

# Chapter 4

## Linear equations and least-squares

### Contents (class version)

---

<b>4.0 Introduction to linear equations</b>	<b>4.2</b>
Linear regression and machine learning	4.4
<b>4.1 Linear least-squares estimation</b>	<b>4.6</b>
Minimization and gradients	4.10
Solving LLS using the normal equations	4.15
Solving LLS problems using the compact SVD	4.17
Uniqueness of LLS solution	4.22
<b>4.2 Moore-Penrose pseudoinverse</b>	<b>4.24</b>
<b>4.3 Linear least-squares estimation: Under-determined case</b>	<b>4.31</b>
Orthogonality principle	4.33
Minimum-norm LS solution via pseudo-inverse	4.36
<b>4.4 Truncated SVD solution</b>	<b>4.40</b>
Low-rank approximation interpretation of truncated SVD solution	4.43
Noise effects and perturbations	4.44
Tikhonov regularization aka ridge regression	4.46

4.5 Summary of LLS solution methods in terms of SVD . . . . .	4.49
4.6 Frames and tight frames . . . . .	4.50
4.7 Projection and orthogonal projection . . . . .	4.56
Projection onto a subspace . . . . .	4.62
Binary classifier design using least-squares . . . . .	4.70
4.8 Summary . . . . .	4.72

4.0 Introduction to linear equations

Source material for this chapter includes [1, §6.1, 8.1–8.4, 4.1, 5.2].

L§6.1

Solving a **system of linear equations** arises in numerous applications. Mathematically, a system of  $M$  equations in  $N$  unknowns is usually written

$$\mathbf{A}\mathbf{x} = \mathbf{y}, \quad \mathbf{A} \in \mathbb{F}^{M \times N}, \quad \mathbf{x} \in \mathbb{F}^N, \quad \mathbf{y} \in \mathbb{F}^M. \tag{4.1}$$

Typically  $\mathbf{A}$  is known from some type of (linear) modeling,  $\mathbf{y}$  corresponds to some type of measurements, and the goal is to solve for  $\mathbf{x}$ .

However, despite the equals sign in the classic expression (4.1), in practice often there does not *exist* any  $\mathbf{x}$  that yields equality, particularly when  $\mathbf{A}$  is tall, so the notation “ $\mathbf{A}\mathbf{x} \approx \mathbf{y}$ ” would be more realistic.

Sometimes there is not a *unique*  $\mathbf{x}$  that satisfies (4.1), in particular whenever  $\mathbf{A}$  is wide. Linear algebra texts focus on examining when a solution  $\mathbf{x}$  exists and is unique for (4.1).

**Solving  $\mathbf{Ax} = \mathbf{y}$**  \_\_\_\_\_ (Read)

[1, Theorem 6.1] reviews the principal existence and uniqueness results for (4.1).

1. There *exists* a solution  $\mathbf{x}$  iff  $\mathbf{y} \in \mathcal{R}(\mathbf{A})$ .
2. There *exists* a solution  $\mathbf{x}$  for all  $\mathbf{y} \in \mathbb{R}^M$  iff  $\mathcal{R}(\mathbf{A}) = \mathbb{F}^M$ .
3. A solution  $\mathbf{x}$  is *unique* iff  $\mathcal{N}(\mathbf{A}) = \{\mathbf{0}\}$ .
4. There exists a unique solution for all  $\mathbf{y} \in \mathbb{R}^M$  iff  $\mathbf{A}$  is  $M \times M$  and non-singular (invertible), *i.e.*, none of the eigenvalues or singular values of  $\mathbf{A}$  are zero.
5. There is at most one solution for all  $\mathbf{y} \in \mathbb{R}^M$  iff  $\mathbf{A}$  has linearly independent columns, *i.e.*,  $\mathcal{N}(\mathbf{A}) = \{\mathbf{0}\}$ , and this is possible only if  $M \geq N$ .
6. The homogeneous system  $\mathbf{Ax} = \mathbf{0}$  has a nontrivial (*i.e.*, nonzero) solution iff  $\text{rank}(\mathbf{A}) < N$ .

(You should read and verify these points.)

To understand #3, suppose  $\mathbf{x}_0 \in \mathcal{N}(\mathbf{A})$  and  $\mathbf{x}_0 \neq \mathbf{0}$  and suppose we have a solution  $\mathbf{Ax}_1 = \mathbf{y}$ . Then  $\mathbf{x}_2 = \mathbf{x}_1 + \mathbf{x}_0$  is also a solution because  $\mathbf{Ax}_2 = \mathbf{Ax}_1 + \mathbf{Ax}_0 = \mathbf{y} + \mathbf{0}$ .

When  $\mathbf{A}$  is wide ( $M < N$ ),  $\mathcal{N}(\mathbf{A})$  is nontrivial and there are (infinitely) many solutions to (4.1). One must design some way to choose among all those solutions.

When  $\mathbf{A}$  is square ( $M = N$ ) and full rank (and hence invertible), then there is a unique solution  $\mathbf{x} = \mathbf{A}^{-1}\mathbf{y}$ . We say that “the number of equations” equals “the number of unknowns.”

---

## Linear regression and machine learning

In **linear regression**, we are given a set of training data consisting of  $M$  input (feature) vectors  $\mathbf{a}_1, \dots, \mathbf{a}_M \in \mathbb{F}^N$  and corresponding responses  $\mathbf{y}_1, \dots, \mathbf{y}_M \in \mathbb{F}$  and we want to find coefficients / weights  $\mathbf{x} \in \mathbb{F}^N$  such that  $\mathbf{a}_m^T \mathbf{x} \approx y_m$ . Stacking up the input vectors into a matrix  $\mathbf{A}$  and the response variables into a vector  $\mathbf{y}$  yields the following matrix-vector form:

$$\begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} \approx \begin{bmatrix} \mathbf{a}_1^T \\ \vdots \\ \mathbf{a}_M^T \end{bmatrix} \mathbf{x}, \quad \text{i.e.,} \quad \mathbf{y} = \mathbf{A}\mathbf{x} + \boldsymbol{\varepsilon},$$

where  $\boldsymbol{\varepsilon} \in \mathbb{F}^M$  denotes noise or other errors.

We want to determine (“learn”) the weight vector, or coefficients,  $\mathbf{x}$  so that given some future feature vector  $\mathbf{a}$ , we can **predict** the corresponding (unknown) response by simply computing  $\mathbf{a}^T \mathbf{x}$ .

This is a classical problem in statistics and it is key “machine learning” method that one usually should consider (for estimation/regression problems) before considering more complicated methods.

The notation for linear regression varies by discipline:

in statistics:  $\mathbf{y} \approx \mathbf{X}\boldsymbol{\beta}$ ; in machine learning:  $\mathbf{y} \approx \mathbf{X}\mathbf{w}$ ; linear algebra:  $\mathbf{y} \approx \mathbf{A}\mathbf{x}$ .

