation using Least-Squares

## The SVD and Matrix Approximation

We begin with the following fundamental question. Let  $\mathbf{A}$  be an  $M \times N$  matrix with rank  $R$ . What is the closest matrix to  $\mathbf{A}$  that has rank  $r$ <sup>1</sup>?

As before, we will use the Frobenius norm to measure the distance between two matrices:

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

Recall that  $\|\mathbf{X}\|_F^2$  is also equal to the sum of the squares of the singular values of  $\mathbf{X}$ .

We can now formulate our problem as

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

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

### Low-rank approximation.

Let  $\mathbf{A}$  be a matrix with SVD

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

---

<sup>1</sup>We 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.

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

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

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, which we will actually use again later in this set of notes.

**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

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

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,

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

In what follows, we use this low-rank matrix approximation result to develop two fundamental tools: total least-squares, and principal components analysis.

## Total Least-Squares

Our fundamental 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

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

### 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  $a_mx$  are as close to the  $y_m$  as possible.

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

$$y_m = a_mx + \text{noise},$$

or in matrix form,

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

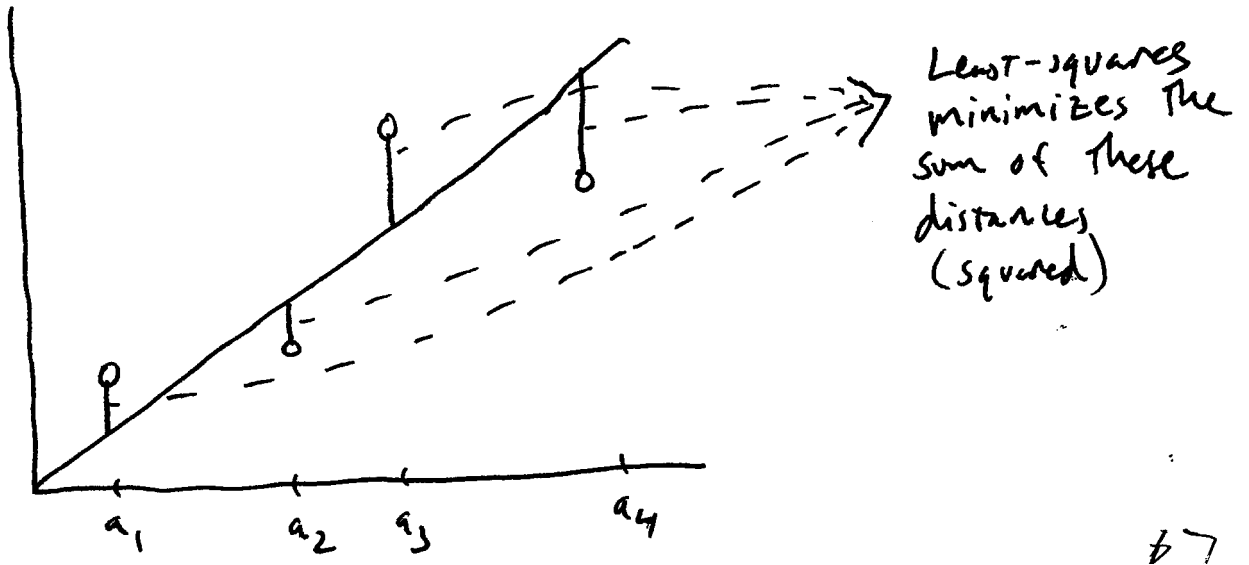
The solution is of course

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

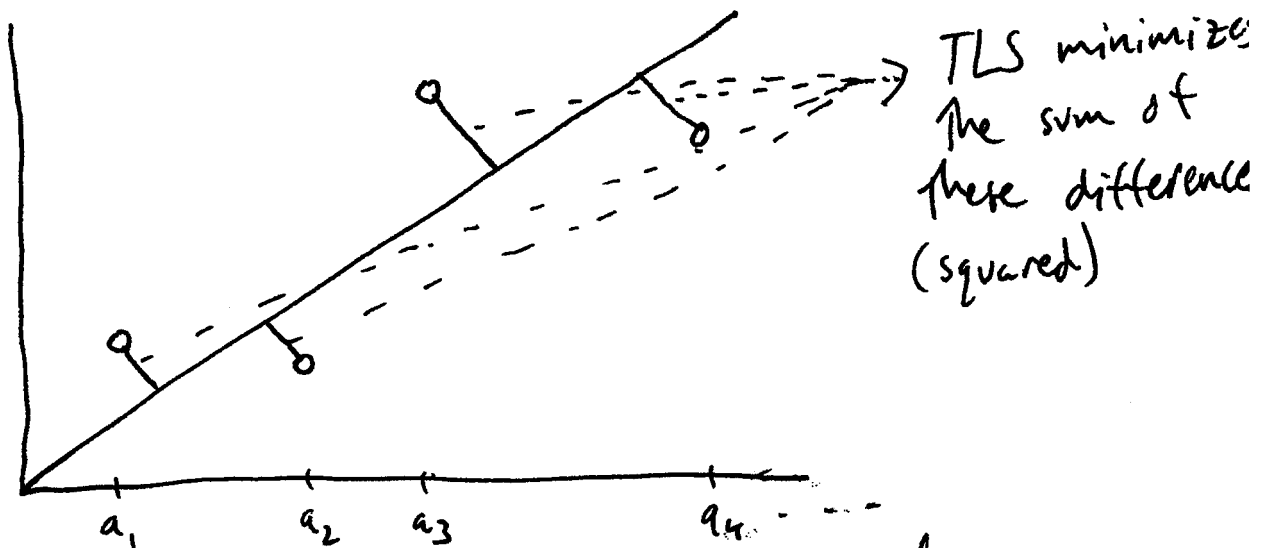
This solution minimizes the size of the residual

