

Chapter 6

Low-rank approximation

Contents (class version)

6.0 Introduction	6.2
6.1 Low-rank approximation via Frobenius norm	6.3
Implementation	6.8
1D example	6.15
Generalization to other norms	6.17
Bases for $\mathbb{F}^{M \times N}$	6.19
Low-rank approximation summary	6.22
Rank and stability	6.23
6.2 Sensor localization application (Multidimensional scaling)	6.24
Practical implementation	6.31
6.3 Proximal operators	6.34
6.4 Alternative low-rank approximation formulations	6.38
6.5 Choosing the rank or regularization parameter	6.46
OptShrink	6.50
6.6 Related methods: autoencoders and PCA	6.55

Relation to autoencoder with linear hidden layer	6.55
Relation to principal component analysis (PCA)	6.57
6.7 Subspace learning	6.59
6.8 Summary	6.63

6.0 Introduction

L§8.1

In many applications, **dimensionality reduction** is important. This chapter focuses on **low-rank approximation** of a matrix. There are theoretical models for why big matrices should be approximately low rank [1]. One of many applications is **image compression**, explored in HW.

Another application is source localization via multidimensional scaling, explored on p. [6.24](#).

Source material for this chapter includes [2, §8.1].

6.1 Low-rank approximation via Frobenius norm

We are given a matrix $\mathbf{A} \in \mathbb{F}^{M \times N}$ (often large), having rank $r \leq \min(M, N)$.

To perform **dimensionality reduction** we want to approximate \mathbf{A} by another matrix $\hat{\mathbf{A}}_K$ having rank $K \leq r$. Often we pick $K \ll r$.

To find the “best” $\hat{\mathbf{A}}_K$ we must define how closely $\hat{\mathbf{A}}_K$ approximates \mathbf{A} . The simplest metric is the **Frobenius norm** of the difference. This criterion leads to the follow **low-rank approximation problem**:

$$\hat{\mathbf{A}}_K \triangleq \mathcal{L}_K^{M \times N} \triangleq \left\{ \mathbf{B} \in \mathbb{F}^{M \times N} : \text{rank}(\mathbf{B}) \leq K \right\}. \quad (6.1)$$

This is a **non-convex** optimization problem because the set $\mathcal{L}_K^{M \times N}$ of rank- K matrices is not convex.

Why rank is a non-convex function

To see why $\text{rank}(\cdot)$ is a non-convex function, and $\mathcal{L}_K^{M \times N}$ is a nonconvex set, consider:

$$\mathbf{A} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}, \quad \alpha \in (0, 1) \implies \text{rank}(\alpha \mathbf{A} + (1 - \alpha) \mathbf{B}) = \text{rank}\left(\begin{bmatrix} \alpha & 0 \\ 0 & 1 - \alpha \end{bmatrix}\right) = 2 \notin \alpha \text{rank}(\mathbf{A}) + (1 - \alpha) \text{rank}(\mathbf{B}) = \alpha \cdot 1 + (1 - \alpha) \cdot 1 = 1.$$

Very remarkably, despite this non-convexity, there is a simple solution for $\hat{\mathbf{A}}_K$ based on any **SVD** of \mathbf{A} [3].

Theorem. The best rank (at most) K approximation uses the first (largest) K singular components of \mathbf{A} :

$$\hat{\mathbf{A}}_K = \boxed{\mathbf{U}_K \Sigma_K \mathbf{V}'_K} \quad \text{where } \mathbf{A} = \mathbf{U} \Sigma \mathbf{V}' = \sum_{k=1}^{\min(M,N)} \sigma_k \mathbf{u}_k \mathbf{v}'_k. \quad (6.2)$$

The approximation error depends on the singular values of the remaining (discarded) terms:

$$\left\| \hat{\mathbf{A}}_K - \mathbf{A} \right\|_{\text{F}} = \boxed{\left\| \mathbf{B} - \mathbf{A} \right\|_{\text{F}}} \leq \left\| \mathbf{B} - \mathbf{A} \right\|_{\text{F}}, \quad \forall \mathbf{B} \in \mathcal{L}_K^{M \times N}.$$

- Clearly the approximation error will be small if the remaining singular values are small compared to the first K singular values.
- Once again we see that the SVD is key tool. Here, the SVD is used both to construct the approximation and to quantify the approximation error.
- The original matrix \mathbf{A} has MN values. The approximation $\hat{\mathbf{A}}_K$ is formed from $K(M + N + 1)$ values. This saves memory if $K \ll \min(M, N)$.

What is the rank of $\sum_{k=1}^K \sigma_k \mathbf{u}_k \mathbf{v}'_k$ when $1 \leq K \leq r = \text{rank}(\mathbf{A})$? (Choose best answer.)

- A: Always 1 B: Always r C: Usually r D: Always K E: Usually K

??

Diagonal case: proof sketch

First consider a $M \times N$ (rectangular) diagonal matrix Σ having rank r , with descending diagonal values, where we want to approximate it by a matrix C of rank at most $K \leq r$, i.e., we want to solve:

$$\begin{aligned} \hat{C} &\triangleq \arg \min_{C \in \mathcal{L}_K^{M \times N}} \|C - \Sigma\|_F^2 = \arg \min_{C \in \mathcal{L}_K^{M \times N}} \left\| \begin{bmatrix} c_{11} & \dots & c_{1N} \\ & \vdots & \\ c_{M1} & \dots & c_{MN} \end{bmatrix} - \begin{bmatrix} \sigma_1 & & & & \\ & \ddots & & & \\ & & \sigma_K & & \\ & & & \sigma_{K+1} & \\ & & & & \ddots \\ & & & & & \sigma_r \\ & & & & & 0 \end{bmatrix} \right\|_F^2 \\ &= \arg \min_{C \in \mathcal{L}_K^{M \times N}} \sum_{k=1}^K (c_{kk} - \sigma_k)^2 + \sum_{k=K+1}^r (c_{kk} - \sigma_k)^2 + \sum_{k=r+1}^{\min(M,N)} (c_{kk} - 0)^2 + \sum_{m \neq n} (c_{mn} - 0)^2. \end{aligned}$$

To minimize this quadruple sum, subject to the constraint $\text{rank}(C) \leq K$, intuitively we prefer the last two terms to be zero. And to minimize the first two terms, because $\sigma_1 \geq \sigma_2 \geq \dots$, it seems natural to keep the first K (using $c_{kk} = \sigma_k$) and set the rest to zero, yielding the following solution:

$$\hat{c}_{mn} = \begin{cases} \sigma_m, & m = n \leq K \\ 0, & \text{otherwise} \end{cases} \implies \hat{C} = \begin{bmatrix} \sigma_1 & & & & \\ & \ddots & & & \\ & & \sigma_K & & \\ & & & \mathbf{0}_{r-K \times r-K} & \\ & & & & \mathbf{0}_{M-r \times N-r} \end{bmatrix} = \sum_{k=1}^K \sigma_k \mathbf{e}_k \tilde{\mathbf{e}}'_k, \quad (6.3)$$

where \mathbf{e}_k is the k th unit vector in \mathbb{F}^M and $\tilde{\mathbf{e}}_k$ is the k th unit vector in \mathbb{F}^N .

The conclusion (6.3) is correct, but the above reasoning is not rigorous. For a rigorous proof, see [wiki] [4].

Proof for general case

Now we assume that (6.3) is correct (it is, though not proven here), and use it to prove the general case.

Denote an SVD of \mathbf{A} by $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}'$. Rewrite any $\mathbf{B} \in \mathbb{F}^{M \times N}$ in terms of the \mathbf{U} and \mathbf{V} bases as follows:

$$\mathbf{B} = \underbrace{(\mathbf{U}\mathbf{U}')}_{\mathbf{I}} \mathbf{B} \underbrace{(\mathbf{V}\mathbf{V}')}_{\mathbf{I}} = \mathbf{U} \underbrace{(\mathbf{U}'\mathbf{B}\mathbf{V})}_{\hookrightarrow \triangleq \mathbf{C} \text{ (not diagonal in general)}} \mathbf{V}' = \mathbf{U}\mathbf{C}\mathbf{V}'. \quad (6.4)$$

Because \mathbf{U} and \mathbf{V} are unitary, $\text{rank}(\mathbf{B}) = \text{rank}(\mathbf{C})$, so (6.1) is equivalent to

$$\begin{aligned} \hat{\mathbf{A}}_K &= \mathbf{U}\hat{\mathbf{C}}\mathbf{V}', \quad \hat{\mathbf{C}} \triangleq \arg \min_{\mathbf{C} \in \mathcal{L}_K^{M \times N}} \left\| \underbrace{\mathbf{U}\mathbf{C}\mathbf{V}'}_{\mathbf{B}} - \underbrace{\mathbf{U}\Sigma\mathbf{V}'}_{\mathbf{A}} \right\|_{\text{F}} = \arg \min_{\mathbf{C} \in \mathcal{L}_K^{M \times N}} \|\mathbf{U}(\mathbf{C} - \Sigma)\mathbf{V}'\|_{\text{F}} \\ &= \arg \min_{\mathbf{C} \in \mathcal{L}_K^{M \times N}} \|\mathbf{C} - \Sigma\|_{\text{F}} = \sum_{k=1}^K \sigma_k \mathbf{e}_k \tilde{\mathbf{e}}'_k, \end{aligned}$$

using the result (6.3) for the diagonal case and the fact that the **Frobenius norm** is **unitarily invariant**.

Thus the final best approximation to \mathbf{A} having rank (at most) K is:

$$\hat{\mathbf{A}}_K = \mathbf{U}\hat{\mathbf{C}}\mathbf{V}' = \sum_{k=1}^K \sigma_k \mathbf{u}_k \mathbf{v}'_k, \text{ where } \mathbf{A} = \mathbf{U}\Sigma\mathbf{V}' = \sum_{k=1}^r \sigma_k \mathbf{u}_k \mathbf{v}'_k,$$

because $\mathbf{U}\mathbf{e}_k = \mathbf{u}_k$. The approximation error depends on the remaining singular values:

$$\left\| \hat{\mathbf{A}}_K - \mathbf{A} \right\|_{\text{F}} =$$

Alternative way to derive the approximation error expression:

(Read)

$$\begin{aligned} \mathbf{A} - \hat{\mathbf{A}}_K &= \sum_{k=1}^r \sigma_k \mathbf{u}_k \mathbf{v}'_k - \sum_{k=1}^K \sigma_k \mathbf{u}_k \mathbf{v}'_k = \sum_{k=K+1}^r \sigma_k \mathbf{u}_k \mathbf{v}'_k \Rightarrow \\ \left\| \hat{\mathbf{A}}_K - \mathbf{A} \right\|_{\text{F}} &= \left\| \sum_{k=K+1}^r \sigma_k \mathbf{u}_k \mathbf{v}'_k \right\|_{\text{F}} = \left\| \mathbf{U}' \sum_{k=K+1}^r \sigma_k \mathbf{u}_k \mathbf{v}'_k \mathbf{V} \right\|_{\text{F}} = \left\| \sum_{k=K+1}^r \sigma_k \mathbf{e}_k \mathbf{e}'_k \right\|_{\text{F}} = \sqrt{\sum_{k=K+1}^r \sigma_k^2}, \end{aligned}$$

which again is related to **Parseval's theorem**.

□

Summary

In words: the best rank (at most) K approximation to a matrix is given by its **truncated SVD** (to K terms).

Implementation

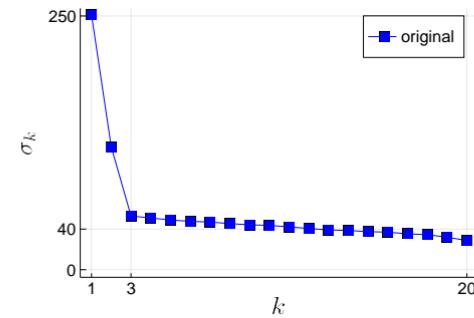
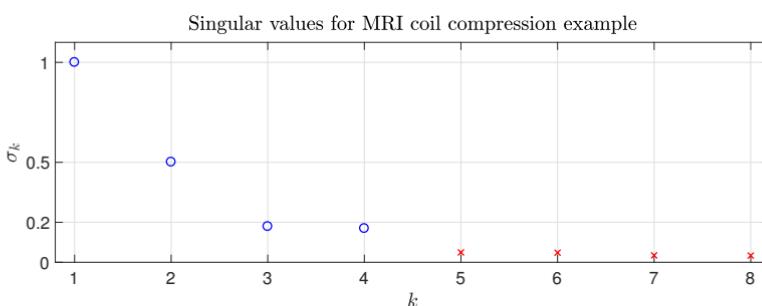
For implementation, we need to choose $K \leq r = \text{rank}(\mathbf{A}) \leq \min(M, N)$.

The approximation error is given by the 2-norm of the excluded singular values ($\sigma_{K+1}, \dots, \sigma_r$), so it is useful to plot all the singular values versus k and look for a “knee in the curve.”

Low-rank approximation is related to **factor analysis**, and a **scree plot** shows the singular values [5].

Example.

(Left figure from MRI example on p. 6.12; right figure from a textbook example with $r = 2$ and noise.)



You hope that your scree plot has a roughly linear section corresponding to the noise.

For moderate sized matrices, low-rank approximation needs just a few lines of JULIA code. The simplest approach uses the sum of outer products solution (6.2) literally:

```
using LinearAlgebra: svd
(U, s, V) = svd(A)
B = zeros(size(A))
for k=1:K
    global B += U[:,k] * s[k] * V[:,k]' # sum of rank-1 outer products
end
```

The `global` declaration is unnecessary if the code is within a function.

To avoid a loop one can use matrix multiplication corresponding to the **compact SVD** form in (6.2):

```
using LinearAlgebra: svd, Diagonal
(U, s, V) = svd(A) # default full=false saves memory
B = U[:,1:K] * Diagonal(s[1:K]) * V[:,1:K]'
```

Why small `s` above? Reminds me that output is a vector, not a matrix.

In both versions, the final `B` is \hat{A}_K .

Quality code would also include a bounds check that $1 \leq K \leq \min(M, N)$.

Yet another option is to compute only the first K SVD components using `svds`:

```
using LinearAlgebra: Diagonal
using Arpack: svds
tmp = svds(A, nsv=K, ritzvec=true)[1] # special SVD data structure
(U, s, Vt) = (tmp.U, tmp.S, tmp.Vt) # extract needed components
B = U * Diagonal(s) * Vt
```

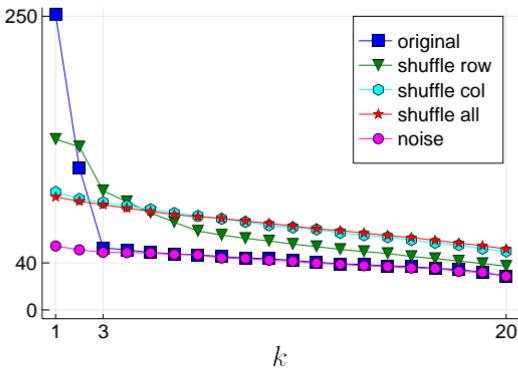
Why `Vt` above? Avoids an unnecessary transpose operation.

Choosing rank via permutation

Looking for a jump in the **scree plot** seems subjective. A more quantitative approach uses a **permutation method** [6].

Consider a matrix that is the sum of a low-rank component and IID random noise. If we permute (randomly) each of the columns (or rows? or all?) of that matrix, the noise part will remain the same (statistically) but the low-rank structure will be destroyed.