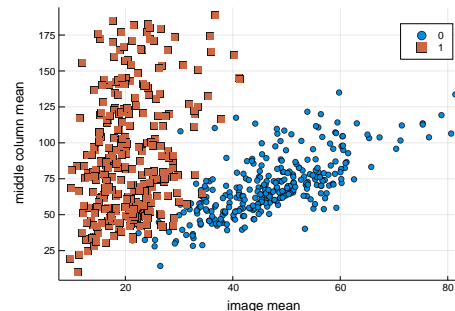
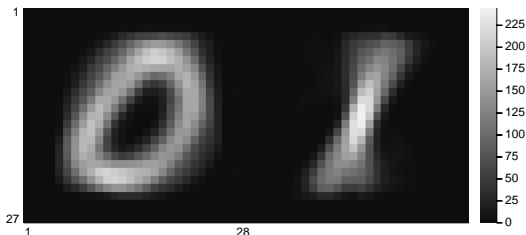
The variables here have many different names:

- $\mathbf{y}$ : response, labels, regressand, endogenous variable, measured variable, criterion variable, dependent variable, predicted variable, ...
- $\mathbf{a}$  (rows of  $\mathbf{A}$ ) : features, regressors, exogenous variables, explanatory variables, covariates, input variables, predictor variables, independent variables, ...
- $\mathbf{x}$ : parameter vector, regression coefficients, unknown, ...
- $\varepsilon$ : noise, disturbance, error, ...

Example. Predict child's height from parent's height. [\[wiki\]](#)

Example. Hand-written digit classification using two “hand-crafted” features (plus a bias or offset).

Here  $\mathbf{A}$  is  $3 \times (28 \cdot 27) = 3 \times 756$ .



### 4.1 Linear least-squares estimation

L§8.1

In a typical situation where  $M > N$ , called “tall” or “over-determined,” one can show that  $\mathcal{R}(\mathbf{A}) \neq \mathbb{F}^M$ . So there will be (infinitely) many  $\mathbf{y}$  (those not in  $\mathcal{R}(\mathbf{A})$ ) for which no solution  $\mathbf{x}$  exists.

Thus, when  $M > N$ , instead of insisting on *exactly* solving  $\mathbf{A}\mathbf{x} = \mathbf{y}$ , usually we look for *approximate* solutions where  $\mathbf{A}\mathbf{x} \approx \mathbf{y}$ . To do this, we must quantify “approximately equal.”

The most important and common approximate solution is to use **linear least-squares (LLS)** estimation or fitting, where we find an estimate  $\hat{\mathbf{x}}$  that “best fits” the data  $\mathbf{y}$  using a Euclidean norm distance as follows:

(4.2)

- $\hat{\mathbf{x}}$  is called the **linear least-squares estimate** (or solution)  
(The “hat” or caret above  $\mathbf{x}$  is often used to denote an estimate.)
- $\mathbf{A}\hat{\mathbf{x}} - \mathbf{y}$  is called the **residual**
- $\|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2$  is the Euclidean norm of the residuals for a candidate solution  $\mathbf{x}$  and is a measure of the “error” of the fit. It is sometimes called the “goodness of fit” even though a larger value means a worse fit!
- $\arg \min_{\mathbf{x}}$  means that we see the argument  $\hat{\mathbf{x}}$  that minimizes the fitting error; the minimum value of the fitting error:  $\min_{\mathbf{x} \in \mathbb{F}^N} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2$  usually is of less interest.  
In words we often say “we minimized the squared error” but we really mean “we found a solution  $\hat{\mathbf{x}}$  that

minimized the squared error.”

Although this technique is ancient, it remains one of the most important linear algebra methods for signal processing and data analysis in general.

---

Example. Suppose we observe noisy samples of signal:

L§8.3

$$y_m = s(t_m) + \epsilon_m, \quad m = 1, \dots, M$$

and we believe that the signal is a cubic polynomial:

$$s(t) = \alpha_0 + \alpha_1 t + \alpha_2 t^2 + \alpha_3 t^3$$

with unknown coefficients. In matrix-vector form:

$$\mathbf{y} \approx \mathbf{A}\mathbf{x}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix}, \quad \mathbf{A} = \begin{bmatrix} 1 & t_1 & t_1^2 & t_1^3 \\ 1 & t_2 & t_2^2 & t_2^3 \\ \vdots & & & \vdots \\ 1 & t_M & t_M^2 & t_M^3 \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} \alpha_0 \\ \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{bmatrix}.$$

A figure on p. 4.9 illustrates this problem. Here the “features” are the time points  $t_m$  raised to different powers, and we call it **polynomial regression**.

To find the coefficient vector  $\mathbf{x}$ , one option would be to take just  $M = 4$  samples. As long as we pick 4 distinct  $t_m$  points then one can show (using the fact that monomials are linearly independent functions) that the  $4 \times 4$  matrix  $\mathbf{A}_4$  has **linearly independent** columns. Thus  $\mathbf{A}_4$  is **invertible** and we could use  $\hat{\mathbf{x}} = \mathbf{A}_4^{-1} \mathbf{y}$  as an estimate of the coefficients.

However, in the presence of noise in the data, this approach would give very noisy and unreliable estimates of the coefficients.

Instead, it is preferable to use all  $M \gg 4$  samples and estimate  $\hat{\mathbf{x}}$  using linear least-squares.

For higher polynomial orders, it is more stable to use **orthogonal polynomials** as the basis, instead of monomials. We use monomials here for simplicity in this example.

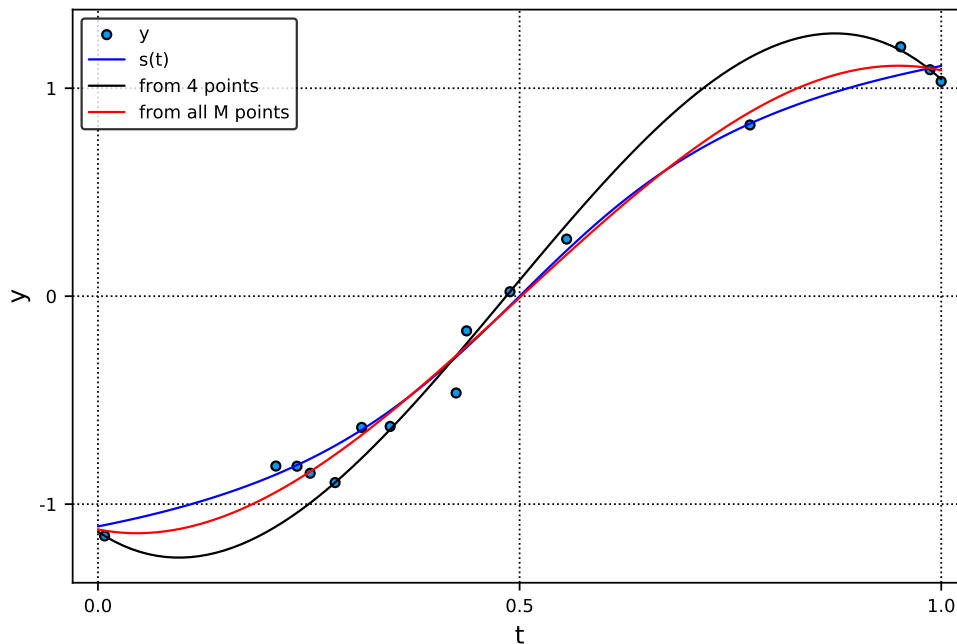
For a demo, see

[https://web.eecs.umich.edu/~fessler/course/551/julia/demo/04\\_ls\\_fit1.html](https://web.eecs.umich.edu/~fessler/course/551/julia/demo/04_ls_fit1.html)