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

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

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 I 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  $\Delta$  are  $M \times (N + 1)$  matrices.

The result of the progression of equations above says that we are looking for a  $\Delta$  (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  $\hat{\mathbf{\Delta}} = \hat{\mathbf{X}} - \mathbf{C}$ .

This is a low-rank approximation problem<sup>2</sup>, and we now know exactly how to solve it. 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  $\hat{\mathbf{X}}$  by leaving out the last term in the sum above<sup>3</sup>:

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

Then

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

---

<sup>2</sup>Or at least a “lower rank” approximation problem.

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

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

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

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

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

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  $\beta$  such that the last entry of  $\mathbf{Z}'\beta$  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 simplifying **linear relationships** in the data. It is used everywhere in signal processing, machine learning, and statistics, with applications including data compression, pattern recognition, and factor analysis.

There are two 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 is geometrical: given a set of vectors, we are trying to find a subspace of a certain dimension that comes closest to containing this set.

## The Karhunen-Loeve Transform

The Karhunen-Loeve (KL) transform is an orthobasis that is tailored to the statistics of a class of random vectors. Suppose that  $\mathbf{x} \in \mathbb{R}^D$

is random and has<sup>4</sup> mean and covariance

$$\mathbb{E}[\mathbf{x}] = \mathbf{0}, \quad \mathbb{E}[\mathbf{x}\mathbf{x}^T] = \mathbf{R}.$$

Then the KL transform (or the KL basis) is simply the eigenvector  $\mathbf{V}$  of  $\mathbf{R} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T$ :

$$\mathbf{x} = \sum_{n=1}^D \alpha_n \mathbf{v}_n, \quad \alpha_n = \langle \mathbf{x}, \mathbf{v}_n \rangle.$$

This transform has the property that if we want to truncate the sum above (i.e. compress the vector by using fewer than  $D$  numbers to represent it), we get an error that is optimal in the mean-square error sense.

Let's set this problem up carefully. We want to find a subspace  $\mathcal{T}$  of dimension  $K$  such that when we project  $\mathbf{x}$  onto  $\mathcal{T}$ , we lose as little of  $\mathbf{x}$  (in expectation) as possible. We want to solve

$$\underset{\mathcal{T}}{\text{minimize}} \mathbb{E} \left[ \underset{\mathbf{t} \in \mathcal{T}}{\text{minimize}} \|\mathbf{x} - \mathbf{t}\|_2^2 \right] \quad \text{subject to} \quad \dim(\mathcal{T}) = K.$$

For a fixed  $\mathcal{T}$ , we know how to solve the inner optimization program if we have an orthobasis, so we can re-write the above as a search of sets of  $K$  orthogonal vectors in  $\mathbb{R}^D$ :

$$\underset{\mathbf{Q}:D \times K}{\text{minimize}} \mathbb{E} \left[ \|\mathbf{x} - \mathbf{Q}\mathbf{Q}^T \mathbf{x}\|_2^2 \right] \quad \text{subject to} \quad \mathbf{Q}^T \mathbf{Q} = \mathbf{I}.$$

---

<sup>4</sup>Modifying this discussion to vectors that are not zero-mean is straightforward.