Now we plot the singular values of this shuffled matrix, which will basically correspond to noise. The singular values of the original data that are larger than those of the shuffled data might be considered to be the “significant” ones and can help guide rank choice. See [6] for theoretical analysis. One might need to permute multiple times and average to get a reliable baseline.



Example. This is from synthetic data that is rank 2 plus noise.

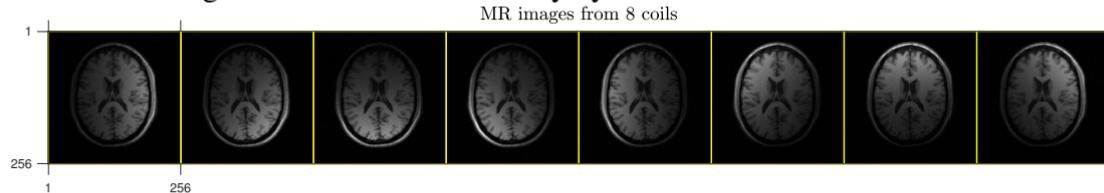
The two largest singular values stand out above the shuffled case.

Here is the key code for permuting each column:

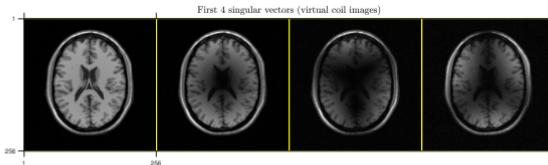
```
using Random: shuffle
Y = mapslices(shuffle, X, dims=1) # shuffle each column
Y = reshape(shuffle(X[:]), size(X)) # or, shuffle all
```

MRI Example

Modern MRI scanners use multiple receive coils to collect more image information at the same time. For example, the fMRI Center on UM's North Campus uses a 32-channel head coil. Processing all 32 sets of 3D images can require undesirably large computation times, and typically the data recorded by 32 coils is close to being linearly dependent. A technique called **coil compression** reduces data size [7, 8], essentially by making a **low-rank** approximation to the data and then just storing and processing the K singular vectors. Here are 8 MRI brain images recorded simultaneously by 8 coils around the head:



The scree plot on p. 6.8 suggests that most of the information is in the first 4 components. We reshape the data into a ($M = 256^2 \times N = 8$) matrix, use an SVD to find the rank-4 approximation, and reshape the first $K = 4$ vectors $\mathbf{u}_1, \dots, \mathbf{u}_4$ into 256×256 “virtual coil” images for display:



(In practice this operation is done on the raw data before making images.)



Non-uniqueness of SVD and low-rank approximation

As noted previously, “the” SVD of a matrix \mathbf{A} is not unique. In particular, if two singular values are the same, then one can swap the corresponding columns of \mathbf{U} and \mathbf{V} (or perform a more general rotation of their subspaces) and have two equally valid SVD expressions. Mathematically, when we write $\mathbf{A} = \sum_{k=1}^{\min(M,N)} \sigma_k \mathbf{u}_k \mathbf{v}'_k$, the order of the terms in the sum is unique if $\sigma_1 > \sigma_2 > \dots > \sigma_{\min(M,N)}$, but if any of the singular values are identical (including two or more zero singular values) then the ordering is not unique, and one can swap corresponding columns of \mathbf{U} and \mathbf{V} .

Example. Here are two different SVDs for a diagonal matrix with $7 = \sigma_1 > \sigma_2 = \sigma_3 = 5$:

$$\underbrace{\begin{bmatrix} 7 & 0 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 5 \end{bmatrix}}_{\mathbf{A}} = \underbrace{\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}}_{\mathbf{U}} \underbrace{\begin{bmatrix} 7 & 0 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 5 \end{bmatrix}}_{\Sigma} \underbrace{\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}}_{\mathbf{V}'}, \quad \underbrace{\begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}}_{\tilde{\mathbf{U}}} \underbrace{\begin{bmatrix} 7 & 0 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 5 \end{bmatrix}}_{\Sigma} \underbrace{\begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}}_{\tilde{\mathbf{V}}'}.$$

The example on the next page shows that there are two distinct rank $K = 2$ approximations to this \mathbf{A} .

Even if all the singular values are distinct, “the” SVD is still not unique because we can multiply both \mathbf{u}_k and \mathbf{v}_k by $e^{i\phi}$ and get a “different” \mathbf{U} and \mathbf{V} . However, the low-rank approximation is still the same in this case because $\sigma_k(e^{i\phi} \mathbf{u}_k)(e^{i\phi} \mathbf{v}_k)' = \sigma_k \mathbf{u}_k \mathbf{v}'_k$.

We have not focused on this issue too much yet because often uniqueness is relatively unimportant.

Now consider the question of whether the rank- K approximation of \mathbf{A} is **unique**. This is simpler to answer.

- If $\sigma_K > \sigma_{K+1}$ then one can extend the proof on p. 6.6 to show that the rank- K approximation is **unique**.
- If $\sigma_K = \sigma_{K+1}$, then the rank- K approximation is *not unique* because we could swap the K th and $(K+1)$ terms.

Example. Here are two distinct rank-1 approximations to a simple diagonal matrix:

$$\mathbf{A} \triangleq \begin{bmatrix} 5 & 0 \\ 0 & 5 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 5 & 0 \\ 0 & 5 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 5 & 0 \\ 0 & 5 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

$$\mathbf{A} \approx \mathbf{B}_1 \triangleq \sigma_1 \mathbf{u}_1 \mathbf{v}'_1 = 5 \begin{bmatrix} 1 \\ 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \end{bmatrix} = \begin{bmatrix} 5 & 0 \\ 0 & 0 \end{bmatrix}, \quad \mathbf{A} \approx \mathbf{B}_2 \triangleq \sigma_1 \tilde{\mathbf{u}}_1 \tilde{\mathbf{v}}'_1 = 5 \begin{bmatrix} 0 \\ 1 \end{bmatrix} \begin{bmatrix} 0 & 1 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 0 & 5 \end{bmatrix}.$$

What is the approximation error $\|\mathbf{A} - \mathbf{B}_1\|_{\text{F}}$ in the preceding example?

A: 0

B: $\sqrt{5}$

C: 5

D: 10

E: 5^2

??

In practice, any rank- K approximation (for suitable K) often is equally useful, so one infrequently worries about non-uniqueness.

Furthermore, in practical finite precision computing, it is dubious to test whether two floating-point values σ_K and σ_{K+1} are “exactly” equal. Avoid code like `if (a == b)` with floating-point numbers!

For analysis of the **numerical stability** of LR matrix approximation, including consideration of the “singular value gap” $\sigma_K - \sigma_{K+1}$, see [9].

1D example

Here we consider a 2×9 matrix where each column is an (x_n, y_n) coordinate for $n = 1, \dots, N = 9$, i.e.,

$$\mathbf{A} = \begin{bmatrix} x_1 & \dots & x_9 \\ y_1 & \dots & y_9 \end{bmatrix}.$$

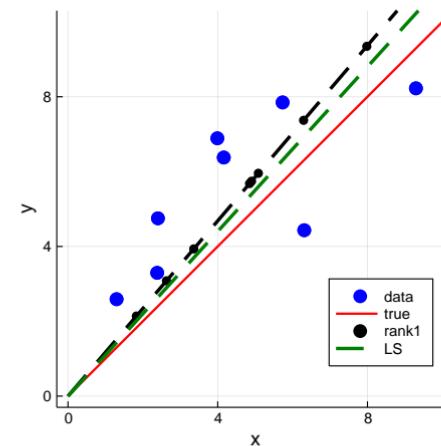
Based on the picture below, what is $\text{rank}(\mathbf{A})$? ??

https://web.eecs.umich.edu/~fessler/course/551/julia/demo/06_low_rank1.html

https://web.eecs.umich.edu/~fessler/course/551/julia/demo/06_low_rank1.ipynb

Suppose we want to perform **dimensionality reduction** of the 2D data by reducing it to a 1D line.

- One option for processing this data is to perform a **rank-1 approximation**, leading to dashed black line in the figure. The black points are the nearest points in the subspace for each of the data points.
- An alternative is to perform a **linear least-squares** fit assuming $y_n \approx \alpha x_n$. This leads to the green line with slope $\hat{\alpha}$ in the figure.

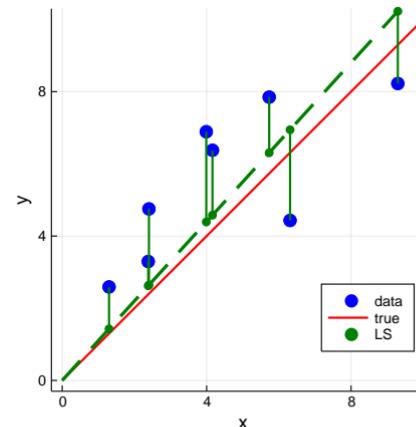
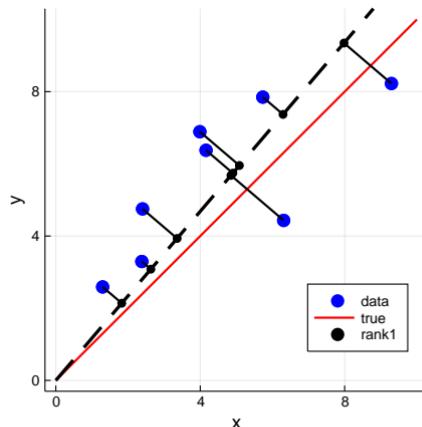


You might find it surprising that these two methods lead to different fits.

The **rank-1 approximation** is $\arg \min_{B \in \mathcal{L}_1^{M \times N}} \left\| \begin{bmatrix} x_1 & \dots & x_N \\ y_1 & \dots & y_N \end{bmatrix} - B \right\|_F^2 = \arg \min_{B \in \mathcal{L}_1^{M \times N}} \sum_{n=1}^N \left\| \begin{bmatrix} x_n \\ y_n \end{bmatrix} - B_{:,n} \right\|_2^2$.

The **linear least-squares** estimator is $\arg \min_{\alpha} \left\| \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} - \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} \alpha \right\|_2^2 = \arg \min_{\alpha} \sum_{n=1}^N |y_n - \alpha x_n|^2$.

These two methods measure approximation error differently, as illustrated in the following two plots.



Generalization to other norms

Theorem (**Eckart-Young-Mirsky**) (See [4] and [10] for a proof.)

For any **unitarily invariant** matrix norm $\|\cdot\|_{\text{UI}}$, the **low-rank approximation problem** has the same solution using the first (largest) K singular components of $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}' = \sum_{k=1}^r \sigma_k \mathbf{u}_k \mathbf{v}'_k$:

$$\hat{\mathbf{A}}_K \triangleq \arg \min_{\substack{\mathbf{B} \in \mathcal{L}_K^{M \times N}}} \underbrace{\|\mathbf{B} - \mathbf{A}\|_{\text{UI}}}_{\hookrightarrow \text{unitarily invariant norm}} = \sum_{k=1}^K \sigma_k \mathbf{u}_k \mathbf{v}'_k.$$

So even though the proof sketch on earlier pages was for the Frobenius norm, remarkably the same result holds (with a different proof) for other unitarily invariant norms such as the **spectral norm**.

To be robust to data outliers, sometimes one prefers other norms that are *not* unitarily invariant [11, 12], especially $\|\cdot\|_1$. Those other norms require different solution methods (typically iterative algorithms). SVD operations are still used as one part of many such algorithms [12].

What is $\|\hat{\mathbf{A}}_K - \mathbf{A}\|_2$ when \mathbf{A} has rank r ?

- A: $\sqrt{\sum_{k=K+1}^r \sigma_k^2}$ B: $\sum_{k=K+1}^r \sigma_k$ C: σ_{K+1} D: 0 E: None of these

??

What is $\|\hat{\mathbf{A}}_K - \mathbf{A}\|_*$? (Choose from same answer list.) ??

Proof for spectral norm

(Read)

Claim: $\text{rank}(\mathbf{B}) \leq K \implies \left\| \mathbf{A} - \hat{\mathbf{A}}_K \right\|_2^2 \leq \|\mathbf{A} - \mathbf{B}\|_2^2$ where $\hat{\mathbf{A}}_K \triangleq \sum_{k=1}^K \sigma_k \mathbf{u}_k \mathbf{v}'_k$.

Proof. Let $\mathbf{W} = [\mathbf{v}_1 \ \dots \ \mathbf{v}_{K+1}]$ so $\dim(\mathcal{R}(\mathbf{W})) = K + 1$ because columns of \mathbf{W} are orthonormal.

$\text{rank}(\mathbf{B}) \leq K \implies \dim(\mathcal{N}(\mathbf{B})) \geq N - K \implies \dim(\mathcal{N}(\mathbf{B})) + \dim(\mathcal{R}(\mathbf{W})) \geq (N - K) + (K + 1) = N + 1$. But both $\mathcal{N}(\mathbf{B})$ and $\mathcal{R}(\mathbf{W})$ are subspaces in \mathbb{F}^N , so they must have a nontrivial intersection, by (3.4).

Thus there exists some $\mathbf{x} \neq \mathbf{0}$ such that $\mathbf{x} \in \mathcal{N}(\mathbf{B})$ and $\mathbf{x} \in \mathcal{R}(\mathbf{W})$. WLOG take $\|\mathbf{x}\|_2 = 1$.

$\mathbf{x} \in \mathcal{R}(\mathbf{W}) \implies \mathbf{x} = \mathbf{W}\mathbf{z}$ for $\mathbf{z} \in \mathbb{F}^{K+1} \implies 1 = \|\mathbf{x}\|_2 = \|\mathbf{W}\mathbf{z}\|_2 = \|\mathbf{z}\|_2 = \|\mathbf{W}'\mathbf{W}\mathbf{z}\|_2 = \|\mathbf{W}'\mathbf{x}\|_2$.

To complete the proof:

$$\begin{aligned} \left\| \mathbf{A} - \hat{\mathbf{A}}_K \right\|_2^2 &= \left\| \sum_{k=K+1}^r \sigma_k \mathbf{u}_k \mathbf{v}'_k \right\|_2^2 = \sigma_{K+1}^2 = \sigma_{K+1}^2 \|\mathbf{x}\|_2^2 = \sigma_{K+1}^2 \|\mathbf{W}'\mathbf{x}\|_2^2 \text{ because } \|\mathbf{x}\|_2 = 1 = \|\mathbf{W}'\mathbf{x}\|_2 \\ &= \sigma_{K+1}^2 \sum_{k=1}^{K+1} \|\mathbf{v}'_k \mathbf{x}\|_2^2 \leq \sum_{k=1}^{K+1} \sigma_k^2 \|\mathbf{v}'_k \mathbf{x}\|_2^2 = \|\mathbf{A}\mathbf{x}\|_2^2 \text{ because } \mathbf{x} \in \mathcal{R}(\mathbf{W}) \\ &= \|(\mathbf{A} - \mathbf{B})\mathbf{x}\|_2^2 \text{ because } \mathbf{x} \in \mathcal{N}(\mathbf{B}) \\ &\leq \|\mathbf{A} - \mathbf{B}\|_2^2 \|\mathbf{x}\|_2^2 = \|\mathbf{A} - \mathbf{B}\|_2^2. \end{aligned}$$

□

Bases for $\mathbb{F}^{M \times N}$

Recall that a set of **linearly independent** vectors $\mathbf{b}_1, \mathbf{b}_2, \dots$ is a **basis** for a vector space \mathcal{V} iff

$$\text{span}(\{\mathbf{b}_1, \mathbf{b}_2, \dots\}) = \mathcal{V},$$

i.e., we can write every vector in \mathcal{V} as a linear combination of the basis vectors: $\mathbf{v} \in \mathcal{V} \implies \mathbf{v} = \sum_k \mathbf{b}_k \alpha_k$ for some $\alpha_k \in \mathbb{F}$. Now we consider two types of bases for the vector space of **matrices** $\mathbb{F}^{M \times N}$.