[https://web.eecs.umich.edu/~fessler/course/551/julia/demo/04\\_ls\\_fit1.ipynb](https://web.eecs.umich.edu/~fessler/course/551/julia/demo/04_ls_fit1.ipynb)



The following curves illustrate that fitting a cubic polynomial using all  $M \gg 4$  points via (4.2) is better than using just 4 points with  $A_4^{-1}$ . (One can prove that statement statistically; see EECS 564.)



A key line in the demo code is `xh = A \ y` and next we explain in detail the very important mathematical foundation behind that computation.

## Minimization and gradients

The cost function in (4.2) is differentiable, so we use derivatives to find a minimizer.

### Review of 1D minimization by example

(Read)

To find the minimizer of a function  $f : \mathbb{R} \mapsto \mathbb{R}$  such as

$$f(x) = x(x - 1)^3,$$

you simply take the **derivative** and set it equal to zero:

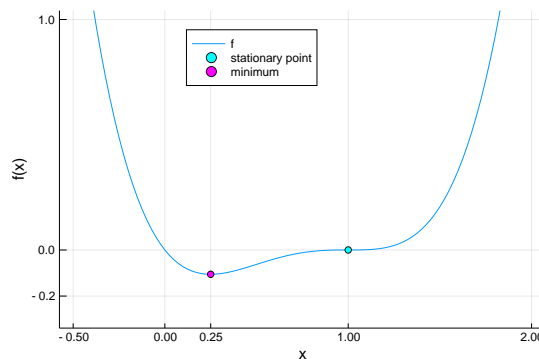
$$0 = \dot{f}(x) = (x - 1)^2(4x - 1).$$

This case has two roots at  $x = 1$  and  $x = 1/4$ . Because  $f(1) = 0$  and  $f(1/4) < 0$  it looks like  $x = 1/4$  is the minimizer.

To be sure that  $x = 1/4$  is a minimizer, we also take the second derivative:

$$\ddot{f}(x) = 6(2x^2 - 3x + 1) \implies \ddot{f}(1/4) = 9/4 > 0,$$

so  $x = 1/4$  is at least a local minimizer. Checking  $\pm\infty$  verifies. This example required extra work because  $f$  is non-convex. For convex differentiable functions, simply finding a point  $x$  where  $\dot{f}(x) = 0$  is sufficient.



## Calculating gradients

---

Similarly, to solve the minimization problem (4.2), we will find a zero of its gradient.

Consider the **affine** function  $f_1 : \mathbb{R}^N \mapsto \mathbb{R}$  defined for a vector  $\mathbf{v} \in \mathbb{R}^N$  and  $b \in \mathbb{R}$  by

$$f_1(\mathbf{x}) \triangleq b + \mathbf{v}'\mathbf{x} = b + \sum_{n=1}^N v_n x_n.$$

The (column) **gradient** of this function is the vector of partial derivatives:

$$\nabla f_1(\mathbf{x}) =$$

---

(Read)

Consider the **quadratic** function  $f_2 : \mathbb{R}^N \mapsto \mathbb{R}$  defined for a matrix  $\mathbf{M} \in \mathbb{R}^{N \times N}$  by

$$f_2(\mathbf{x}) \triangleq \frac{1}{2} \mathbf{x}' \mathbf{M} \mathbf{x} = \frac{1}{2} \sum_{k=1}^N \sum_{l=1}^N M_{kl} x_k x_l.$$

The partial derivative of this function w.r.t.  $x_n$  is

$$\frac{\partial}{\partial x_n} f_2(\mathbf{x}) = \frac{\partial}{\partial x_n} \frac{1}{2} \sum_{k=1}^N \sum_{l=1}^N M_{kl} x_k x_l.$$

I deliberately chose different letters  $k, l$  for the summations different from  $n$  to avoid mishaps.

Using linearity of differentiation:

$$\begin{aligned}
 \frac{\partial}{\partial x_n} f_2(\mathbf{x}) &= \frac{1}{2} \sum_{k=1}^N \sum_{l=1}^N M_{kl} \left( \frac{\partial}{\partial x_n} x_k x_l \right) \\
 &= \frac{1}{2} \sum_{k=1}^N \sum_{l=1}^N M_{kl} (x_l \mathbb{I}_{\{n=k \neq l\}} + x_k \mathbb{I}_{\{n=l \neq k\}} + 2x_n \mathbb{I}_{\{n=k=l\}}) \\
 &= \frac{1}{2} \left( \sum_{k=1}^N \sum_{l=1}^N M_{kl} x_l \mathbb{I}_{\{n=k \neq l\}} + \sum_{k=1}^N \sum_{l=1}^N M_{kl} x_k \mathbb{I}_{\{n=l \neq k\}} + 2 \sum_{k=1}^N \sum_{l=1}^N M_{kl} x_n \mathbb{I}_{\{n=k=l\}} \right) \\
 &= \frac{1}{2} \left( \sum_{l=1}^N M_{nl} x_l \mathbb{I}_{\{n \neq l\}} + \sum_{k=1}^N M_{kn} x_k \mathbb{I}_{\{n \neq k\}} + 2M_{nn} x_n \right) \\
 &= \frac{1}{2} \sum_{l=1}^N M_{nl} x_l + \frac{1}{2} \sum_{k=1}^N M_{kn} x_k \\
 &= \frac{1}{2} [\mathbf{M}\mathbf{x}]_n + \frac{1}{2} [\mathbf{M}'\mathbf{x}]_n,
 \end{aligned}$$

where  $[\mathbf{A}\mathbf{x}]_i$  denotes the  $i$ th element of the vector  $\mathbf{A}\mathbf{x}$ .