Now notice that

$$\begin{aligned}
 \mathbb{E} \left[ \|\mathbf{x} - \mathbf{Q}\mathbf{Q}^T\mathbf{x}\|_2^2 \right] &= \mathbb{E} \left[ \|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^T)\mathbf{x}\|_2^2 \right] \\
 &= \mathbb{E}[\text{trace}((\mathbf{I} - \mathbf{Q}\mathbf{Q}^T)\mathbf{x}\mathbf{x}^T(\mathbf{I} - \mathbf{Q}\mathbf{Q}^T))] \\
 &= \text{trace}((\mathbf{I} - \mathbf{Q}\mathbf{Q}^T) \mathbb{E}[\mathbf{x}\mathbf{x}^T](\mathbf{I} - \mathbf{Q}\mathbf{Q}^T)) \\
 &= \text{trace}((\mathbf{I} - \mathbf{Q}\mathbf{Q}^T)\mathbf{R}(\mathbf{I} - \mathbf{Q}\mathbf{Q}^T)),
 \end{aligned}$$

where in the second step above we have used the fact that for any vector  $\mathbf{v}$ ,  $\|\mathbf{v}\|_2^2 = \text{trace}(\mathbf{v}\mathbf{v}^T)$ . Now notice that

$$\begin{aligned}
 \text{trace}((\mathbf{I} - \mathbf{Q}\mathbf{Q}^T)\mathbf{R}(\mathbf{I} - \mathbf{Q}\mathbf{Q}^T)) &= \\
 \text{trace}(\mathbf{R}) - 2\text{trace}(\mathbf{Q}\mathbf{Q}^T\mathbf{R}) + \text{trace}(\mathbf{Q}\mathbf{Q}^T).
 \end{aligned}$$

We now apply three facts:  $\text{trace}(\mathbf{R})$  does not depend on  $\mathbf{Q}$ ,  $\text{trace}(\mathbf{Q}\mathbf{Q}^T) = \text{trace}(\mathbf{Q}^T\mathbf{Q}) = K$  also does not depend on  $\mathbf{Q}$ , and

$$\text{trace}(\mathbf{Q}\mathbf{Q}^T\mathbf{R}) = \text{trace}(\mathbf{Q}^T\mathbf{R}\mathbf{Q}),$$

to transform into the equivalent program

$$\underset{\mathbf{W}:D \times K}{\text{maximize}} \quad \text{trace}(\mathbf{W}^T\mathbf{R}\mathbf{W}) \quad \text{subject to} \quad \mathbf{W}^T\mathbf{W} = \mathbf{I}.$$

In the Technical Details section below, we show that this expression is maximized by taking

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

where the  $\mathbf{v}_k$  correspond to the  $K$  eigenvectors of  $\mathbf{R}$  corresponding to the  $K$  largest eigenvalues.

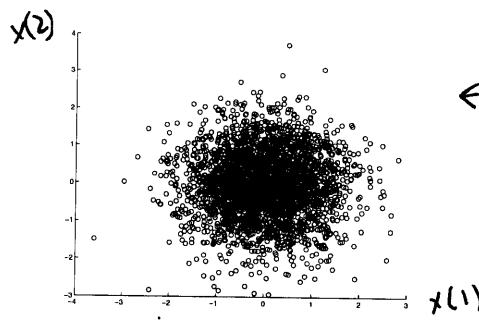
**Moral:** The best (in terms of mean-squared error) way to get a  $K$  term approximation of random data is to transform into the orthonobasis formed by the eigenvectors of the covariance matrix, then

truncating the coefficients to  $K$  terms. This set of eigenvectors  $\mathbf{V}$  is called the **KL transform**.

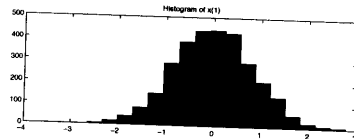
In some sense, the  $\mathbf{v}_1, \dots, \mathbf{v}_K$  are the  $K$  most important “features” of  $\mathbf{x}$  — they are completely determined by the covariance matrix  $\mathbf{R}$ .

Examples in  $\mathbb{R}^2$ :

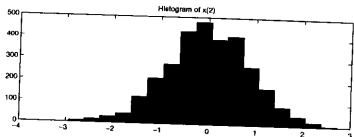
①  $\mathbf{R} = \mathbf{I}$ . A typical scatter plot will look like:



← 3000 realizations



← Histogram of  $x$  projected onto horizontal axis



← Histogram of  $x$  projected onto vertical axis

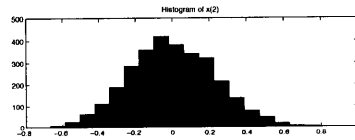
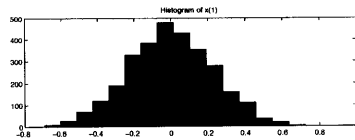
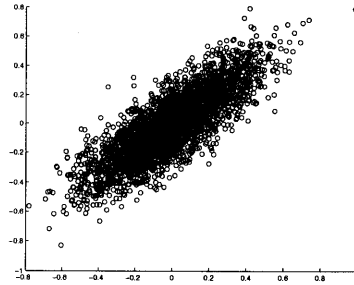
②  $R = V\Lambda V^T$

$\Lambda = \begin{bmatrix} \lambda_1 & \\ & \lambda_2 \end{bmatrix}$

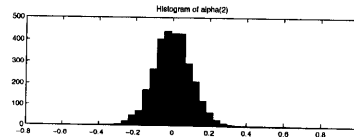
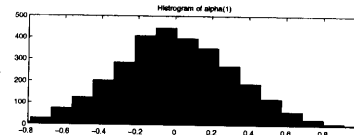
with  $\lambda_1 = 10\lambda_2$

$V = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$

← scatter plot, 3000 realizations



→ Histograms projected onto horiz. (top) and vert. (bottom) axes.