Canonical basis for $\mathbb{F}^{M \times N}$

First, the obvious (canonical) basis for $\mathbb{F}^{M \times N}$ is: $\mathcal{B} = \{\mathbf{e}_m \tilde{\mathbf{e}}'_n : m = 1, \dots, M, n = 1, \dots, N\}$

where \mathbf{e}_m denotes the m th unit vector of length M and $\tilde{\mathbf{e}}_n$ denotes the n th unit vector of length N . These are

clearly linearly independent and $\text{span}(\mathcal{B}) = \mathbb{F}^{M \times N}$ because $\mathbf{A} \in \mathbb{F}^{M \times N} \implies \mathbf{A} = \sum_{m=1}^M \sum_{n=1}^N a_{mn} \mathbf{e}_m \tilde{\mathbf{e}}'_n$.

To attempt “**dimensionality reduction**” using this basis, we could seek an approximation of the form

$\hat{\mathbf{B}} = \sum_{m=1}^M \sum_{n=1}^N b_{mn} \mathbf{e}_m \tilde{\mathbf{e}}'_n$, where we limit the number of nonzero coefficients b_{mn} .

Let $\|\text{vec}(\mathbf{B})\|_0$ denote the number of nonzero elements of matrix \mathbf{B} . Then a natural optimization problem is:

$$\hat{\mathbf{B}} = \underset{\mathbf{B} \in \mathbb{F}^{M \times N}, \|\text{vec}(\mathbf{B})\|_0 \leq K}{\arg \min} \|\mathbf{A} - \mathbf{B}\|_{\text{F}}.$$

This problem has a trivial solution: choose the K elements of \mathbf{A} having the largest magnitudes $|a_{mn}|$ and set

all other elements of $\hat{\mathbf{B}}$ to zero. The approximation error is the square root of the sum of squared magnitudes of all those other elements.

Example.

(Read)

$$\hat{\mathbf{B}} = \arg \min_{\mathbf{B} \in \mathbb{F}^{M \times N}, \|\mathbf{B}\|_0 \leq 3} \left\| \begin{bmatrix} 3 & 5 & 7 \\ 6 & 0 & 2 \end{bmatrix} - \mathbf{B} \right\|_{\text{F}} = \begin{bmatrix} 0 & 5 & 7 \\ 6 & 0 & 0 \end{bmatrix}.$$

The approximation error is $\left\| \mathbf{A} - \hat{\mathbf{B}} \right\|_{\text{F}} = \left\| \begin{bmatrix} 3 & 0 & 0 \\ 0 & 0 & 2 \end{bmatrix} \right\|_{\text{F}} = \sqrt{2^2 + 3^2}$.

How large must K be here for this approximation to have zero error?

K = 5 because \mathbf{A} has 5 nonzero elements.

How large must K be here for the rank- K approximation to have zero error?

K = 2 because \mathbf{A} has rank 2.

The **low-rank approximation** seems more parsimonious than approximation using \mathcal{B} .

To state this rigorously, we must first show that our low-rank approximation also uses a **basis** for $\mathbb{F}^{M \times N}$.

Orthonormal bases for $\mathbb{F}^{M \times N}$

Let \mathbf{U} denote any **unitary** $M \times M$ matrix and \mathbf{V} denote any **unitary** $N \times N$ matrix.

Now endow $\mathbb{F}^{M \times N}$ with the (Frobenius) **inner product**: $\langle \mathbf{A}, \mathbf{B} \rangle =$

and corresponding (Frobenius) norm: $\|\mathbf{A}\|_{\text{F}} = \sqrt{\langle \mathbf{A}, \mathbf{A} \rangle} = \sqrt{\text{trace}\{\mathbf{A}'\mathbf{A}\}}$.

Claim: the following set of MN rank-1 $M \times N$ outer-product matrices is an **orthonormal basis** for $\mathbb{F}^{M \times N}$:

$$\{\mathbf{u}_m \mathbf{v}'_n : m = 1, \dots, M, n = 1, \dots, N\}.$$

Proof of orthonormality (for the Frobenius inner product):

$$\begin{aligned}\langle \mathbf{u}_m \mathbf{v}'_n, \mathbf{u}_k \mathbf{v}'_l \rangle &= \text{trace}\{(\mathbf{u}_k \mathbf{v}'_l)' (\mathbf{u}_m \mathbf{v}'_n)\} = \text{trace}\{\mathbf{v}_l \mathbf{u}'_k (\mathbf{u}_m \mathbf{v}'_n)\} = \text{trace}\{\mathbf{u}'_k \mathbf{u}_m \mathbf{v}'_n \mathbf{v}_l\} \\ &= (\mathbf{u}'_k \mathbf{u}_m)(\mathbf{v}'_n \mathbf{v}_l) = \begin{cases} 1, & k = m, l = n \\ 0, & \text{otherwise.} \end{cases}\end{aligned}$$

Being orthonormal, the set $\{\mathbf{u}_m \mathbf{v}'_n\}$ spans a MN -dimensional space and is a basis for $\mathbb{F}^{M \times N}$. Thus one can write any $\mathbf{X} \in \mathbb{F}^{M \times N}$ in terms of that basis, i.e., $\mathbf{X} = \mathbf{U} \mathbf{C} \mathbf{V}' = \sum_{m=1}^M \sum_{n=1}^N c_{mn} \mathbf{u}_m \mathbf{v}'_n$, where $\mathbf{C} \triangleq \mathbf{U}' \mathbf{X} \mathbf{V}$, consistent with (6.4).

When working with a matrix \mathbf{A} , using the basis for $\mathbb{F}^{M \times N}$ that is formed from its own left and right **singular vectors**, i.e., \mathbf{U} and \mathbf{V} , respectively, turns out to be provide the most parsimonious representation. Note that the original low-rank problem formulation (6.1) did not involve any SVD, but an **SVD** arose in the solution.

Does the canonical basis $\{\mathbf{e}_m \tilde{\mathbf{e}}'_n : m = 1, \dots, M, n = 1, \dots, N\}$ form an orthonormal basis for $\mathbb{F}^{M \times N}$ when using the Frobenius inner product?

A: Yes

B: No

C: Insufficient information

??

Low-rank approximation summary

The low-rank approximation problem (6.1) and its SVD-based solution (6.2) provides the best low-rank *representation* (or approximation) of a given matrix (in the Frobenius norm sense). Specifically, if \mathbf{B} is any rank- K matrix having the same size as a matrix $\mathbf{A} = \sum_{k=1}^r \mathbf{u}_k \mathbf{v}'_k$, then:

$$\sqrt{\sum_{k=K+1}^r \sigma_k^2} = \left\| \mathbf{A} - \hat{\mathbf{A}}_K \right\|_{\text{F}} = \left\| \mathbf{A} - \sum_{k=1}^K \sigma_k \mathbf{u}_k \mathbf{v}'_k \right\|_{\text{F}} \leq \left\| \mathbf{A} - \mathbf{B} \right\|_{\text{F}}, \quad \forall \mathbf{B} \in \mathcal{L}_K^{M \times N}.$$

Generalizations

(Read)

What if the data is corrupted by noise? What if

$$\mathbf{Y} = \mathbf{X} + \boldsymbol{\varepsilon}$$

where \mathbf{X} is a low-rank matrix (or close to low rank). Due to the noise matrix $\boldsymbol{\varepsilon}$, typically \mathbf{Y} is not low-rank. We really want to recover \mathbf{X} here, not just represent or approximate \mathbf{Y} . How do we do that? An answer called OptShrink is given in [13] and discussed on p. 6.50.

Low-rank approximation is related closely to **principal component analysis (PCA)**, discussed on p. 6.57. There are many interesting **PCA generalizations** including robust methods, nonlinear models, multi-linear (tensor) models, **nonnegative matrix factorizations (NMF)** [14], methods that use sparsity, ...

Rank and stability

(Read)

Although **rank** is a simple mathematical concept, it is not stable numerically. Precisely, it is not a continuous function of the elements of a matrix.

Example. Consider the following two matrices:

$$\mathbf{A} = \begin{bmatrix} 1.7 & 1.7 \\ 1.7 & 1.7 \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} 1.7 & 1.7 \\ 1.7 + \varepsilon & 1.7 \end{bmatrix}, \quad 0 < |\varepsilon| \ll 1.$$

These two matrices are nearly indistinguishable numerically, yet $\text{rank}(\mathbf{A}) = 1$ and $\text{rank}(\mathbf{B}) = 2$.

An alternative matrix property used in some situations is the **stable rank** [15], defined (for $\mathbf{A} \neq 0$) as:

$$\frac{\|\mathbf{A}\|_{\text{F}}^2}{\|\mathbf{A}\|_2^2}, \text{ where } 1 \leq \frac{\|\mathbf{A}\|_{\text{F}}^2}{\|\mathbf{A}\|_2^2} = \frac{1}{\sigma_1^2} \sum_{k=1}^{\min(M,N)} \sigma_k^2 = 1 + \sum_{k=2}^{\min(M,N)} \frac{\sigma_k^2}{\sigma_1^2} \leq r,$$

because (see HW) $\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_{\text{F}} \leq \sqrt{r} \|\mathbf{A}\|_2$. For the above example, $\|\mathbf{B}\|_{\text{F}}^2 / \|\mathbf{B}\|_2^2 \approx 1 + \varepsilon/2$.

Example. For $\mathbf{A} = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$, the stable rank is $\frac{\|\mathbf{A}\|_{\text{F}}^2}{\|\mathbf{A}\|_2^2} = \frac{2^2 + 1^2}{2^2} = 1.25 \in [1, 2]$.

What is the **stable rank** of $\mathbf{A} = \mathbf{x}\mathbf{y}'$?

A: 0

B: 1

C: $\|\mathbf{x}\|_2 \|\mathbf{y}\|_2$ D: $\|\mathbf{x}\|_2^2 \|\mathbf{y}\|_2^2$

??

6.2 Sensor localization application (Multidimensional scaling)

We now turn to another application of matrix factorization: **sensor localization** via **multidimensional scaling**. Suppose $J \geq d$ sensors are located at unknown locations $\mathbf{c}_1, \dots, \mathbf{c}_J \in \mathbb{R}^d$, where typically $d = 2$ or $d = 3$. All we are given is the $J \times J$ **distance matrix** \mathbf{D} having elements

$$d_{ij} = d_{ji} =$$

Clearly by definition the diagonal elements of \mathbf{D} are zero: $d_{jj} = 0$.

Goal: given \mathbf{D} , determine the locations $\{\mathbf{c}_j\}$.

Limitation: there is a fundamental ambiguity because translating or rotating coordinates does not change \mathbf{D} . Specifically, if $\tilde{\mathbf{c}}_j = \mathbf{Q}\mathbf{c}_j + \mathbf{d}$ where \mathbf{Q} is a $d \times d$ orthogonal matrix (such as a rotation matrix) and $\mathbf{d} \in \mathbb{R}^d$ is a displacement vector, then

$$\|\tilde{\mathbf{c}}_i - \tilde{\mathbf{c}}_j\|_2 = \|(\mathbf{Q}\mathbf{c}_i + \mathbf{d}) - (\mathbf{Q}\mathbf{c}_j + \mathbf{d})\|_2 = \|\mathbf{Q}(\mathbf{c}_i - \mathbf{c}_j)\|_2 = \|\mathbf{c}_i - \mathbf{c}_j\|_2$$

because the Euclidean norm is **unitarily invariant**. So \mathbf{D} would be the same for $\{\tilde{\mathbf{c}}_j\}$ and $\{\mathbf{c}_j\}$.

Thus we must be content with determining locations $\{\mathbf{c}_j\}$ to within some translation and rotation factor.

Note the choice to use a matrix representation of the data!

And of course the locations are vectors.

Derivation (analysis) ---

To derive a solution (for noiseless data), define another matrix S by the square of the elements of D :

$$s_{ij} \triangleq d_{ij}^2 = \|\mathbf{c}_i - \mathbf{c}_j\|_2^2 = \|\mathbf{c}_i\|_2^2 + \|\mathbf{c}_j\|_2^2 - 2 \langle \mathbf{c}_i, \mathbf{c}_j \rangle.$$

Naturally we write S in matrix form:

$$\begin{aligned} S &= \begin{bmatrix} s_{11} & \dots & s_{1J} \\ \vdots & \dots & \vdots \\ s_{J1} & \dots & s_{JJ} \end{bmatrix} = \begin{bmatrix} \|\mathbf{c}_1\|_2^2 & \dots & \|\mathbf{c}_1\|_2^2 \\ \vdots & \dots & \vdots \\ \|\mathbf{c}_J\|_2^2 & \dots & \|\mathbf{c}_J\|_2^2 \end{bmatrix} + \begin{bmatrix} \|\mathbf{c}_1\|_2^2 & \dots & \|\mathbf{c}_J\|_2^2 \\ \vdots & \dots & \vdots \\ \|\mathbf{c}_1\|_2^2 & \dots & \|\mathbf{c}_J\|_2^2 \end{bmatrix} - 2 \begin{bmatrix} \langle \mathbf{c}_1, \mathbf{c}_1 \rangle & \dots & \langle \mathbf{c}_1, \mathbf{c}_J \rangle \\ \vdots & \dots & \vdots \\ \langle \mathbf{c}_J, \mathbf{c}_1 \rangle & \dots & \langle \mathbf{c}_J, \mathbf{c}_J \rangle \end{bmatrix} \\ &= \begin{bmatrix} \|\mathbf{c}_1\|_2^2 \\ \vdots \\ \|\mathbf{c}_J\|_2^2 \end{bmatrix} \mathbf{1}'_J + \mathbf{1}_J [\|\mathbf{c}_1\|_2^2 \ \dots \ \|\mathbf{c}_J\|_2^2] - 2 \begin{bmatrix} \mathbf{c}'_1 \\ \vdots \\ \mathbf{c}'_J \end{bmatrix} [\mathbf{c}_1 \ \dots \ \mathbf{c}_J] \\ &= \mathbf{r} \mathbf{1}'_J + \mathbf{1}_J \mathbf{r}' - 2 \mathbf{C}' \mathbf{C}, \quad \mathbf{r} \triangleq \begin{bmatrix} \|\mathbf{c}_1\|_2^2 \\ \vdots \\ \|\mathbf{c}_J\|_2^2 \end{bmatrix}, \quad \mathbf{C} \triangleq \underbrace{[\mathbf{c}_1 \ \dots \ \mathbf{c}_J]}_{d \times J \text{ unknown locations}}. \end{aligned} \tag{6.5}$$

Unfortunately we do not know the values in the vector \mathbf{r} .

So we know the $J \times J$ matrix \mathbf{S} and we want to solve for the $d \times J$ unknown location matrix \mathbf{C} using

$$\mathbf{S} = \mathbf{r}\mathbf{1}'_J + \mathbf{1}_J\mathbf{r}' - 2\mathbf{C}'\mathbf{C}. \quad (6.6)$$

Because \mathbf{S} on the left is symmetric and zero along its diagonal, it contains $(J^2 - J)/2$ distinct known values. The right-hand side depends on the dJ unknown values in \mathbf{C} .