Now write the **gradient** of this function in column form by stacking up all the partial derivatives into a vector:

$$\nabla f_2(\mathbf{x}) = \begin{bmatrix} \frac{\partial}{\partial x_1} f_2 \\ \vdots \\ \frac{\partial}{\partial x_N} f_2 \end{bmatrix} = \frac{1}{2} \begin{bmatrix} [\mathbf{M}\mathbf{x}]_1 + [\mathbf{M}'\mathbf{x}]_1 \\ \vdots \\ [\mathbf{M}\mathbf{x}]_N + [\mathbf{M}'\mathbf{x}]_N \end{bmatrix} = \frac{1}{2} (\mathbf{M}\mathbf{x} + \mathbf{M}'\mathbf{x}).$$

**Summary of key gradients needed**

- constant:  $\nabla c = \mathbf{0}$
- linear:  $\nabla \mathbf{v}'\mathbf{x} = \mathbf{v}$
- quadratic:  $\nabla \frac{1}{2} \mathbf{x}' \mathbf{M} \mathbf{x} = \frac{1}{2} (\mathbf{M} + \mathbf{M}') \mathbf{x}$  and if  $\mathbf{M}$  is symmetric then  $\nabla \frac{1}{2} \mathbf{x}' \mathbf{M} \mathbf{x} = \mathbf{M} \mathbf{x}$ .

**The LLS problem**

Now consider the LLS cost function corresponding to (4.2) for the real-valued case with  $f : \mathbb{R}^N \mapsto \mathbb{R}$ :

$$\text{[Yellow box containing equation (4.3)]} \tag{4.3}$$

Applying the linearity of differentiation and the two gradients derived above, the (column) **gradient** of the cost function  $f$  w.r.t.  $\mathbf{x}$  is:

$$\nabla f(\mathbf{x}) = \text{[Yellow box]} \tag{4.4}$$

**Complex case** (Read)

The above derivation is applicable only when  $\mathbf{x}$  and  $\mathbf{y}$  and  $\mathbf{A}$  are all real-valued. When they are complex-valued, then strictly speaking the LLS cost function  $f$  is not differentiable! Instead of working with the **gradient**, in the complex case we work with the **conjugate cogradient**, using methods from **Wirtinger calculus** [2] [3] [4]. One intermediate step above is different:



$$f(\mathbf{x}) = \frac{1}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2 = \frac{1}{2} (\mathbf{x}' \mathbf{A}' \mathbf{Ax} - 2 \operatorname{real}\{\mathbf{y}' \mathbf{Ax}\} + \mathbf{y}' \mathbf{y}),$$

but conveniently the final expression for the **conjugate cogradient** is identical to (4.4):

$$\nabla f(\mathbf{x}) = \mathbf{A}'(\mathbf{Ax} - \mathbf{y}).$$

The only difference is that here  $\mathbf{A}'$  truly denotes a **Hermitian transpose**, whereas in (4.4) it is simply an ordinary matrix transpose.

We use (4.4) hereafter, knowing that it applies to both real and complex-valued problems.

And we will call it the “**gradient**” even though technically it is the “**conjugate cogradient**” in the complex case.

## Solving LLS using the normal equations

L§8.2

One can show that the LLS cost function corresponding to (4.2), *i.e.*,

$$f(\mathbf{x}) = \frac{1}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2,$$

where  $f : \mathbb{F}^N \mapsto \mathbb{R}$ , is a **convex function**. Because  $f$  is convex and differentiable, a vector  $\hat{\mathbf{x}}$  is a minimizer of  $f$  if and only if  $\hat{\mathbf{x}}$  is a point where the **gradient** of  $f$  is zero. We have shown that the gradient of  $f$  w.r.t.  $\mathbf{x}$  is

$$\nabla f(\mathbf{x}) = \mathbf{A}'(\mathbf{Ax} - \mathbf{y}).$$

Setting this gradient to zero, *i.e.*,  $\nabla f(\mathbf{x})|_{\mathbf{x}=\hat{\mathbf{x}}} = \mathbf{0}$ , and rearranging yields the **normal equations**:

(4.5)

The set of minimizers of  $f$  is the set of points  $\hat{\mathbf{x}}$  satisfying (4.5).

Note that we started with a  $M \times N$  matrix  $\mathbf{A}$  in (4.1) and (4.2), but the normal equations *always* have a square  $N \times N$  **Gram matrix**.

- If  $\mathbf{A}$  is a wide matrix ( $M < N$ ), then the rank of the  $N \times N$  Gram matrix  $\mathbf{A}'\mathbf{A}$  is at most  $M$ , so  $\mathbf{A}'\mathbf{A}$  is singular. We will use the SVD of  $\mathbf{A}$  to analyze this case **under-determined** case on p. 4.31.
- If  $\mathbf{A}$  is tall or square *and* has linearly independent columns, *i.e.*, has rank  $N$ , then  $\mathbf{A}'\mathbf{A}$  is invertible.

## Practical LLS solution using backslash

When  $\mathbf{A}$  has full column rank, the Gram matrix  $\mathbf{A}'\mathbf{A}$  is invertible and the LLS cost function has a unique minimizer given by the solution to the normal equations (4.5). Specifically:

$$\mathbf{A} \in \mathbb{F}^{M \times N} \text{ and } \underbrace{\text{rank}(\mathbf{A}) = N}_{\Rightarrow M \geq N} \implies \hat{\mathbf{x}} = \arg \min_{\mathbf{x} \in \mathbb{F}^N} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 = \text{ } \quad (4.6)$$

This expression is useful for analysis, but it is not the best computational approach!

1. It might be tempting to implement (4.6) in JULIA using one of the two following code snippets.

Practical question: which of the following is likely to require less computation?

A: `xh = inv(A' * A) * (A' * y)`

B: `xh = inv(A' * A) * A' * y`

C: Both always use the same computation.

Using `inv(A' * A)` is computationally inefficient. Do not use it (at least not for large problem sizes)!

(An exception is applications where we must solve LLS problems for many  $\mathbf{y}$  with the same  $\mathbf{A}$  (as in HW).

A more efficient approach uses the **QR decomposition** of the matrix  $\mathbf{A}$ , without ever forming the Gram matrix  $\mathbf{A}'\mathbf{A}$ . In JULIA code it is simply: `xh = A \ y` or `xh = \ (A, y)`, spoken as “ $\mathbf{A}$  backslash  $\mathbf{y}$ .” It is also more efficient than applying **Gaussian elimination** to (4.5).

Next we turn to the SVD for more insight.



## Solving LLS problems using the compact SVD

L§8.4

The **compact SVD** provides an insightful way to analyze the LLS problem (4.2). It is also a reasonable way to solve LLS problems where  $M$  and  $N$  are not too large. (Large problems require iterative methods, discussed later.) Using a compact SVD  $\mathbf{A} = \mathbf{U}_r \Sigma_r \mathbf{V}_r'$  from (3.12), the LLS cost function (4.3) becomes:

$$\begin{aligned} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 &= \\ &= \\ &= \\ &= \end{aligned}$$

where we **completed the square** and used the fact that  $\|\mathbf{U}_r \mathbf{z}_r\|_2 = \|\mathbf{z}_r\|_2$  because  $\mathbf{U}$  is unitary.

Term 2 is independent of  $\mathbf{x}$ , so to minimize the LLS cost function we must minimize Term 1, *i.e.*,

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\Sigma_r \mathbf{V}_r' \mathbf{x} - \mathbf{U}_r' \mathbf{y}\|_2^2. \quad (4.7)$$

If  $r = 0$  (*i.e.*, if  $\mathbf{A} = \mathbf{0}$ ) then Term 1 vanishes, so we focus on the usual case where  $1 \leq r \leq N$  hereafter.