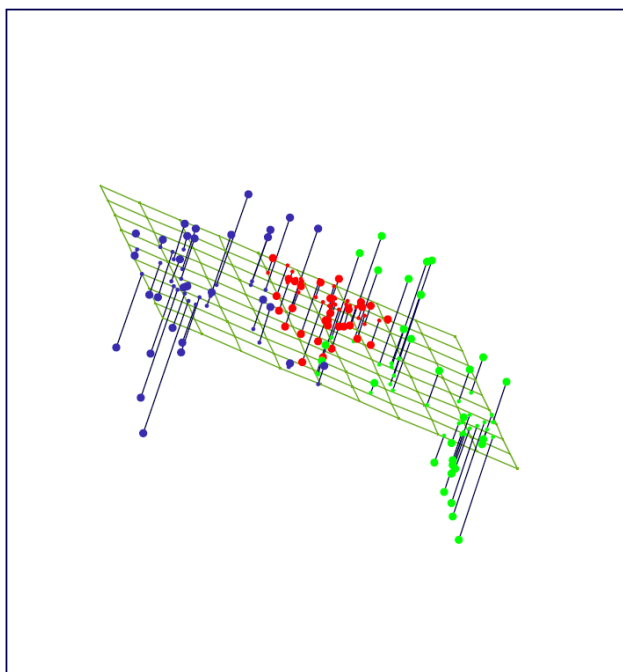
← projected onto the principal axis [1]

← projected onto  $v_2 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$

//

## PCA on observed data

A very similar procedure to the above solves a common geometrical problem. Suppose that I have a bunch of data points  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N \in \mathbb{R}^D$ , and I want to find the  $K$ -dimensional affine space (subspace plus offset) that comes closest to containing them. We don't even need to think of the data as random here; they are just points that we want to fit with a hyperplane.<sup>5</sup>



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

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

---

<sup>5</sup>This is pulled from Chapter 14 of Tibshirani and Hastie's *Elements of Statistical Learning*.



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

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

If we fix  $\boldsymbol{\mu}$  and  $\mathbf{Q}$ , then by arguments very similar to those we have made before, the optimal  $\boldsymbol{\theta}_n$  are given by

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

This means our objective reduces to solving

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

The offset  $\boldsymbol{\mu}$  is unconstrained; if we again fix  $\mathbf{Q}$ , we can solve for the optimal  $\boldsymbol{\mu}$  by taking a gradient and setting it equal to zero:

$$\begin{aligned} \nabla_{\boldsymbol{\mu}} \left( \sum_{n=1}^N \|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^T)(\mathbf{x}_n - \boldsymbol{\mu})\|_2^2 \right) &= -2 \sum_{n=1}^N (\mathbf{I} - \mathbf{Q}\mathbf{Q}^T)(\mathbf{x}_n - \boldsymbol{\mu}) \\ &= -2(\mathbf{I} - \mathbf{Q}\mathbf{Q}^T) \left( \left( \sum_{n=1}^N \mathbf{x}_n \right) - N\boldsymbol{\mu} \right). \end{aligned}$$

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

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

All that remains is solving for  $\mathbf{Q}$ . We have

$$\underset{\mathbf{Q}:D \times K}{\text{minimize}} \sum_{n=1}^N \|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^T)(\mathbf{x}_n - \hat{\boldsymbol{\mu}})\|_2^2 \quad \text{subject to} \quad \mathbf{Q}^T\mathbf{Q} = \mathbf{I}.$$

Again, using an argument that perfectly parallels that in the Technical Details section below, this program is solved by forming

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

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

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

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

So even though we posed this problem as being purely geometrical, the answer parallels the statistical KL transform — we simply replace the “true” covariance matrix  $\mathbf{R}$  with the sample covariance  $N^{-1}\mathbf{S}$ .

## Technical Details: Subspace Approx. Lemma

We prove the subspace approximation lemma from page 2. First, with  $\mathbf{Q}$  fixed, we can break the optimization over  $\Theta$  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  $\Theta$ . 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

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

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

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

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

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

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

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  $\sigma_m^2$  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  $\hat{\mathbf{W}}$  are equal to one, while the rest are zero. There might be many such matrices that fit this bill, but one of them is

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

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

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

satisfies

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