Recall that source localization has a translation ambiguity, so we may as well use a coordinate system where the centroid is $\mathbf{0}$, *i.e.*, we can assume $\mathbf{C}\mathbf{1}_J = \sum_{j=1}^J \mathbf{c}_j = \mathbf{0}$ without losing any further generality.

Define the following “de-meaning” operator that projects onto the **orthogonal complement** of $\mathcal{R}(\mathbf{1}_J)$:

$$\mathbf{P}^\perp \triangleq$$

Multiplying \mathbf{P}^\perp by any vector produces a new vector whose mean value is zero, so we say it “de-means” the input vector. (The English word “demean” has a very different definition.) In winter, sometimes an airplane must be “de-iced” by removing ice from the wings. Here **de-mean** has similar use: removing the mean.

Clearly $\mathbf{P}^\perp \mathbf{1}_J = \mathbf{0}$, so $\mathbf{1}'_J \mathbf{P}^\perp = \mathbf{0}'$ and $\mathbf{C}\mathbf{P}^\perp = \mathbf{C}(\mathbf{I} - \frac{1}{J}\mathbf{1}_J\mathbf{1}'_J) = \mathbf{C}$ so from (6.6):

$$\mathbf{P}^\perp \mathbf{S} \mathbf{P}^\perp =$$

Rearranging leads to the following simple expression for the $J \times J$ **Gram matrix**:

$$\mathbf{G} \triangleq \mathbf{C}'\mathbf{C} =$$

In words, we remove the row and column means of \mathbf{S} and divide by -2 .

Now we have the **Gram matrix** $\mathbf{G} = \mathbf{C}'\mathbf{C}$, but we really want the coordinates matrix \mathbf{C} itself.

Although $\mathbf{C}'\mathbf{C}$ is a $J \times J$ matrix, $d \leq J$ so typically $\text{rank}(\mathbf{C}'\mathbf{C}) = d$ (in the absence of noise). If the locations happen to be linearly dependent, *i.e.*, if the sensors are along one line through the origin in 2D or in the same plane through the origin in 3D, then $\text{rank}(\mathbf{C}'\mathbf{C}) < d$.

To elaborate, consider $d = 2$ where

$$\mathbf{C} \triangleq [\mathbf{c}_1 \ \dots \ \mathbf{c}_J] = \begin{bmatrix} x_1 & \dots & x_J \\ y_1 & \dots & y_J \end{bmatrix},$$

where $\mathbf{c}_j = (x_j, y_j)$. Because \mathbf{C} has two rows, $\text{rank}(\mathbf{C}'\mathbf{C}) = \text{rank}(\mathbf{C}) \leq 2$.

When does $\text{rank}(\mathbf{C}) = 1$? Only if the two rows of \mathbf{C} are linearly dependent, *i.e.*, if $\mathbf{y} = \alpha\mathbf{x}$ for some $\alpha \in \mathbb{R}$, *i.e.*, if $y_j = \alpha x_j$ for all j , which is the equation for points along a line through the origin with slope α .

The matrix $\mathbf{C}'\mathbf{C}$ consists of inner products, and it is invariant to rotations of the source locations:

$$\tilde{\mathbf{C}} = \mathbf{Q}\mathbf{C} \implies \tilde{\mathbf{C}}'\tilde{\mathbf{C}} =$$

Given the $J \times J$ Gram matrix $\mathbf{G} \triangleq \mathbf{C}'\mathbf{C} = -\frac{1}{2}\mathbf{P}^\perp S\mathbf{P}^\perp$, we determine its **compact SVD**:

$$\mathbf{G} = \mathbf{V}\Sigma\mathbf{V}' = \sum_{k=1}^J \sigma_k \mathbf{v}_k \mathbf{v}_k' =$$

$$\Sigma_d \triangleq$$

$$\mathbf{V}_d \triangleq$$

- Because \mathbf{G} is (symmetric) positive semidefinite (in the absence of noise), we know $\mathbf{U} = \mathbf{V}$.
- Because \mathbf{G} has rank at most d (in the absence of noise), the above **low-rank** form is *exact*, not an approximation like we used earlier in this chapter.

We define our estimate of the source locations to be the J columns of this $d \times J$ matrix:

$$\hat{\mathbf{C}} = \Sigma_d^{1/2} \mathbf{V}_d' = \underbrace{\begin{bmatrix} \sqrt{\sigma_1} & & \\ & \ddots & \\ & & \sqrt{\sigma_d} \end{bmatrix}}_{d \times d} \underbrace{\begin{bmatrix} - & \mathbf{v}_1' & - \\ \vdots & \vdots & \vdots \\ - & \mathbf{v}_d' & - \end{bmatrix}}_{d \times J},$$

because this estimate $\hat{\mathbf{C}}$ is consistent with \mathbf{C} to within an unknown rotation and translation:

$$\hat{\mathbf{C}}'\hat{\mathbf{C}} = (\Sigma_d^{1/2} \mathbf{V}_d')' (\Sigma_d^{1/2} \mathbf{V}_d) = \mathbf{V}_d \Sigma_d \mathbf{V}_d' =$$

Now it seems like we are done because we have an expression for \mathbf{G} that is computable from \mathbf{S} , so we can use its eigendecomposition to find the source location estimates $\hat{\mathbf{C}}$. However, we assumed that $\mathbf{C}\mathbf{1}_J = \mathbf{0}$ so we must verify that $\hat{\mathbf{C}}\mathbf{1}_J = \mathbf{0}$ to ensure that our approach is self consistent.

Because $\mathbf{P}^\perp \mathbf{1}_J = \mathbf{0}$, it follows that $\mathbf{G}\mathbf{1}_J = \mathbf{0}$, so $\mathbf{V}\Sigma\mathbf{V}'\mathbf{1}_J = \mathbf{V}_d\Sigma_d\mathbf{V}'_d\mathbf{1}_J = \mathbf{0}$, which in turn implies that $\hat{\mathbf{C}}\mathbf{1}_J = \Sigma_d^{1/2}\mathbf{V}'_d\mathbf{1}_J = \mathbf{0}$.

What is the size of the final $\mathbf{0}$ in the immediately preceding line?

A: $J \times J$ B: $J \times 1$ C: $J \times d$ D: $d \times J$ E: $d \times 1$

??

Multidimensional scaling

Summary of method for finding locations from distances (in the possible presence of noise):

- Given distances arranged in a matrix \mathbf{D} .
- Compute \mathbf{S} from \mathbf{D} by squaring elements.
- Compute $\hat{\mathbf{G}} \triangleq -\frac{1}{2}\mathbf{P}^\perp \mathbf{S} \mathbf{P}^\perp \approx \mathbf{C}'\mathbf{C}$ by de-meaning.
- Optional: force $\hat{\mathbf{G}}$ to be symmetric by replacing it with $(\hat{\mathbf{G}} + \hat{\mathbf{G}}')/2$
- Compute the **compact SVD** $\hat{\mathbf{G}} = \mathbf{U}_d\Sigma_d\mathbf{V}'_d$. (If $\hat{\mathbf{G}}$ is symmetric then $\mathbf{U}_d = \mathbf{V}_d$.)
- Use rank- d terms for source location estimate: $\hat{\mathbf{C}} = \Sigma_d^{1/2}\mathbf{V}'_d$.

This remarkably simple algorithm for finding $\hat{\mathbf{C}}$ is called (classical) **multidimensional scaling** [16, Ch. 12], and using the low-rank compact SVD is a key step. Because $\mathbf{G} = \mathbf{C}'\mathbf{C}$, \mathbf{G} is positive semidefinite (in the absence of noise). If we force $\hat{\mathbf{G}}$ to be likewise symmetric, then its eigenvalues are all real and nonnegative so alternatively one could use an eigendecomposition of $\hat{\mathbf{G}}$ with suitably ordered eigenvalues.

Questions

Consider a given noiseless distance matrix D with $J > d$. The position estimates \hat{C} returned by the SVD-based MDS algorithm are unique (to within numerical precision), *i.e.*, are the same for any SVD of the Gram matrix.

A: True

B: False

??

In the possible presence of noise, the product $\hat{C}'\hat{C}$ is unique (to within numerical precision) for any SVD of the calculated Gram matrix \hat{G} .

A: never

B: usually

C: always

??

Practical implementation

The `MultivariateStats` package of JULIA has a `classical_mds` command.

But the method is so simple that it is just as easy to code our own.

Example. Consider $J = 3$ sources that are all equally distant: $S = D .^2 = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix}$.

Just a few lines of JULIA suffices (see notebook linked on next page):

```
using Plots; using LinearAlgebra: svd, Diagonal; using Statistics: mean
D = [0 1 1; 1 0 1; 1 1 0] + 0e-9 * randn(3,3) # given distances
S = D.^2
tmp = S .- mean(S, dims=1)
G = -0.5 * (tmp .- mean(tmp, dims=2))
(_, s, V) = svd(G)
C = Diagonal(sqrt.(s[1:2])) * V[:,1:2]'
scatter(C[1,:], C[2,:], aspect_ratio=1, label="",
       xlabel="x", ylabel="y", title="s=$(s)",
       xtick=-.5:.5:.5, ytick=-.5:.5:.5, xlim=(-0.6,0.6), ylim=(-0.6,0.6))
```

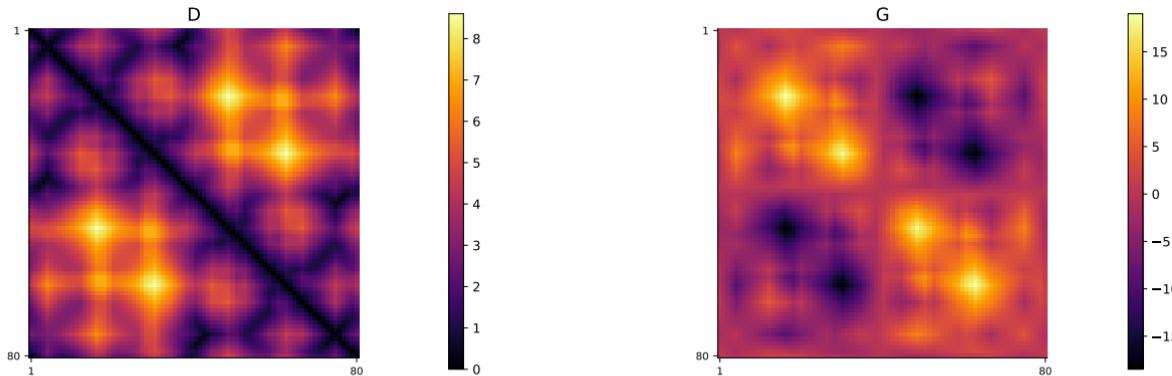
What “shape” should appear? ??

Example. An unknown collection of $J = 80$ source points:

https://web.eecs.umich.edu/~fessler/course/551/julia/demo/06_source_local1.html

https://web.eecs.umich.edu/~fessler/course/551/julia/demo/06_source_local1.ipynb

Images of D and G :



Run demo files to show \hat{C} and consider noisy case to examine robustness.

Extensions

- When the distance measurements are noisy then D may not be perfectly symmetric and $P^\perp S P^\perp$ may have (hopefully small) negative eigenvalues.
Use **truncated SVD**: truncate rank to dimension d of embedding.
- What if D is **geodesic distances** or distances along roads between locations?
Distortion in location estimates.
- What if some sensors cannot reach some other sensors, *i.e.*, if D has some missing elements?
- Is there a way to formulate the problem that considers the possibility of noise (errors in the distance values) at the outset, instead of simply empirically investigating the effects of noise on our solution.
- What if there are outliers in the data, *e.g.*, some sensors that give incorrect values (either due to errors or because the sensor was hacked or otherwise compromised). Is MDS robust to such errors, or is a different formulation needed, perhaps based on a robust distance measure like $\|\cdot\|_1$?
- If we want to compare the location estimates \hat{C} to some ground truth, we can use the Procrustes method to align them, compensating for the inherent rotation/translation ambiguity.

6.3 Proximal operators

Our next topic is some alternative low-rank approximation methods, also formulated as optimization problems. To describe those alternatives, we first need to introduce an important operation.

Define. The **proximal operator** [17, 18] (also called the **proximal mapping** [19, Ch. 6]) associated with a (typically **convex**) function $f : \mathbb{F}^N \mapsto \mathbb{R}$, is defined, for any $\mathbf{v} \in \mathbb{F}^N$, as the following minimizer:

$$\text{prox}_f(\mathbf{v}) \triangleq \arg \min_{\mathbf{x} \in \mathbb{F}^N} \frac{1}{2} \|\mathbf{v} - \mathbf{x}\|_2^2 + f(\mathbf{x}). \quad (6.7)$$

- The norm squared is **strictly convex**, so when f is convex the “arg min” is unique for any $\mathbf{v} \in \mathbb{F}^N$.
- Usually the function f is additively separable, in which case the “arg min” separates into N individual 1D minimization problems. See example below.
- The proximal operator transforms one function $f(\cdot)$ into another function $\text{prox}_f(\cdot)$.

In general, the functions f and prox_f have the same (domain?, range?)

A: true,true

B: true,false

C: false,true

D: false,false

??

??

Example. The main case used in EECS 551 is when $\mathbb{F} = \mathbb{R}$ and $f(\mathbf{x}) = \beta \|\mathbf{x}\|_p^p = \beta \sum_{n=1}^N |x_n|^p$, for $\beta \geq 0$. The most important case is when $p = 1$, for which

$$\text{prox}_f(\mathbf{v}) = \arg \min_{\mathbf{x} \in \mathbb{R}^N} \frac{1}{2} \|\mathbf{v} - \mathbf{x}\|_2^2 + \beta \|\mathbf{x}\|_1 = \arg \min_{\mathbf{x} \in \mathbb{R}^N} \sum_{n=1}^N \left(\frac{1}{2} (v_n - x_n)^2 + \beta |x_n| \right).$$

This sum is additively separable, meaning each x_n appears in just one term of the sum. So we can solve the minimization problem separately for each x_n term by solving

$$\hat{x}_n \triangleq \arg \min_{x_n \in \mathbb{R}} g(x_n), \quad g_n(x_n) \triangleq \frac{1}{2} (v_n - x_n)^2 + \beta |x_n|, \quad n = 1, \dots, N.$$

Often for minimization problems we take the derivative and set it to zero. That approach does not quite work here because $|\cdot|$ is not differentiable at 0. So we need a “braces and cases” approach.

Suppose the minimizer satisfies $\hat{x}_n > 0$. In this regime the function g is differentiable with derivative

$$\dot{g}_n(x_n) = x_n - v_n + \beta 1.$$

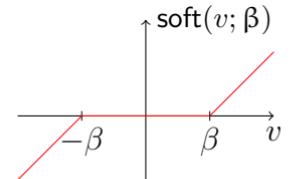
Equating to zero yields $\hat{x}_n = v_n - \beta$. However, we must keep in mind that we assumed here that $\hat{x}_n > 0$, so this solution is valid only when $v_n > \beta$.