Although (4.7) initially might look more complicated than the original LLS problem, it is actually simpler because  $\Sigma_r$  is invertible and  $V_r$  has orthonormal columns. By inspection (without taking any gradients!) one possible solution to (4.7) is:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\Sigma_r V_r' \mathbf{x} - U_r' \mathbf{y}\|_2^2 = \quad (4.8)$$

because this solution makes Term 1 identically zero.

$$\text{Verification: } \Sigma_r V_r' \hat{\mathbf{x}} = \Sigma_r \underbrace{V_r' (V_r \Sigma_r^{-1} U_r' \mathbf{y})}_I = \underbrace{\Sigma_r \Sigma_r^{-1}}_I U_r' \mathbf{y} = U_r' \mathbf{y}.$$

Recall that  $\Sigma_r$  is  $r \times r$  and contains the  $r$  *nonzero* singular values of  $\mathbf{A}$  along its diagonal, so it is invertible.

If  $r = N$ , i.e., if  $\mathbf{A}$  has full column rank, then  $V_r = V$  and the solution (4.8) is the *unique* minimizer.

However, if  $r < N$ , then there are multiple minimizers. All minimizers are given by

$$\hat{\mathbf{x}} = \quad (4.9)$$

where  $V = [V_r \ V_0]$  and  $z_0$  is *any* vector of the appropriate length!

Any such solution makes Term 1 identically zero because:

$$\Sigma_r V_r' \hat{\mathbf{x}} = \Sigma_r V_r' (V_r \Sigma_r^{-1} U_r' \mathbf{y} + V_0 z_0) = U_r' \mathbf{y} + \underbrace{\Sigma_r V_r' V_0}_{\mathbf{0}} z_0 = U_r' \mathbf{y}.$$

2. When  $\mathbf{A}$  is  $M \times N$ , what is the length of the vector  $z_0$  in (4.9)?

A:  $r$

B:  $N - r$

C:  $M - r$

D:  $N$

E:  $M$

??

Letting  $\mathbf{x} = \mathbf{V}\mathbf{z} = [\mathbf{V}_r \quad \mathbf{V}_0] \begin{bmatrix} \mathbf{z}_r \\ \mathbf{z}_0 \end{bmatrix}$ , so  $\mathbf{z}$  denotes the coordinates of  $\mathbf{x}$  in the  $\mathbf{V}$  coordinate system, another way of writing the general solution in (4.9) is

$$\hat{\mathbf{x}} = [\mathbf{V}_r \quad \mathbf{V}_0] \begin{bmatrix} \Sigma_r^{-1} \mathbf{U}_r' \mathbf{y} \\ \mathbf{z}_0 \end{bmatrix} = \mathbf{V} \hat{\mathbf{z}}, \quad \hat{\mathbf{z}} = \begin{bmatrix} \Sigma_r^{-1} \mathbf{U}_r' \mathbf{y} \\ \mathbf{z}_0 \end{bmatrix} = \begin{bmatrix} [\mathbf{U}_r' \mathbf{y}]_1 / \sigma_1 \\ \vdots \\ [\mathbf{U}_r' \mathbf{y}]_r / \sigma_r \\ \mathbf{z}_0 \end{bmatrix}. \quad (4.10)$$

This is a general expression for “the” LLS solution in terms of a compact SVD of  $\mathbf{A}$  and the data  $\mathbf{y}$ .

The arbitrary choice of  $\mathbf{z}_0$  may seem unsettling. We will address that concern shortly.

There are two ways that  $\mathbf{A}$  can have rank  $r < N$ , leading to a non-unique solution:

- If  $\mathbf{A}$  is square or tall but has
- If  $\mathbf{A}$  is wide, because then

When  $\mathbf{A}$  is tall with full rank  $r = N$ , the (unique) LLS solution is

(Read)

Example. Return to the example of fitting a cubic polynomial on p. 4.7. In this case  $M \gg N = 4$ . As long as at least 4 of the  $\{t_m\}$  values are distinct, the fact that monomials are linearly independent functions implies that  $\mathbf{A}$  has full rank, *i.e.*,  $r = N = 4$ . In this (typical) case, the SVD simplifies to

$$\underbrace{\mathbf{A}}_{M \times 4} = \underbrace{\mathbf{U}}_{M \times M} \underbrace{\mathbf{\Sigma}}_{M \times 4} \underbrace{\mathbf{V}'}_{4 \times 4} = \underbrace{\mathbf{U}_4}_{M \times 4} \underbrace{\mathbf{\Sigma}_4}_{4 \times 4} \underbrace{\mathbf{V}'_4}_{4 \times 4}, \quad \mathbf{\Sigma} = \begin{bmatrix} \mathbf{\Sigma}_4 \\ \mathbf{0}_{(M-4) \times 4} \end{bmatrix} = \begin{bmatrix} \sigma_1 & 0 & 0 & 0 \\ 0 & \sigma_2 & 0 & 0 \\ 0 & 0 & \sigma_3 & 0 \\ 0 & 0 & 0 & \sigma_4 \\ \hline \mathbf{0}_{(M-4) \times 4} \end{bmatrix}.$$

The solution  $\hat{\mathbf{x}}$  in (4.9) simplifies to

$$\hat{\mathbf{x}} = \mathbf{V} \mathbf{\Sigma}_4^{-1} \mathbf{U}'_4 \mathbf{y}. \quad (4.11)$$

Because  $r = N = 4$  here, this is a final and unique LLS solution; there is no arbitrary vector  $\mathbf{z}_0$  to select.

## Practical implementation

(Read)

Returning to the JULIA demo for this example, where  $\mathbf{A}$  is tall and  $r = N$ , we can implement this SVD-based solution using the following code.

Start with an SVD:

```
U, s, V = svd(A)
```

Recall that for this **economy SVD** version:

- $\mathbf{U}$  is  $M \times N$  instead of the usual  $M \times M$ , i.e., is  $\mathbf{U}_r = \mathbf{U}_N$  in (3.12)
- $\mathbf{s}$  is a vector of the  $(\sigma_1, \dots, \sigma_N)$  values, so  $\mathbf{\Sigma}_r = \mathbf{\Sigma}_N = \text{Diagonal}(\mathbf{s})$  because  $r = N$  here.

Using this economy SVD, two mathematically equivalent forms of the LS solution are:

- $\mathbf{xh} = \mathbf{V} * \text{Diagonal}(1 ./ \mathbf{s}) * (\mathbf{U}' * \mathbf{y})$  which looks like  $\hat{\mathbf{x}} = \mathbf{V}\mathbf{\Sigma}_4^{-1}\mathbf{U}_4'\mathbf{y}$  in (4.11)
- $\mathbf{xh} = \mathbf{V} * ((1 ./ \mathbf{s}) .* (\mathbf{U}' * \mathbf{y}))$  which looks like (4.10).