Now suppose the minimizer satisfies $\hat{x}_n < 0$. In this regime the function g is differentiable with derivative

$$\dot{g}_n(x_n) = x_n - v_n + \beta(-1).$$

Equating to zero yields $\hat{x}_n = v_n + \beta$. We must keep in mind that we assumed here that $\hat{x}_n < 0$, so this solution is valid only when $v_n < -\beta$. Because we have considered all cases where the minimizer is nonzero, for any other input value v_n , the minimizer must be zero. Thus we have shown that

$$\arg \min_{x_n \in \mathbb{R}} \frac{1}{2}(v_n - x_n)^2 + \beta |x_n| = \text{soft}(v_n; \beta), \quad \text{soft}(v_n; \beta) \triangleq \begin{cases} v_n - \beta, & v_n > \beta \\ v_n + \beta, & v_n < -\beta \\ 0, & \text{otherwise.} \end{cases}$$



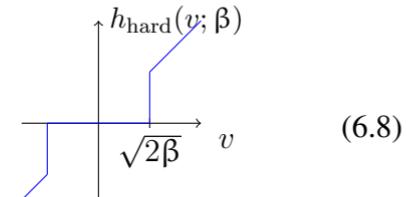
This operation is called **soft thresholding** and it induces **sparsity** because of the range of values set to zero. Putting it all together, we have shown that the proximal operator for the 1-norm scaled by β is element-wise soft thresholding:

$$\text{prox}_{\beta \|\cdot\|_1}(\mathbf{v}) = \text{soft} . (\mathbf{v}; \beta).$$

Example. Another main case of interest here is the (nonconvex) 0-norm $f(\mathbf{x}) = \beta \|\mathbf{x}\|_0$, for $\beta \geq 0$, for which the **proximal operator** is (element-wise) **hard thresholding**. Here is the derivation:

$$\text{prox}_{\beta \|\cdot\|_0}(\mathbf{v}) = \arg \min_{\mathbf{x} \in \mathbb{F}^N} \frac{1}{2} \|\mathbf{v} - \mathbf{x}\|_2^2 + \beta \|\mathbf{x}\|_0 = \arg \min_{\mathbf{x} \in \mathbb{F}^N} \sum_{n=1}^N \left(\frac{1}{2} |v_n - x_n|^2 + \beta \mathbb{I}_{\{x_n \neq 0\}} \right) = h_{\text{hard}}(\mathbf{v}; \beta)$$

$$h_{\text{hard}}(v; \beta) \triangleq \arg \min_{x \in \mathbb{R}} \frac{1}{2} |v - x|^2 + \beta \mathbb{I}_{\{x \neq 0\}} = \begin{cases} v, & \frac{1}{2} |v|^2 > \beta \\ 0, & \text{otherwise.} \end{cases} \quad (6.8)$$



Again we cannot derive the minimizer by simple differentiation, because the indicator function is not differentiable at zero. So we must apply “braces and cases.”

If the minimizer is nonzero, then the second term is β and the minimizer of the first term is v , and the overall cost is $0 + \beta = \beta$. If the minimizer is 0, then the first term is $(v - 0)^2/2$ and the second term is 0, so the cost is $v^2/2$. So 0 is the minimizer when $v^2/2 < \beta$, otherwise v is the minimizer.

The minimizer is not unique if $|v| = \sqrt{2\beta}$; in such cases, both v and 0 are global minimizers. The threshold is $\sqrt{2\beta}$ due to the $1/2$ in front of the norm in the proximal operator definition.

Now we return to low-rank approximation where we will use these results.

6.4 Alternative low-rank approximation formulations

If $\mathbf{Y} \in \mathbb{F}^{M \times N}$ is a given data matrix with SVD $\mathbf{Y} = \mathbf{U}\Sigma\mathbf{V}'$ and $\mathbf{Y} = \mathbf{X} + \boldsymbol{\varepsilon}$ where we believe $\text{rank}(\mathbf{X}) \leq K$ and $\boldsymbol{\varepsilon}$ denotes a $M \times N$ noise matrix, then the constrained formulation discussed so far is natural [3]:

$$\hat{\mathbf{X}} = \arg \min_{\mathbf{X} \in \mathcal{L}_K^{M \times N}} \|\mathbf{Y} - \mathbf{X}\|_F^2 = \mathbf{U}_K \Sigma_K \mathbf{V}'_K = \sum_{k=1}^K \sigma_k \mathbf{u}_k \mathbf{v}'_k = \sum_{k=1}^r \sigma_k \mathbb{I}_{\{k \leq K\}} \mathbf{u}_k \mathbf{v}'_k.$$

This problem is nonconvex, but still has a nice solution (6.2), reiterated above in three different forms. The final form reminds us that we have “set to 0” all the singular values for $k > K$.

Unconstrained / regularized formulation

When the rank is unknown, an alternative approach is to consider an *unconstrained* cost function that discourages (penalizes) high rank (high complexity) rather than constraining it [20]:

$$\hat{\mathbf{X}} = \arg \min_{\mathbf{X} \in \mathbb{F}^{M \times N}} h_{\text{hard}}(\|\mathbf{Y} - \mathbf{X}\|_F; \beta) \quad (6.9)$$

where (see proof sketch below) $h_{\text{hard}}(\cdot; \beta)$ is the hard-thresholding function defined in (6.8).

General unitarily invariant formulations

There are several ways to form low-rank approximations that balance between data-fit and model complexity, all having the general form:

$$\hat{\mathbf{X}} = \arg \min_{\mathbf{X} \in \mathbb{F}^{M \times N}} \quad \text{(6.10)}$$

for some **unitarily invariant** matrix norm $\|\cdot\|$ and for some **regularizer** $R(\mathbf{X})$, where $R : \mathbb{F}^{M \times N} \mapsto \mathbb{R}$, such as $R(\mathbf{X}) = \text{rank}(\mathbf{X})$. The regularization parameter (aka **hyperparameter**) $\beta > 0$ controls the trade-off between fit to the data \mathbf{Y} , as measured by $\|\mathbf{Y} - \mathbf{X}\|$, and model complexity, as quantified by $R(\mathbf{X})$.

We assume hereafter that the regularizer is also a unitarily invariant function, *i.e.*, $R(\mathbf{U}\mathbf{X}\mathbf{V}) = R(\mathbf{X})$ for *any* suitably sized unitary matrices \mathbf{U} and \mathbf{V} . Because both the data-fit norm and the regularizer are unitarily invariant, using the symmetric gauge principles of [4, 10], one can show that any minimizer has the form

$$\hat{\mathbf{X}} = \sum_{k=1}^r \hat{w}_k \mathbf{u}_k \mathbf{v}'_k = \mathbf{U}_r \hat{\Sigma}_r \mathbf{V}'_r \text{ where } \hat{\Sigma}_r = \text{diag}\{\hat{w}_k\}, \quad \hat{w}_k = h_k(\sigma_1, \dots, \sigma_r; \beta), \quad \mathbf{Y} = \mathbf{U}_r \Sigma_r \mathbf{V}'_r, \quad \text{(6.11)}$$

for some **shrinkage** or **thresholding** function $h_k(\cdot; \beta)$ that depends on the data-fit norm and the regularizer.

What is the rank of the $\hat{\mathbf{X}}$ in (6.11) when $\text{rank}(\mathbf{Y}) = r$? (Choose best answer.)

- A: Usually 1 B: Always 1 C: Usually r D: Always r E: None of these

??

Accepting that (6.11) is the form of the solution, we can rewrite the general optimization problem (6.10) as

$$\hat{\mathbf{X}} = \mathbf{U}_r \hat{\Sigma}_r \mathbf{V}'_r, \quad \hat{\Sigma}_r = \arg \min_{\mathbf{S}=\text{diag}\{s_1, \dots, s_r\}} \quad (6.12)$$

Singular value hard thresholding

Now return to the special case (6.9) with the Frobenius norm and the rank penalty. The equivalent minimization problem in terms of the singular values requires minimizing over $\{s_1, \dots, s_r\}$:

(6.13)

We started with the cost function (6.9), simplified it to the equivalent problem with diagonal matrices (6.12), and now we have r separate 1D problems that we solved earlier when deriving the proximal operator for the 0-norm that resulted in hard thresholding; see (6.8).

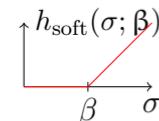
This analysis leads to $\hat{w}_k = h_{\text{hard}}(\sigma_k; \beta)$ in the method (6.9) that is called **singular value hard thresholding**.

Singular value soft thresholding

Any formulation involving $\text{rank}(\cdot)$ is non-convex, and it is somewhat remarkable that a nice solution to (6.9) exists despite that non-convexity. Often it can be preferable to have a **convex** formulation. The **convex relaxation** [21] of rank is the **nuclear norm** that leads to **soft thresholding** of singular values:

$$\hat{\mathbf{X}} = \quad \quad \quad (6.14)$$

$$\text{So } \hat{w}_k = h_{\text{soft}}(\sigma_k; \beta), \quad h_{\text{soft}}(\sigma; \beta) \triangleq [\sigma - \beta]_+ = \begin{cases} \sigma - \beta, & \sigma > \beta, \\ 0, & \text{otherwise.} \end{cases}$$



This approach is called **singular value soft thresholding (SVST)**.

Proof sketch: $\frac{1}{2}\|\Sigma_r - \mathbf{S}\|_{\text{F}}^2 + \beta\|\mathbf{S}\|_* = \sum_{k=1}^r \frac{1}{2}(\sigma_k - s_k)^2 + \beta s_k \implies \hat{w}_k = \arg \min_s \frac{1}{2}(\sigma_k - s)^2 + \beta |s|$.

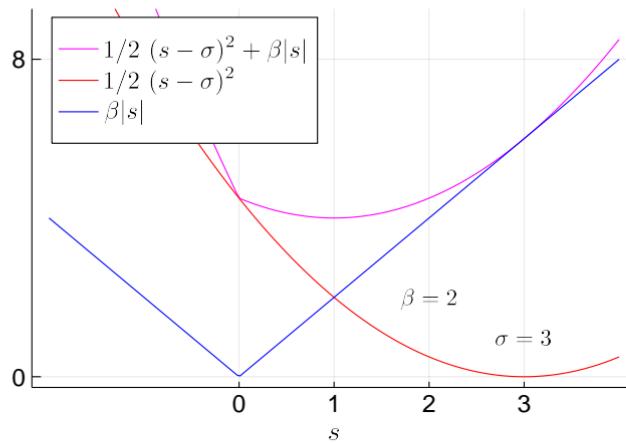
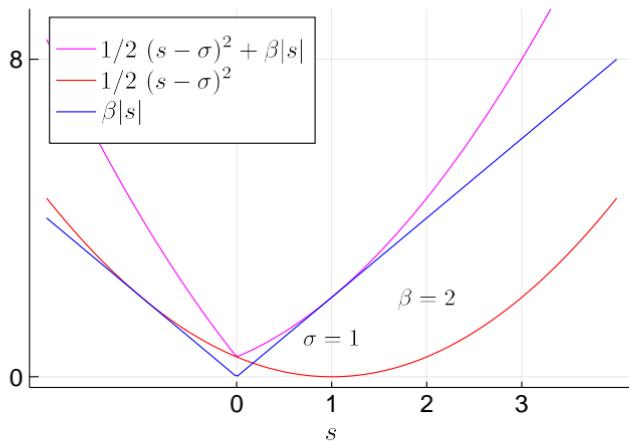
Sketching we see that the minimizer will be some $s \geq 0$ because $\sigma_k \geq 0$.

Differentiating: $0 = (s - \sigma_k) + \beta \Big|_{s=\hat{w}_k} \implies \hat{w}_k = \sigma_k - \beta$ but we need $\hat{w}_k \geq 0$ so $\hat{w}_k = [\sigma_k - \beta]_+$.

The following figures illustrate the cases where $\sigma_k < \beta$ and $\sigma_k > \beta$, by plotting the following functions:

$$\frac{1}{2}(s - \sigma)^2 + \beta |s|.$$

- When $0 < \beta \leq \sigma$ the minimizer is at $s = \sigma - \beta$.
- When $0 \leq \sigma \leq \beta$, the minimizer is at $s = 0$.



One can generalize this to complex numbers and show that

$$\arg \min_{s \in \mathbb{C}} \frac{1}{2} |s - z|^2 + \beta |z| = [|\bar{z}| - \beta]_+ e^{i\angle z}.$$

If \mathbf{Y} has nonzero singular values 3, 5, 6, 7, 9 and $\beta = 6$, what is the rank of $\hat{\mathbf{X}}$ here?

- A: 1 B: 2 C: 3 D: 4 E: 5

??

This $\hat{\mathbf{X}}$ equals the rank- K approximation to \mathbf{Y} , where K is answer to previous problem. (?)

- A: True B: False

??

In general, what is $\|\hat{\mathbf{X}}\|_2$ in terms of the notation used throughout here?

- A: β B: σ_1 C: $\sigma_1 - \beta$ D: $\sqrt{\sigma_r - \beta}$ E: $[\sigma_1 - \beta]_+$

??

In general, if $\|\mathbf{Y}\|_2 \geq \beta$, then what is $\|\mathbf{Y} - \hat{\mathbf{X}}\|_2$ in terms of the notation used throughout here?

- A: β B: σ_1 C: $\sigma_1 - \beta$ D: $\sqrt{\sigma_r - \beta}$ E: $[\sigma_1 - \beta]_+$

??

One motivation for these unconstrained formulations is for solving **matrix completion** or **matrix sensing** problems where we must replace $\|\mathbf{Y} - \mathbf{X}\|_F^2$ with something like $\|\mathbf{y} - \mathbf{A} \text{vec}(\mathbf{X})\|_2^2$ where \mathbf{A} is a known matrix. These problems do not have closed-form solutions and require iterative methods, and convergence of iterative methods is better understood for convex problems.

Other extensions of low-rank approximation

- Sometimes we might want to retain some columns of the data exactly, leading to a different type of constrained low-rank approximation [22].
- Another variation replaces the Frobenius norm by the spectral norm: (HW)

$$\arg \min_{\mathbf{X}} \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\|_F^2 + \beta \|\mathbf{X}\|_*.$$

- Another variation replaces the Frobenius norm by the nuclear norm:

$$\arg \min_{\mathbf{X}} \|\mathbf{Y} - \mathbf{X}\|_* + \beta \|\mathbf{X}\|_*.$$

The solution to this variation is left as an exercise. It seems not to be useful.

- For further generalizations using weighted norms, see [23–25].

If we use two Frobenius norms, the solution has the following general form:

$$\hat{\mathbf{X}} = \arg \min_{\mathbf{X}} \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\|_{\text{F}}^2 + \beta \frac{1}{2} \|\mathbf{X}\|_{\text{F}}^2 = \sum_{k=1}^r \hat{w}_k \mathbf{u}_k \mathbf{v}'_k, \quad \hat{w}_k = h_{\text{FF}}(\sigma_k, \beta).$$

What is $h_{\text{FF}}(\sigma, \beta)$?

- A: $[\sigma - \beta]_+$ B: $\sigma H(\sigma - \sqrt{\beta})$ C: $\max(\sigma, \beta)$ D: $\frac{\sigma}{1 + \beta}$ E: $(\sigma - \beta)^2$

??

If \mathbf{Y} is $M \times N$ with rank r , then what is the rank of the solution $\hat{\mathbf{X}}$ here?

- A: 0 B: 1 C: r D: $\min(N, M)$ E: None of these.

??

6.5 Choosing the rank or regularization parameter

In all the above low-rank approximation formulations, we must either

- choose the rank K directly, or
- choose the regularization parameter β that influences the rank of $\hat{\mathbf{X}}$.

This selection process is non-trivial because when $\mathbf{Y} = \mathbf{X} + \varepsilon$ and the **latent** matrix \mathbf{X} has low-rank:

- $\left\| \hat{\mathbf{X}} - \mathbf{Y} \right\|_{\text{F}}$: the difference between the noisy data \mathbf{Y} and the low-rank approximation $\hat{\mathbf{X}}$ always
 - decreases monotonically as K increases (discretely)
 - increases monotonically as β increases (continuously).
- $\left\| \hat{\mathbf{X}} - \mathbf{X} \right\|_{\text{F}}$: the difference between the low-rank matrix \mathbf{X} and the low-rank approximation $\hat{\mathbf{X}}$:
 - usually decreases initially as K increases, then increases
 - usually decreases initially as β increases, then increases.

The figures on the next page illustrate these properties.

Often we normalize the difference between $\hat{\mathbf{X}}$ and \mathbf{Y} or that between $\hat{\mathbf{X}}$ and \mathbf{X} and define the (unitless) **normalized root mean-squared difference (NRMSD)** and **normalized root mean-squared error (NRMSE)**:

$$\text{NRMSD} = \frac{\left\| \hat{\mathbf{X}} - \mathbf{Y} \right\|_{\text{F}}}{\left\| \mathbf{Y} \right\|_{\text{F}}} \cdot 100\%, \quad \text{NRMSE} = \frac{\left\| \hat{\mathbf{X}} - \mathbf{X} \right\|_{\text{F}}}{\left\| \mathbf{X} \right\|_{\text{F}}} \cdot 100\%.$$

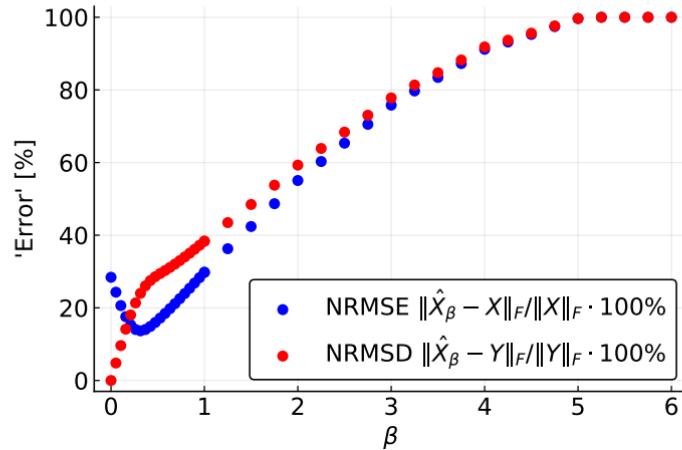
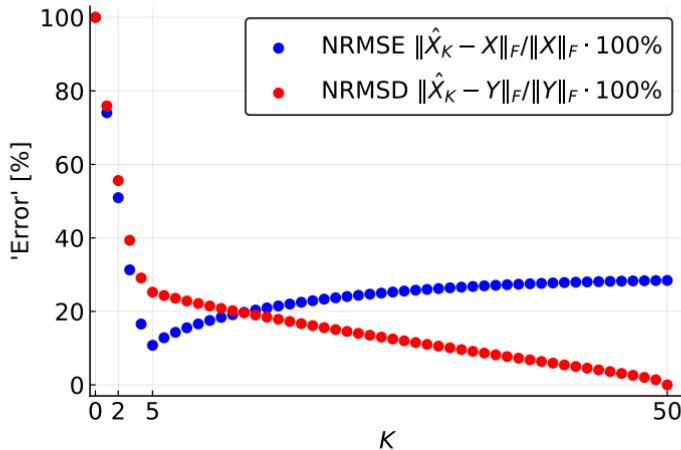
(There are many variations of these definitions in the literature.)

Example. In this case the 100×50 matrix \mathbf{X} has rank 5 and $\mathbf{Y} = \mathbf{X} + \boldsymbol{\varepsilon}$ where $\boldsymbol{\varepsilon}$ is $M \times N$ matrix with IID Gaussian noise.

See JULIA demo:

https://web.eecs.umich.edu/~fessler/course/551/julia/demo/06_lr_sure1.html

https://web.eecs.umich.edu/~fessler/course/551/julia/demo/06_lr_sure1.ipynb



Ideally we would like to find K_* or β_* that minimizes $\left\| \hat{\mathbf{X}} - \mathbf{X} \right\|_F$, the minimizers of the blue curves above. Finding either of the choices exactly is impossible because \mathbf{X} is unknown in practice!

All we have is $\mathbf{Y} = \mathbf{X} + \boldsymbol{\varepsilon}$, and usually we are *hoping* that \mathbf{X} is low rank, without knowing for sure.

Researchers have developed surrogates for $\left\| \hat{\mathbf{X}} - \mathbf{X} \right\|_F$ that help choose β .

One notable method is **Stein's unbiased risk estimate (SURE)**. If $\hat{\mathbf{x}}(\mathbf{y}; \beta)$ is a **weakly differentiable** estimate of a parameter $\mathbf{x} \in \mathbb{R}^d$ from data with Gaussian noise: $\mathbf{y} \sim N(\mathbf{x}, \sigma_0^2 \mathbf{I})$, then an unbiased estimate of the **risk**, *i.e.*, the mean-squared error (**MSE**), is:

$$\text{MSE}\{\beta\} = E[\|\hat{\mathbf{x}}(\mathbf{y}; \beta) - \mathbf{x}\|^2] = \text{SURE}\{\beta\} \triangleq \underbrace{\|\hat{\mathbf{x}}(\mathbf{y}; \beta) - \mathbf{y}\|_2^2 - d\sigma_0^2 + 2\sigma_0^2 \sum_i \frac{\partial}{\partial y_i} \hat{\mathbf{x}}_i(\mathbf{y}; \beta)}_{\text{Independent of } \mathbf{x}!}.$$

Divergence of $\hat{\mathbf{x}}(\cdot; \beta)$

For a proof, see [\[wiki\]](#). Note that σ_0^2 is a noise variance, not a singular value!

The SURE approach uses the approximation to select β as follows:

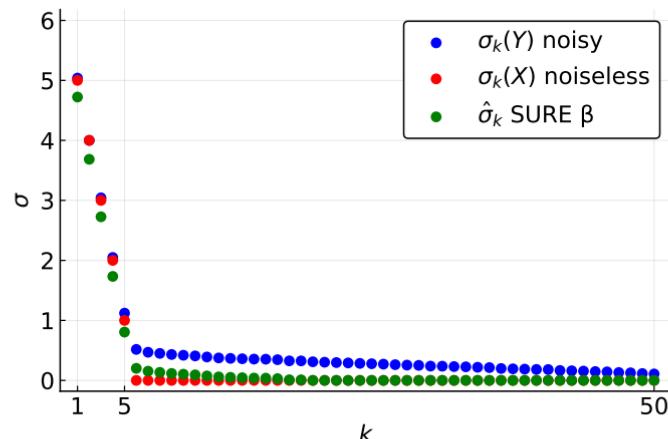
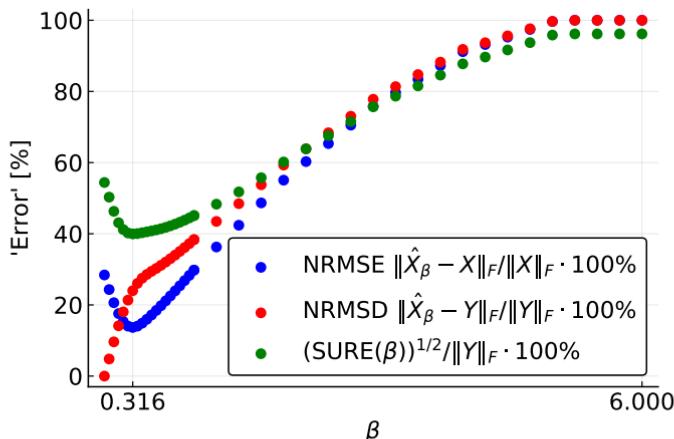
$$\beta_* \triangleq \arg \min_{\beta} \|\hat{\mathbf{x}}_{\beta} - \mathbf{x}\|^2 \approx \arg \min_{\beta} E[\|\hat{\mathbf{x}}_{\beta} - \mathbf{x}\|^2] = \arg \min_{\beta} \text{SURE}\{\beta\}.$$

The hard thresholding function is *not* weakly differentiable, so we cannot apply the SURE method to rank constrained or rank regularized cases. The soft thresholding function *is* weakly differentiable. Remarkably the divergence expression derived in [\[20\]](#) for any method of the form $\hat{\mathbf{X}} = \sum_k \mathbf{u}_k h_k(\sigma_k; \beta) \mathbf{v}'_k$ turns out to

depend only on the singular values so it is practical to evaluate. When the singular values of \mathbf{Y} are distinct:

$$\begin{aligned} \text{SURE}\{\beta\} &= \left\| \hat{\mathbf{X}} - \mathbf{Y} \right\|_{\text{F}}^2 - MN\sigma_0^2 \\ &\quad + 2\sigma_0^2 \left(|M-N| \sum_{i=1}^{\min(M,N)} \frac{h(\sigma_i; \beta)}{\sigma_i} + \sum_{i=1}^{\min(M,N)} h_i(\sigma_i; \beta) + 2 \sum_{i \neq j}^{\min(M,N)} \frac{\sigma_i h_i(\sigma_i; \beta)}{\sigma_i^2 - \sigma_j^2} \right). \end{aligned}$$

See the demo notebook. A practical challenge is that one must know σ_0^2 to apply this method.



The minimizer β_* of the SURE function is remarkably close to the minimum MSE choice of β .

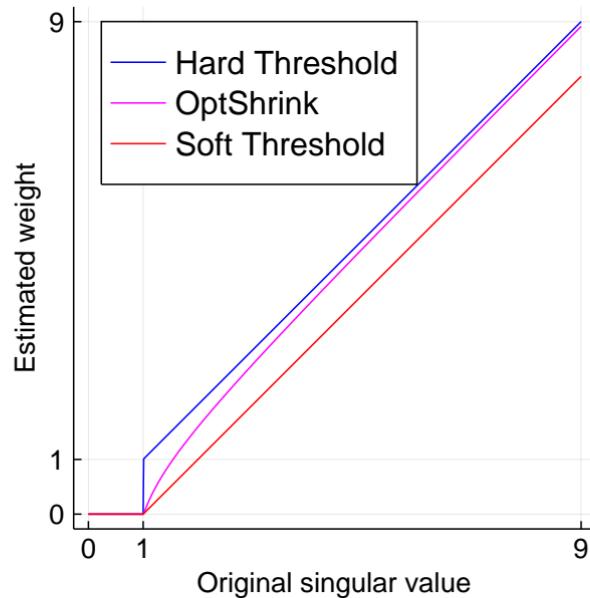
OptShrink

3002

IEEE TRANSACTIONS ON INFORMATION THEORY, VOL. 60, NO. 5, MAY 2014

OptShrink: An Algorithm for Improved Low-Rank Signal Matrix Denoising by Optimal, Data-Driven Singular Value Shrinkage

Raj Rao Nadakuditi, Member, IEEE



Model: $\mathbf{Y} = \mathbf{X} + \boldsymbol{\varepsilon}$. Random matrix theory says, for quite general noise models (bi-unitarily invariant), the best estimate of \mathbf{X} is $\hat{\mathbf{X}} = \sum_{k=1}^r \hat{w}_k \mathbf{u}_k \mathbf{v}'_k$, where \hat{w}_k is the best “weight” for the k th SVD component [13].

Algorithm 1 OptShrink: A New Algorithm for Low-Rank Matrix Denoising by Optimal, Data-Driven Singular Value Shrinkage

- 1: Input: $\tilde{X} = n \times m$ signal-plus-noise matrix
 - 2: Input: \hat{r} = Estimate of the effective rank of the latent low-rank signal matrix
 - 3: Compute $\tilde{X} = \sum_{i=1}^q \hat{\sigma}_i \hat{u}_i \hat{v}_i^H$
 - 4: Compute $\hat{\Sigma}_{\hat{r}} = \text{diag}(\hat{\sigma}_{\hat{r}+1}, \dots, \hat{\sigma}_q) \in \mathbb{R}^{(n-\hat{r}) \times (m-\hat{r})}$
 - 5: **for** $i = 1, \dots, \hat{r}$ **do**
 - 6: Compute $\hat{D}(\hat{\sigma}_i; \hat{\Sigma}_{\hat{r}})$ using (16a) and $\hat{D}'(\hat{\sigma}_i; \hat{\Sigma}_{\hat{r}})$ using (16b)
 - 7: Compute $\hat{w}_{i,\hat{r}}^{\text{opt}} = -2 \frac{\hat{D}(\hat{\sigma}_i; \hat{\Sigma}_{\hat{r}})}{\hat{D}'(\hat{\sigma}_i; \hat{\Sigma}_{\hat{r}})}$
 - 8: **end for**
 - 9: **return** $\hat{S}_{\text{opt}} = \sum_{i=1}^{\hat{r}} \hat{w}_{i,\hat{r}}^{\text{opt}} \hat{u}_i \hat{v}_i^H$ = denoised estimate of the rank \hat{r} signal matrix
 - 10: **return** (optional) Compute estimate of MSE using (17a)
 - 11: **return** (optional) Compute estimate of relative MSE using (17b)
-

For a matrix $X \in \mathbb{K}^{n \times m}$. Define

$$\begin{aligned}\hat{D}(z; X) &:= \frac{1}{n} \text{Tr} \left(z (z^2 I - XX^H)^{-1} \right) \\ &\quad + \frac{1}{m} \text{Tr} \left(z (z^2 I - X^H X)^{-1} \right)\end{aligned}\quad (16a)$$

and

$$\begin{aligned}\hat{D}'(z; X) &:= \frac{1}{n} \text{Tr} \left[z (z^2 I - XX^H)^{-1} \right] \\ &\quad + \frac{1}{m} \text{Tr} \left[-2z^2 (z^2 I - X^H X)^{-2} + (z^2 I - X^H X)^{-1} \right] \\ &\quad + \frac{1}{m} \text{Tr} \left[z (z^2 I - X^H X)^{-1} \right] \\ &\quad + \frac{1}{n} \text{Tr} \left[-2z^2 (z^2 I - XX^H)^{-2} + (z^2 I - XX^H)^{-1} \right].\end{aligned}\quad (16b)$$

Note: “ X ” in (16a) and (16b) is the $(n - \hat{r}) \times (m - \hat{r})$ (rectangular) diagonal matrix $\hat{\Sigma}_{\hat{r}}$.

The online code <http://web.eecs.umich.edu/~rajnrao/optshrink/> and the equations above are practical only if both n and m are sufficiently small. Here we adapt the approach to allow one of n or m (but not both) to be large.