Why the parentheses  $(\mathbf{U}' * \mathbf{y})$  ? Saves computation because  $\mathbf{U}'_N$  is wide.

When  $\mathbf{A}$  is tall with full rank ( $r = N$ ), the compact SVD (used in our mathematical analysis) and the economy SVD (returned by `svd`) are identical, i.e.,  $\mathbf{U}_r = \mathbf{U}_N$ ,  $\mathbf{\Sigma}_r = \mathbf{\Sigma}_N$ , and  $\mathbf{V}_r = \mathbf{V}_N = \mathbf{V}$ .

**Uniqueness of LLS solution** (reprise)

(Read)

If  $\mathbf{A}$  is  $M \times N$ , then by the definition of rank,  $\text{rank}(\mathbf{A}) = N$  iff  $\mathbf{A}$  has linearly independent columns. Having  $M \geq N$  is a necessary condition for linear independence of the columns of  $\mathbf{A}$ .

Conversely, if  $M < N$  then  $\text{rank}(\mathbf{A}) \neq N$  because  $\text{rank}(\mathbf{A}) \leq \min(M, N) = M < N$ .

**Over-determined (tall) case**

In the frequent case where  $M > N$  we have (graphically):

$$\underbrace{\mathbf{A}}_{M \times N} \underbrace{\mathbf{x}}_N = \underbrace{\mathbf{y}}_M.$$

In words, there are more equations than unknowns, called an **over-determined system**. Rarely is there an exact solution in this case due to noise in the data  $\mathbf{y}$ , so “best fit” solutions like LLS (4.2) are used instead.

When  $M \geq N$  and  $\mathbf{A}$  has rank  $r = N$ , the solution (4.9), (4.10) has no arbitrary terms and simplifies to:

$$\underbrace{\hat{z}_n = \frac{[\mathbf{U}'\mathbf{y}]_n}{\sigma_n}}_{n=1, \dots, N} \implies \hat{\mathbf{z}} = \Sigma_N^{-1} \mathbf{U}_N' \mathbf{y}, \quad \hat{\mathbf{x}} = \mathbf{V} \hat{\mathbf{z}} \implies \hat{\mathbf{x}} = \mathbf{V} \Sigma_N^{-1} \mathbf{U}_N' \mathbf{y}. \quad (4.12)$$

This is the *unique* solution to (4.2) when  $\text{rank}(\mathbf{A}) = N$ .

**Over-determined full-rank case using SVD**

The expression (4.12) uses the **compact SVD** and it is useful to also write it in terms of the **full SVD**:

$$\hat{\mathbf{x}} = \mathbf{V} \Sigma_N^{-1} \mathbf{U}_N' \mathbf{y} = \mathbf{V} \left[ \Sigma_N^{-1} \mid \mathbf{0}_{N \times (M-N)} \right] \begin{bmatrix} \mathbf{U}_N' \\ \mathbf{U}_0' \end{bmatrix} \mathbf{y} = \mathbf{V} \left[ \Sigma_N^{-1} \mid \mathbf{0}_{N \times (M-N)} \right] \mathbf{U}' \mathbf{y}.$$

This form is suboptimal for *computation* because the term  $\mathbf{U}_0' \mathbf{y}$  is computed but then multiplied by 0.

An even more concise expression is

$$M \geq N \text{ and } \text{rank}(\mathbf{A}) = N \implies \text{[redacted]} \quad (4.13)$$

where  $\Sigma^+ \triangleq \left[ \Sigma_N^{-1} \mid \mathbf{0}_{N \times (M-N)} \right]$  denotes the **Moore-Penrose pseudoinverse** of  $\Sigma = \begin{bmatrix} \Sigma_N \\ \mathbf{0}_{(M-N) \times N} \end{bmatrix}$ .

The next page discusses this new term in detail and explains  $\mathbf{A}^+$ .

## 4.2 Moore-Penrose pseudoinverse

L§4.1

Define. In general, for any matrix  $\mathbf{A} \in \mathbb{F}^{M \times N}$ , the **Moore-Penrose pseudoinverse** of  $\mathbf{A}$ , denoted  $\mathbf{A}^+ \in \mathbb{F}^{N \times M}$ , is a generalization of the usual notation of matrix inverse that satisfies the following four properties:

- (weaker than  $\mathbf{A}\mathbf{A}^{-1} = \mathbf{I}$ )
- 
- (symmetry)
- 

Properties. (For proofs see [\[wiki\]](#).)

(Another notation is  $\mathbf{A}^\dagger$ .) L§4.3

- $\mathbf{A}^+$  is unique and  $(\mathbf{A}^+)^+ = \mathbf{A}$
- If  $\mathbf{A}$  is invertible, then  $\mathbf{A}^+ = \mathbf{A}^{-1}$
- $\mathbf{A}^+ = (\mathbf{A}'\mathbf{A})^+\mathbf{A}'$
- If  $\mathbf{A}$  has full column rank (linearly independent columns), then  $\mathbf{A}^+ =$    
 In this case,  $\mathbf{A}^+$  is a **left inverse** because  $\mathbf{A}^+\mathbf{A} = \mathbf{I}_N$ .
- $\mathbf{A}^+ = \mathbf{A}'(\mathbf{A}\mathbf{A}')^+$
- If  $\mathbf{A}$  has full row rank, (linearly independent rows), then  $\mathbf{A}^+ = \mathbf{A}'(\mathbf{A}\mathbf{A}')^{-1}$ .  
 In this case,  $\mathbf{A}^+$  is a **right inverse** because  $\mathbf{A}\mathbf{A}^+ = \mathbf{I}_M$ .
- $\mathbf{0}_{M \times N}^+ = \mathbf{0}_{N \times M}$  so in particular  $\mathbf{0}^+ = \mathbf{0}$
- $(\mathbf{A}')^+ = (\mathbf{A}^+)'$



## Pseudo-inverse and matrix products

Caution: in general  $(AB)^+ \neq B^+A^+$ , so be careful with matrix products. However, for  $B \in \mathbb{F}^{K \times L}$ :

- (P1) If  $Q$  is a  $M \times K$  matrix with **orthonormal columns**, i.e.,  $Q'Q = I_{K \times K}$ , then  $(QB)^+ = B^+Q'$ .



Partial proof:

- $(QB)(QB)^+(QB) = QB(B^+Q')QB = QBB^+B = QB$
  - $(QB)^+QB(QB)^+ = (B^+Q')QB(B^+Q') = B^+BB^+Q' = B^+Q' = (QB)^+$ , etc.
- (P2) If  $Q$  is a  $L \times N$  matrix with **orthonormal rows**, i.e.,  $QQ' = I_{L \times L}$ , then  $(BQ)^+ = Q'B^+$ .
  - $(BQ)(BQ)^+(BQ) = BQ(Q'B^+)BQ = BB^+BQ = BQ$
  - $(BQ)^+(BQ)(BQ)^+ = (Q'B^+)BQ(Q'B^+) = Q'B^+BB^+ = Q'B^+ = (BQ)^+$ , etc.
- (P3) If  $A \in \mathbb{F}^{M \times N}$  has full column rank (linearly independent columns) and  $B \in \mathbb{F}^{N \times K}$  has full row rank (linearly independent rows), then  $(AB)^+ = B^+A^+$ .