Let \mathbf{S} denote a $K \times L$ rectangular diagonal matrix with entries $s_1, \dots, s_{\min(K,L)}$. We want to evaluate:

$$D(z, \mathbf{S}) = \frac{1}{K} \text{trace}\{z(z^2 \mathbf{I} - \mathbf{S}\mathbf{S}')^{-1}\} \frac{1}{L} \text{trace}\{z(z^2 \mathbf{I} - \mathbf{S}'\mathbf{S})^{-1}\}.$$

Note that $D(z, \mathbf{S}) = D(z, \mathbf{S}')$ so WLOG assume $K \geq L$ (tall). Then $\mathbf{S} = \begin{bmatrix} \mathbf{S}_L \\ \mathbf{0} \end{bmatrix}$ and

$$\begin{aligned} z^2 \mathbf{I} - \mathbf{S}\mathbf{S}' &= z^2 \mathbf{I} - \begin{bmatrix} \mathbf{S}_L^2 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} = \begin{bmatrix} z^2 \mathbf{I} - \mathbf{S}_L^2 & \mathbf{0} \\ \mathbf{0} & z^2 \mathbf{I} \end{bmatrix} \\ \implies \text{trace}\{z(z^2 \mathbf{I} - \mathbf{S}\mathbf{S}')^{-1}\} &= \text{trace}\left\{z \begin{bmatrix} z^2 \mathbf{I} - \mathbf{S}_L^2 & \mathbf{0} \\ \mathbf{0} & z^2 \mathbf{I} \end{bmatrix}^{-1}\right\} = \sum_{k=1}^L \frac{z}{z^2 - s_k^2} + \frac{K-L}{z}. \end{aligned}$$

Similarly, as long as $z \neq s_k$:

$$\text{trace}\{z(z^2 \mathbf{I} - \mathbf{S}'\mathbf{S})^{-1}\} = \sum_{k=1}^L \frac{z}{z^2 - s_k^2}, \quad \implies D(z, \mathbf{S}) = \frac{1}{KL} \left(\sum_{k=1}^L \frac{z}{z^2 - s_k^2} + \frac{K-L}{z} \right) \left(\sum_{k=1}^L \frac{z}{z^2 - s_k^2} \right).$$

A similar simplification is feasible for $D'(z, \mathbf{S})$. (HW)

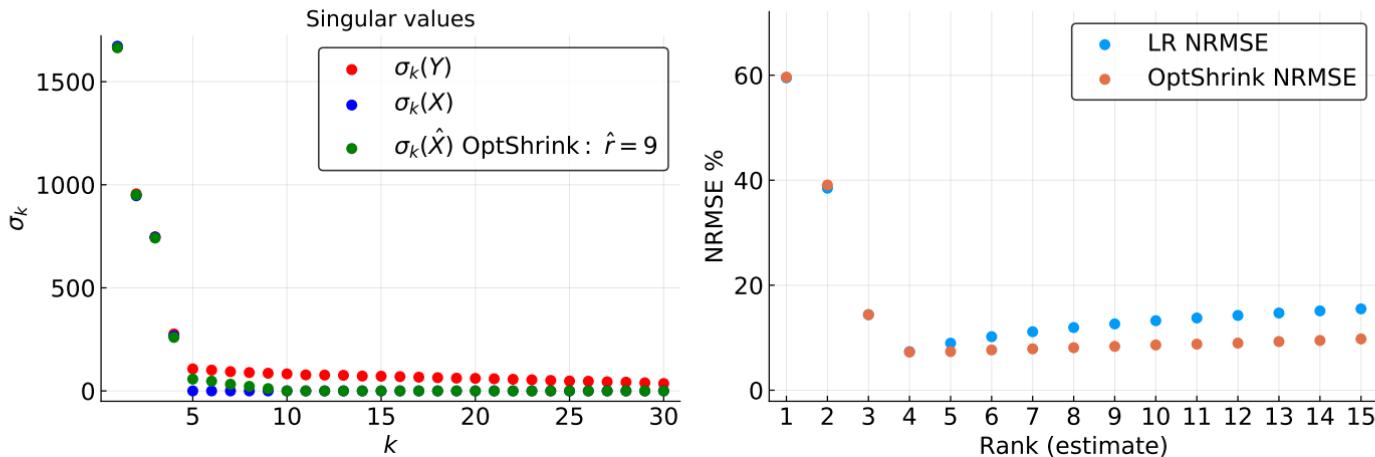
Combining, one can implement OptShrink to calculate the optimal SVD component weights when at most one of m or n is large.

Example showing OptShrink's robustness to rank estimate

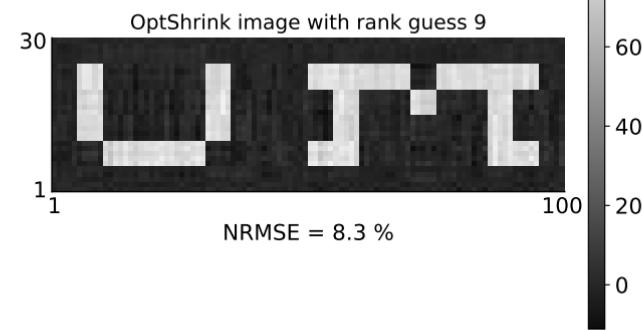
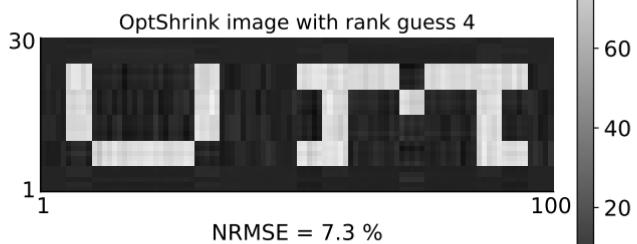
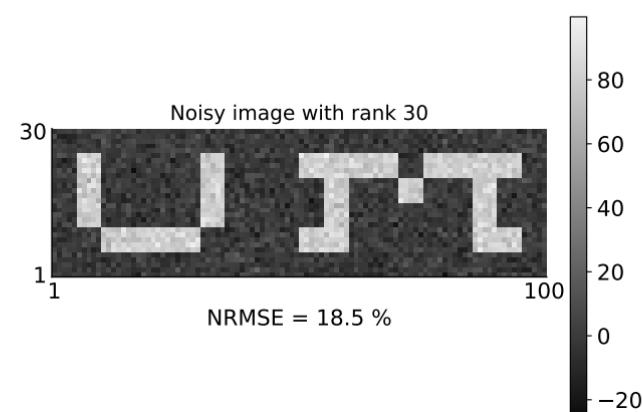
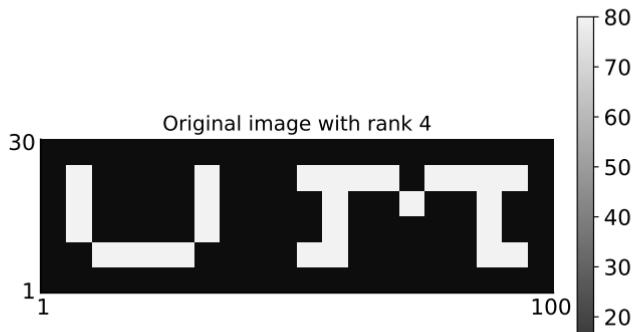
\mathbf{X} is a 100×30 logo image (aka 2D array aka matrix) with true rank 4. The noisy data is $\mathbf{Y} = \mathbf{X} + \boldsymbol{\varepsilon}$.

https://web.eecs.umich.edu/~fessler/course/551/julia/demo/06_optshrink1.html

https://web.eecs.umich.edu/~fessler/course/551/julia/demo/06_optshrink1.ipynb



The error of $\hat{\mathbf{X}}$ for conventional low-rank approximation increases faster (due to over-fitting noise) than that of OptShrink. OptShrink finds the “threshold” for singular value shrinkage *adaptively* from the data, namely from the tail singular values, using random matrix theory.



6.6 Related methods: autoencoders and PCA

Relation to autoencoder with linear hidden layer

An **autoencoder** is a type of **artificial neural network (ANN)** that is designed to perform **dimensionality reduction**. An autoencoder consists of an **encoder** and a **decoder**:

$$\mathbf{x} \in \mathbb{F}^M \rightarrow \boxed{\text{encoder } \phi} \rightarrow \mathbf{z} \in \mathbb{F}^K \rightarrow \boxed{\text{decoder } \psi} \rightarrow \hat{\mathbf{x}} = \psi(\phi(\mathbf{x})) \in \mathbb{F}^M,$$

where typically $K \ll M$, *i.e.*, the dimension of the **code** or **latent variable** \mathbf{z} is much less than that of a data point \mathbf{x} .

The simplest possible case is where the encoder and decoder are each linear, and hence each can be represented via a **matrix**:

$$\phi(\mathbf{x}) = \mathbf{A}\mathbf{x}, \quad \mathbf{A} \in \mathbb{F}^{K \times M} \text{ (wide)}$$

$$\psi(\mathbf{z}) = \mathbf{B}\mathbf{z}, \quad \mathbf{B} \in \mathbb{F}^{M \times K} \text{ (tall).}$$

In this case, the output is $\hat{\mathbf{x}} = \mathbf{B}\mathbf{A}\mathbf{x}$ and the natural approach to **training** this ANN given training data $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{F}^M$ is

$$(\hat{\mathbf{A}}, \hat{\mathbf{B}}) = \arg \min_{\mathbf{A} \in \mathbb{F}^{K \times M}, \mathbf{B} \in \mathbb{F}^{M \times K}} \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{B}\mathbf{A}\mathbf{x}_n\|_2^2 = \arg \min_{\mathbf{A} \in \mathbb{F}^{K \times M}, \mathbf{B} \in \mathbb{F}^{M \times K}} \|\mathbf{X} - \mathbf{B}\mathbf{A}\mathbf{X}\|_{\text{F}}^2, \quad \mathbf{X} \triangleq [\mathbf{x}_1 \ \dots \ \mathbf{x}_N].$$

Note that the product \mathbf{BAX} has at most rank K due to the dimensions of \mathbf{A} and \mathbf{B} . So we can use the low-rank approximation work to observe that one solution to this training problem is

$$\mathbf{A} = \mathbf{U}'_K, \quad \mathbf{B} = \mathbf{U}_K,$$

where \mathbf{U}_K is the first K left singular vectors of \mathbf{U} , and an **SVD** is $\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}'$. For that choice of \mathbf{A} and \mathbf{B} , and assuming $K \leq \min(M, N)$, the product becomes

$$\begin{aligned} \mathbf{BAX} &= \mathbf{U}_K \mathbf{U}'_K \underbrace{\mathbf{U}\Sigma\mathbf{V}'}_{\mathbf{X}} = \mathbf{U}_K \mathbf{U}'_K [\mathbf{U}_K \quad \mathbf{U}_{:, (K+1):M}] \Sigma \mathbf{V}' \\ &= \mathbf{U}_K [\mathbf{I}_K \quad \mathbf{0}_{K \times M}] \begin{bmatrix} \Sigma_K & \mathbf{0}_{K \times (N-K)} \\ \mathbf{0}_{(M-K) \times K} & \Sigma_{(M-K) \times (N-K)} \end{bmatrix} [\mathbf{V}_K \quad \mathbf{V}_{:, (K+1):N}]' = \mathbf{U}_K \Sigma_K \mathbf{V}'_K = \mathbf{X}_K, \end{aligned}$$

namely the optimal rank at most K approximation to \mathbf{X} .

So performing **low-rank approximation** of training data is comparable to learning an autoencoder with a single linear hidden layer of reduced dimension.

Nonlinear autoencoders are simply generalizations with nonlinear operations such as $\phi(\mathbf{x}) = \sigma(\mathbf{W}\mathbf{x} + \mathbf{b})$ where σ is some nonlinear function, \mathbf{W} is a $M \times K$ matrix and $\mathbf{b} \in \mathbb{F}^M$ is **bias vector**.

Exercise. Revisit the training derivation above in the case where the encoder and decoder are **affine** instead of linear, i.e., $\phi(\mathbf{x}) = \mathbf{Ax} + \mathbf{b}_1$, $\psi(\mathbf{z}) = \mathbf{Bz} + \mathbf{b}_2$.

Relation to principal component analysis (PCA)

So far we have focused on low-rank approximation of a matrix. Often the low-dimensional subspaces $\text{span}(\mathbf{U}_K)$ and $\text{span}(\mathbf{V}_K)$ are more important to us than the overall approximation $\hat{\mathbf{A}}_K = \mathbf{U}_K \Sigma_K \mathbf{V}'_K$.

The **principal component analysis (PCA)** approach focuses on learning a linear transform of the data (having **orthonormal** columns) that maximizes variance. PCA starts with a set of N training vectors $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{F}^M$ and first subtracts from each the mean of all training vectors to form a first de-meaned data matrix:

$$\mathbf{X}_1 \triangleq [\mathbf{x}_1 - \boldsymbol{\mu} \quad \dots \quad \mathbf{x}_N - \boldsymbol{\mu}] = \mathbf{X} - \mathbf{X} \frac{1}{N} \mathbf{1}_N \mathbf{1}'_N, \quad \text{where } \mathbf{X} \triangleq [\mathbf{x}_1 \quad \dots \quad \mathbf{x}_N], \quad \boldsymbol{\mu} \triangleq \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n.$$

Note that $\mathbf{X}_1 \mathbf{1}_N = \mathbf{0}_M$. Now PCA seeks the (unit norm) linear combination of rows of \mathbf{X} (elements of \mathbf{x}) that maximizes (empirical) variance as follows:

$$\arg \max_{\mathbf{u}: \|\mathbf{u}\|=1} \frac{1}{N} \sum_{n=1}^N |\mathbf{u}' \mathbf{X}_1[:, n]|^2 = \arg \max_{\mathbf{u}: \|\mathbf{u}\|=1} \|(\mathbf{u}' \mathbf{X}_1)'\|_2 = \arg \max_{\mathbf{u}: \|\mathbf{u}\|=1} \|\mathbf{X}'_1 \mathbf{u}\|_2 = \mathbf{u}_1,$$

where an **SVD** of \mathbf{X}_1 is $\mathbf{X}_1 = \mathbf{U} \Sigma \mathbf{V}'$. Another way of writing this uses a **Rayleigh quotient**:

$$\arg \max_{\mathbf{u} \in \mathbb{F}^M - \{\mathbf{0}\}} \frac{\mathbf{u}' \mathbf{X}_1 \mathbf{X}'_1 \mathbf{u}}{\|\mathbf{u}\|_2^2},$$

where we recognize $\mathbf{X}_1 \mathbf{X}'_1 / N \in \mathbb{F}^{M \times M}$ as an **empirical covariance matrix**.

Having found that \mathbf{u}_1 is the (unit norm) linear combination that maximizes the (empirical) variance, we can now remove the component along that direction and then seek another linear combination that maximizes the variance of what is left. To remove the component along \mathbf{u}_1 we construct a new matrix as follows:

$$\mathbf{X}_2 \triangleq \mathbf{P}_{\mathcal{R}(\mathbf{u}_1)}^\perp \mathbf{X}_1 = (\mathbf{I} - \mathbf{P}_{\mathcal{R}(\mathbf{u}_1)}) \mathbf{X}_1 = (\mathbf{I} - \mathbf{u}_1 \mathbf{u}'_1) \mathbf{X}_1.$$

Note that $\mathbf{X}_2 \mathbf{1} = \mathbf{0}$ so we still have a zero-mean array. Now we want to find

$$\arg \max_{\mathbf{u}: \|\mathbf{u}\|=1} \|(\mathbf{u}' \mathbf{X}_2)'\|_2 = \arg \max_{\mathbf{u}: \|\mathbf{u}\|=1} \|\mathbf{X}'_1 (\mathbf{I} - \mathbf{u}_1 \mathbf{u}'_1) \mathbf{u}\|_2 = \mathbf{u}_2.$$

To see why \mathbf{u}_2 is the solution here, use an SVD of \mathbf{X}_1 to write

$$\mathbf{X}_2 = (\mathbf{I} - \mathbf{u}_1 \mathbf{u}'_1) \mathbf{X}_1 = (\mathbf{I} - \mathbf{u}_1 \mathbf{u}'_1) \mathbf{U} \Sigma \mathbf{V}' = \mathbf{U} \operatorname{diag}\{0, \sigma_2, \sigma_3, \dots\} \mathbf{V}',$$

which is almost an SVD of \mathbf{X}_2 except for a permutation. Clearly the nonzero singular values of \mathbf{X}_2 are $(\sigma_2, \sigma_3, \dots, \sigma_r)$ and the spectral norm of \mathbf{X}_2 is σ_2 and the corresponding left singular vector, which is its principal left singular vector, is \mathbf{u}_2 .

One can continue this reasoning to show that \mathbf{U}_K , the first K columns of \mathbf{U} , provides a linear transform with orthonormal columns that maximizes the resulting variances of the **scores** $\mathbf{U}'_K \mathbf{X}_1$.

Using the same SVD of \mathbf{X}_1 , we have $\mathbf{U}'_K \mathbf{X}_1 = \mathbf{U}'_K \mathbf{U} \Sigma \mathbf{V}' = \Sigma_K \mathbf{V}'_K$, where \mathbf{V}_K denotes the first K columns of \mathbf{V} . Thus instead of storing all MN elements of \mathbf{X}_1 it suffices to store the MK and NK elements of \mathbf{U}_K and \mathbf{V}_K , reflecting dimensionality reduction.

6.7 Subspace learning

This section discusses an application of low-rank matrix approximation to subspace learning. In this setting, the low-dimensional subspace $\mathcal{R}(\mathbf{U}_K)$ (or $\mathcal{R}(\mathbf{V}_K)$) is more important to us than the overall approximation $\hat{\mathbf{A}}_K = \mathbf{U}_K \Sigma_K \mathbf{V}_K'$.

In supervised classification, we are given N labeled training data samples $\mathbf{x}_{j,1}, \dots, \mathbf{x}_{j,N}$ for each of J classes of objects, where each feature vector $\mathbf{x}_{j,n} \in \mathbb{F}^M$. We would like to learn something about the nature of each class so that later when we get a new sample vector \mathbf{x}_0 we can assess which class it is most like.

One basic approach to classification is the **K-nearest-neighbors method**. When $K = 1$, this method simply chooses the class of the nearest neighbor among the training samples:

$$\hat{j} = \text{[redacted]}$$

This basic approach requires JN comparisons, so complexity grows with training data size. Also, its performance may not improve as N increases when the marginal distributions overlap.

Instead, often we try to learn a model from the training data, and then use fit to the model as the basis for classification. Here we focus on learning a subspace model for each class. Then to classify a test point \mathbf{x}_0 we simply compute the distance of \mathbf{x}_0 to each of the J subspaces and see which is closest.

To simplify notation, we drop the class subscript j and focus on learning a subspace for a set of samples for one class type, $\mathbf{x}_1, \dots, \mathbf{x}_N$ where each $\mathbf{x}_n \in \mathbb{F}^M$.

Saying that all the \mathbf{x}_n vectors lie in a subspace of dimension K means that there is some $M \times K$ matrix \mathbf{Q} having orthonormal columns such that $\mathbf{x}_n \in \mathcal{R}(\mathbf{Q})$. The columns of \mathbf{Q} form an **orthonormal basis** for the subspace. In practice the data will only *approximately* lie in a subspace, so $\mathbf{x}_n \approx \mathbf{Q}\mathbf{z}_n$ for some coefficient vector $\mathbf{z}_n \in \mathbb{F}^K$. We want to *learn* the subspace basis matrix \mathbf{Q} from the training data $\mathbf{X} = [\mathbf{x}_1 \ \dots \ \mathbf{x}_N]$, a $M \times N$ matrix. Writing $\mathbf{x}_n \approx \mathbf{Q}\mathbf{z}_n$ in matrix form:

$$\underbrace{\mathbf{X}}_{M \times N} \approx$$

To find \mathbf{Q} and \mathbf{Z} we could pursue the following optimization problem:

$$\hat{\mathbf{Q}} =$$
(6.15)

In this setting $K < \min(M, N)$ so the product $\mathbf{Q}\mathbf{Z}$ is matrix with (at most) rank K .

Thus this problem is essentially a low-rank approximation problem except that here we really care only about the subspace basis \mathbf{Q} and not the coefficients \mathbf{Z} (nor their product). The low-rank solution is

$$\hat{\mathbf{X}}_K = \sum_{k=1}^K \sigma_k \mathbf{u}_k \mathbf{v}'_k =$$

So to learn the subspace \mathbf{Q} we simply use the first K singular vectors of the training data \mathbf{X} .

After learning a subspace Q_j for each of the J classes, we perform classification of a test point \mathbf{x}_0 by finding the nearest subspace:

$$\hat{j} =$$

$$f_j(\mathbf{x}_0) \triangleq$$

where $P_Q \mathbf{x}$ denotes the orthogonal projection of \mathbf{x} onto the subspace $\mathcal{R}(Q)$.

A Discussion section task sheet will explore this method with hand-written digits from the **MNIST** data.

Learning subspaces via (6.15) is optimal in the Frobenius norm sense of low-rank approximation, but is *not* necessarily optimal for the purposes of classification. An extension called **supervised PCA** aims to find subspaces that are good for both approximation and classification [26–28].

See [29] for a special issue on subspace learning and a generalization called submanifold learning.

Which form above is the most compute efficient?

A: (a)

B: (b)

C: (c)

??

For that efficient form, how many multiplies are needed per test sample?

A: $J2M(K + 1)$

B: $JM(M + 1)$

C: JM^2

D: $J2MK$

??

Subspace clustering

(Read)

The **subspace learning** approach described above, especially (6.15), assumes that the training data is **labeled** by class type, so it is **supervised learning**. There are also **unsupervised** subspace learning methods that perform **subspace clustering** as part of the subspace learning process [30, 31] [10].

Given unlabeled training data $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{F}^M$, in a J -subspace clustering approach we want to learn a set of J orthogonal bases $\mathbf{Q}_1, \dots, \mathbf{Q}_J \in \mathcal{V}_K(\mathbb{F}^M) \subset \mathbb{F}^{M \times K}$, such that each training point \mathbf{x}_n is close to one of the J subspaces $\{\mathcal{R}(\mathbf{Q}_j)\}$. One possible formulation is

$$\arg \min_{\mathbf{Q}_1, \dots, \mathbf{Q}_J \in \mathcal{V}_K(\mathbb{F}^M)} \min_{j_1, \dots, j_N \in \{1, \dots, J\}} \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{P}_{\mathbf{Q}_{j_n}} \mathbf{x}_n\|_2^2.$$

An alternative way of writing this is to use weights where w_{nj} indicates whether the n th training point is assigned to the j class:

$$\arg \min_{\mathbf{Q}_1, \dots, \mathbf{Q}_J \in \mathcal{V}_K(\mathbb{F}^M)} \min_{\mathbf{W} \in \mathcal{W}} \sum_{n=1}^N \sum_{j=1}^J w_{nj} \|\mathbf{x}_n - \mathbf{P}_{\mathbf{Q}_j} \mathbf{x}_n\|_2^2, \quad \mathcal{W} \triangleq \left\{ \mathbf{W} \in \mathbb{R}^{N \times J} : w_{nj} \in \{0, 1\}, \sum_{j=1}^J w_{nj} = 1 \right\}.$$

These are challenging discrete optimization problems called **integer programming** problems. One can alternate between updating the basis $\{\mathbf{Q}_j\}$ and the cluster assignments ($\{j_n\}$ or \mathbf{W}).

An alternative formulation called **sparse subspace clustering (SSC)** relaxes the discrete problem [30].

A related concept for unlabeled data is called **generalized principal component analysis (GPCA)** [32].

6.8 Summary

This chapter described several methods for computing a **low-rank approximation** to a matrix, all of which use a **SVD**. It also discussed methods for choosing the **rank** or **regularization parameter**. It applied the methods to sensor localization (**multidimensional scaling**) and to **subspace learning**, which is useful for **classification** problems.

There are two main perspectives in this chapter.

- We started by taking an arbitrary matrix \mathbf{A} and finding a low-rank approximation $\hat{\mathbf{A}}_K \approx \mathbf{A}$.

This perspective has applications in data compression and in dimensionality reduction.

In these applications, we want K to be small, but we also want the approximation error $\|\hat{\mathbf{A}}_K - \mathbf{A}\|_{\text{F}}$ to be small, and there is a trade-off between the two desires as we vary K .

- We then considered matrix **denoising** applications where we are given a noisy matrix \mathbf{Y} that we model as $\mathbf{Y} = \mathbf{X} + \boldsymbol{\varepsilon}$ where we think \mathbf{X} is a low-rank matrix. Then we process \mathbf{Y} to make an estimate $\hat{\mathbf{X}}$ of the latent matrix \mathbf{X} .

In this setting, we want the estimation error $\|\hat{\mathbf{X}} - \mathbf{X}\|_{\text{F}}$ to be small, and we do not care much about the approximation error $\|\hat{\mathbf{X}} - \mathbf{Y}\|_{\text{F}}$.

Bibliography

-
- [1] M. Udell and A. Townsend. “Why are big data matrices approximately low rank?” In: *SIAM J. Math. of Data Sci.* 1.1 (2019), 144–60 (cit. on p. 6.2).
 - [2] A. J. Laub. *Matrix analysis for scientists and engineers*. Soc. Indust. Appl. Math., 2005 (cit. on p. 6.2).

- [3] C. Eckart and G. Young. “The approximation of one matrix by another of lower rank”. In: *Psychometrika* 1.3 (1936), 211–8 (cit. on pp. 6.3, 6.38).
- [4] L. Mirsky. “Symmetric gauge functions and unitarily invariant norms”. In: *The quarterly journal of mathematics* 11.1 (Jan. 1960), 50–9 (cit. on pp. 6.6, 6.17, 6.39).
- [5] R. B. Cattell. “The scree test for the number of factors”. In: *Multivariate Behavioral Research* 1.2 (1966), 245–76 (cit. on p. 6.8).
- [6] E. Dobriban. *Factor selection by permutation*. 2017 (cit. on p. 6.11).
- [7] M. Buehrer, K. P. Pruessmann, P. Boesiger, and S. Kozerke. “Array compression for MRI with large coil arrays”. In: *Mag. Res. Med.* 57.6 (June 2007), 1131–9 (cit. on p. 6.12).
- [8] T. Zhang, J. M. Pauly, S. S. Vasanawala, and M. Lustig. “Coil compression for accelerated imaging with Cartesian sampling”. In: *Mag. Res. Med.* 69.2 (Feb. 2013), 571–82 (cit. on p. 6.12).
- [9] P. Drineas and I. C. F. Ipsen. “Low-rank matrix approximations do not need a singular value gap”. In: *SIAM J. Matrix. Anal. Appl.* 40.1 (Jan. 2019), 299–319 (cit. on p. 6.14).
- [10] Y-L. Yu and D. Schuurmans. “Rank/norm regularization with closed-form solutions: application to subspace clustering”. In: *Proc. 27th Conf. Uncertainty in AI*. 2011, 778–85 (cit. on pp. 6.17, 6.39, 6.62).
- [11] P. Verboon and W. J. Heiser. “Resistant lower rank approximation of matrices by iterative majorization”. In: *Comp. Stat. Data Anal.* 18.4 (Nov. 1994), 457–67 (cit. on p. 6.17).
- [12] E. J. Candes, X. Li, Y. Ma, and J. Wright. “Robust principal component analysis?” In: *J. Assoc. Comput. Mach.* 58.3 (May 2011), 1–37 (cit. on p. 6.17).
- [13] R. R. Nadakuditi. “OptShrink: an algorithm for improved low-rank signal matrix denoising by optimal, data-driven singular value shrinkage”. In: *IEEE Trans. Info. Theory* 60.5 (May 2014), 3002–18 (cit. on pp. 6.22, 6.50).
- [14] Y. Shitov. “The nonnegative rank of a matrix: hard problems, easy solutions”. In: *SIAM Review* 59.4 (2017), 794–800 (cit. on p. 6.22).
- [15] M. B. Cohen, J. Nelson, and D. P. Woodruff. *Optimal approximate matrix product in terms of stable rank*. 2016 (cit. on p. 6.23).
- [16] I. Borg and P. J. F. Groenen. *Modern multidimensional scaling: Theory and applications*. Springer, 2005 (cit. on p. 6.29).
- [17] J. J. Moreau. “Proximité et dualité dans un espace hilbertien”. In: *Bulletin de la Société Mathématique de France* 93 (1965), 273–99 (cit. on p. 6.34).
- [18] N. Parikh and S. Boyd. “Proximal algorithms”. In: *Found. Trends in Optimization* 1.3 (2013), 123–231 (cit. on p. 6.34).
- [19] A. Beck. *First-order methods in optimization*. Soc. Indust. Appl. Math., 2017 (cit. on p. 6.34).

- [20] E. J. Candes, C. A. Sing-Long, and J. D. Trzasko. “Unbiased risk estimates for singular value thresholding and spectral estimators”. In: *IEEE Trans. Sig. Proc.* 61.19 (Oct. 2013), 4643–57 (cit. on pp. [6.38](#), [6.48](#)).
- [21] E. J. Candes and T. Tao. “The power of convex relaxation: near-optimal matrix completion”. In: *IEEE Trans. Info. Theory* 56.5 (May 2010), 2053–80 (cit. on p. [6.41](#)).
- [22] G. Golub, A. Hoffman, and G. Stewart. “A generalization of the Eckart-Young-Mirsky matrix approximation theorem”. In: *Linear Algebra and its Applications* 88 (Apr. 1987), 317–27 (cit. on p. [6.44](#)).
- [23] S. Gu, L. Zhang, W. Zuo, and X. Feng. “Weighted nuclear norm minimization with application to image denoising”. In: *Proc. IEEE Conf. on Comp. Vision and Pattern Recognition*. 2014, 2862–9 (cit. on p. [6.44](#)).
- [24] Y. Xie, S. Gu, Y. Liu, W. Zuo, W. Zhang, and L. Zhang. “Weighted Schatten p-norm minimization for image denoising and background subtraction”. In: *IEEE Trans. Im. Proc.* 25.10 (Oct. 2016), 4842–57 (cit. on p. [6.44](#)).
- [25] X. Liu, X-Y. Jing, G. Tang, F. Wu, and Q. Ge. “Image denoising using weighted nuclear norm minimization with multiple strategies”. In: *Signal Processing* 135 (June 2017), 239–52 (cit. on p. [6.44](#)).
- [26] E. Bair, T. Hastie, D. Paul, and R. Tibshirani. “Prediction by supervised principal components”. In: *J. Am. Stat. Assoc.* 101.473 (2006), 119–37 (cit. on p. [6.61](#)).
- [27] E. Barshan, A. Ghodsi, Z. Azimifar, and M. Z. Jahromi. “Supervised principal component analysis: Visualization, classification and regression on subspaces and submanifolds”. In: *Pattern Recognition* 44.7 (July 2011), 1357–71 (cit. on p. [6.61](#)).
- [28] A. Ritchie, C. Scott, L. Balzano, D. Kessler, and C. S. Sripada. “Supervised principal component analysis via manifold optimization”. In: *IEEE Data Science Workshop (DSW)*. 2019, 6–10 (cit. on p. [6.61](#)).
- [29] Y. Ma, P. Niyogi, G. Sapiro, and R. Vidal. “Dimensionality reduction via subspace and submanifold learning [From the guest editors]”. In: *IEEE Sig. Proc. Mag.* 28.2 (Mar. 2011), 14–126 (cit. on p. [6.61](#)).
- [30] R. Vidal. “Subspace clustering”. In: *IEEE Sig. Proc. Mag.* 28.2 (Mar. 2011), 52–68 (cit. on p. [6.62](#)).
- [31] J. Lipor, D. Hong, D. Zhang, and L. Balzano. *Subspace clustering using ensembles of K-subspaces*. 2017 (cit. on p. [6.62](#)).
- [32] R. Vidal, Y. Ma, and S. Sastry. “Generalized principal component analysis (GPCA)”. In: *IEEE Trans. Patt. Anal. Mach. Int.* 27.12 (Dec. 2005), 1945–59 (cit. on p. [6.62](#)).