Fortunately, these product properties serve our needs.

A special case of the first two preceding results is when  $B = I$ , leading to:

- (P4) If  $Q$  is a matrix with **orthonormal columns**, then  $Q^+ = Q'$ .
- (P5) If  $Q$  is a matrix with **orthonormal rows**, then  $Q^+ = Q'$ .

In the context of the **compact SVD**, (P4) implies that  $U_r^+ = U_r'$  and  $V_r^+ = V_r'$ .

Because  $U_r$  has orthonormal columns and  $V_r'$  has orthonormal rows and  $\Sigma_r$  is invertible, it follows that:

$$A = U_r \Sigma_r V_r' \implies A^+ = (U_r \Sigma_r V_r')^+ = (U_r (\Sigma_r V_r'))^+ \stackrel{\text{P1}}{=} (\Sigma_r V_r')^+ U_r' \stackrel{\text{P2}}{=} V_r \Sigma_r^+ U_r' = V_r \Sigma_r^{-1} U_r'.$$

Example. A pseudo-inverse example of particular interest when working with the **SVD** is that of a **rectangular diagonal matrix**:

$$\Sigma = \underbrace{\begin{bmatrix} \Sigma_r & \mathbf{0}_{r \times (N-r)} \\ \mathbf{0}_{(M-r) \times r} & \mathbf{0}_{(M-r) \times (N-r)} \end{bmatrix}}_{M \times N} \implies \Sigma^+ = \underbrace{\begin{bmatrix} \text{ } & \mathbf{0}_{r \times (M-r)} \\ \mathbf{0}_{(N-r) \times r} & \text{ } \end{bmatrix}}_{\text{ } }.$$

Exercise. Verify that this  $\Sigma^+$  satisfies the four conditions for a pseudo-inverse on p. 4.24.

An important special case is when  $r = N$ , i.e.,  $\mathbf{A}$  has full column rank, in which case:

$$\Sigma = \underbrace{\begin{bmatrix} \Sigma_N \\ \mathbf{0}_{(M-N) \times N} \end{bmatrix}}_{M \times N} \implies \Sigma^+ = \text{ }$$

Noting that  $\Sigma^+ \Sigma = \mathbf{I}_N$  in the tall full-rank case, it is trivial to verify the four defining properties for this special case.

Caution: in general  $\mathbf{A}^+ \mathbf{A} \neq \mathbf{I}_N$  and  $\mathbf{A} \mathbf{A}^+ \neq \mathbf{I}_M$ .



## Pseudo-inverse and SVD

L§5.2

Using the orthogonal matrix product properties of the pseudo-inverse yields the following **SVD** property:

$$\underbrace{\mathbf{A}}_{M \times N} = \underbrace{\mathbf{U}}_{M \times M} \underbrace{\mathbf{\Sigma}}_{M \times N} \underbrace{\mathbf{V}'}_{N \times N} \implies \underbrace{\mathbf{A}^+}_{N \times M} = \quad (4.14)$$

The **compact SVD** version of pseudo-inverse is also useful:

$$\mathbf{A}^+ = \mathbf{V} \mathbf{\Sigma}^+ \mathbf{U}' = \left[ \begin{array}{c|c} \mathbf{V}_r & \mathbf{V}_0 \end{array} \right] \mathbf{\Sigma}^+ \begin{bmatrix} \mathbf{U}'_r \\ \mathbf{U}'_0 \end{bmatrix} \implies \quad (4.15)$$

Here is one sanity check for the pseudo-inverse expression (4.14):

$$\mathbf{A} \mathbf{A}^+ \mathbf{A} = \mathbf{U} \mathbf{\Sigma} \underbrace{\mathbf{V}' \mathbf{V}} \mathbf{\Sigma}^+ \underbrace{\mathbf{U}' \mathbf{U}} \mathbf{\Sigma} \mathbf{V}' = \mathbf{U} \underbrace{\mathbf{\Sigma} \mathbf{\Sigma}^+ \mathbf{\Sigma}} \mathbf{V}' = \mathbf{U} \mathbf{\Sigma} \mathbf{V}' = \mathbf{A}.$$

The key expression (4.14) simplifies to the usual matrix inverse in the rare case where  $\mathbf{A}$  is square and invertible:  $\mathbf{A}^{-1} = \mathbf{V} \mathbf{\Sigma}^{-1} \mathbf{U}'$ . (However, for invertible matrices often we do not need to use the SVD.)

3.

LLS solution  $\hat{\mathbf{x}} = \mathbf{A}^+ \mathbf{y}$  lies in which of the four fundamental spaces of  $\mathbf{A}$ ?

A:  $\mathcal{N}(\mathbf{A})$ B:  $\mathcal{N}^\perp(\mathbf{A})$ C:  $\mathcal{R}(\mathbf{A})$ D:  $\mathcal{R}^\perp(\mathbf{A})$ 

E: None of these.

??

### Warm-up questions

4. If  $A$  has linearly independent columns, then  $\|A^+Ax\| = \|x\|$ .  
 A: True B: False ??
5. Let  $A$  have **full SVD**  $A = U\Sigma V'$  and **compact SVD**  $A = U_r\Sigma_r V_r'$ , where we partition unitary matrix  $V$  as usual as  $V = \begin{bmatrix} V_r & V_0 \end{bmatrix}$ . Which of the following is the tuple  $(\mathcal{N}(V'), \mathcal{N}(V_r'))$ ?  
 A:  $(0, 0)$  B:  $(0, \mathcal{R}(V_0))$  C:  $(\mathcal{R}(V_0), 0)$  D:  $(\mathcal{R}(V_0), \mathcal{R}(V_0))$  E: None of these. ??
6. If matrix  $A$  has **SVD**  $A = U\Sigma V' = U_r\Sigma_r V_r'$ , then two expressions for its **pseudo-inverse** are:  

$$A^+ = V\Sigma^+U' = V_r\Sigma_r^{-1}U_r'.$$
 Which of the following is the **null space** of  $A^+$ , i.e.,  $\mathcal{N}(A^+)$ ?  
 A:  $\mathcal{N}(A)$  B:  $\mathcal{N}^\perp(A)$  C:  $\mathcal{R}(A)$  D:  $\mathcal{R}^\perp(A)$  E: None of these. ??
7. If  $A = U_r\Sigma_r V_r'$  then  $A^+ = V_r\Sigma_r^{-1}U_r'$  is an SVD of  $A^+$ .  
 A: True B: False ??

**Projector / idempotent preview**

(Read)

The following matrix is called an **orthogonal projector** (see p. 4.56):

$$P_{\mathcal{R}(A')} = P_{\mathcal{N}^\perp(A)} \triangleq A^+ A = V \Sigma^+ U' U \Sigma V' = V \Sigma^+ \Sigma V' = V_r V_r'.$$

Because  $P_{\mathcal{R}(A')} P_{\mathcal{R}(A')} = P_{\mathcal{R}(A')}$ , it is called **idempotent**.

Likewise the following matrix is also a projector (and also idempotent):

$$P_{\mathcal{R}(A)} = P_{\mathcal{N}^\perp(A')} \triangleq A A^+ = U \Sigma V' V \Sigma^+ U' = U \Sigma \Sigma^+ U' = U_r U_r'.$$

**Relating pseudo-inverse and normal equations**

(Read)

The following equalities (properties of the pseudo-inverse) follow:

$$\begin{aligned} A' A A^+ &= A' \text{ because } A' A A^+ = V_r \Sigma_r U_r' (U_r U_r') = V_r \Sigma_r U_r' = A' \\ A^+ A A' &= A' \text{ because } A^+ A A' = (V_r V_r') V_r \Sigma_r U_r' = V_r \Sigma_r U_r' = A'. \end{aligned}$$

Multiplying the first of these two equalities by  $y$  yields:

$$A' A A^+ y = A' y \implies A' A \hat{x} = A' y,$$

so the pseudo-inverse solution  $\hat{x} = A^+ y$  satisfies the **normal equations** always, regardless of the rank of  $A$ .

## LLS solution using pseudo-inverse

---

Using (4.13), we write the LLS estimate (4.13) particularly concisely in the full-rank case as follows:

$$\mathbf{A} \in \mathbb{F}^{M \times N} \text{ and } \underbrace{\text{rank}(\mathbf{A}) = N}_{\Rightarrow M \geq N} \implies \hat{\mathbf{x}} = \arg \min_{\mathbf{x} \in \mathbb{F}^N} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 = \text{ } \quad (4.16)$$

This is a wonderfully elegant form on paper, but it is not the best computational approach!

We could implement (4.16) in JULIA using the pseudo-inverse function: `xh = pinv(A) * y;`

Alternatively we could use the `inv` form discussed on p. 4.16: `xh = inv(A' * A) * (A' * y);`

However, both these approaches are computationally inefficient!

We rarely use `pinv` or `inv` in practice (at least for large problem sizes) unless we must solve LLS problems for many different  $\mathbf{y}$  with the same  $\mathbf{A}$  (as in HW).

Otherwise, the backslash approach `xh = A \ y` discussed on p. 4.16 with its more efficient **QR decomposition** is preferable.

Nevertheless, the SVD is very helpful conceptually, especially in the under-determined case.

If  $\mathbf{A}$  does not have full column rank, then `A \ y` will produce an error message like `SingularException`. Then one should first think about why  $\mathbf{A}$  does not have linearly dependent columns, and then perhaps consider using `pinv`, as discussed in the next section, or one of the regularized solutions discussed after that.

### 4.3 Linear least-squares estimation: Under-determined case

We focused previously on cases where  $M \geq N$  and  $\mathbf{A}$  has full rank, and derived the concise LLS solution based on the pseudo-inverse (4.16). Now we examine cases where  $M < N$ , called **under-determined**, as well as cases where  $\mathbf{A}$  is tall but rank deficient *i.e.*,  $\text{rank}(\mathbf{A}) < N$ . In these cases there is not a unique minimizer to the LLS cost function  $\|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2$ . Using the **compact SVD** of  $\mathbf{A}$ , we showed in (4.9) that any minimizer has the form

$$\hat{\mathbf{x}} = \underbrace{\mathbf{V}_r \Sigma_r^{-1} \mathbf{U}_r' \mathbf{y}}_{\hat{\mathbf{x}}_{\mathcal{R}} \in \mathcal{R}(\mathbf{V}_r) = \mathcal{N}^\perp(\mathbf{A})} + \underbrace{\mathbf{V}_0 \hat{\mathbf{z}}_0}_{\hat{\mathbf{x}}_{\mathcal{N}} \in \mathcal{N}(\mathbf{A})} = \quad (4.17)$$

The choice of  $\hat{\mathbf{z}}_0$  is arbitrary because the residual is invariant to the nullspace component:

$$\mathbf{A}\hat{\mathbf{x}} - \mathbf{y} = \mathbf{A}(\mathbf{V}_r \hat{\mathbf{z}}_r + \mathbf{V}_0 \hat{\mathbf{z}}_0) - \mathbf{y} = \mathbf{A}(\hat{\mathbf{x}}_{\mathcal{R}} + \hat{\mathbf{x}}_{\mathcal{N}}) - \mathbf{y} = \mathbf{A}\hat{\mathbf{x}}_{\mathcal{R}} + \mathbf{0} - \mathbf{y} = \mathbf{A}\hat{\mathbf{x}}_{\mathcal{R}} - \mathbf{y}.$$

In other words, for any  $\hat{\mathbf{x}} = \hat{\mathbf{x}}_{\mathcal{R}} + \hat{\mathbf{x}}_{\mathcal{N}}$  where  $\hat{\mathbf{x}}_{\mathcal{N}} \in \mathcal{N}(\mathbf{A})$ , the error  $\|\mathbf{A}\hat{\mathbf{x}} - \mathbf{y}\|_2^2$  is identical.

When  $r < N$ , the LLS solution is *not* unique because we can choose  $\hat{\mathbf{z}}_0$  arbitrarily, or equivalently we can choose any  $\hat{\mathbf{x}}_{\mathcal{N}} \in \mathcal{N}(\mathbf{A})$ . The LLS criterion by itself is insufficient to identify a unique best  $\hat{\mathbf{x}}$  in under-determined (or rank-deficient) problems.

Conversely, in the full rank case where  $\text{rank}(\mathbf{A}) = r = N$ , the matrix  $\mathbf{V}_0$  does not exist and  $\mathcal{N}(\mathbf{A}) = \{\mathbf{0}\}$  so the LLS solution is unique.

If  $\mathbf{A}$  is wide, *i.e.*,  $M < N$ , can we have  $\text{rank}(\mathbf{A}) = N$ ?

No, because  $\text{rank}(\mathbf{A}) \leq \min(M, N) = M < N$  per (3.9). So wide cases are always under-determined.

To summarize:

- When  $M \geq N$  and  $\mathbf{A}$  has full rank  $N$ , the LLS solution is unique and is  $\hat{\mathbf{x}} = \mathbf{A}^+ \mathbf{y}$ .
- Otherwise (*i.e.*, if  $M < N$  or if  $\mathbf{A}$  has rank less than  $N$ ), the LLS solution is not unique and any  $\hat{\mathbf{x}}$  of the form  $\hat{\mathbf{x}} = \mathbf{A}^+ \mathbf{y} + \hat{\mathbf{x}}_{\mathcal{N}}$  where  $\hat{\mathbf{x}}_{\mathcal{N}} \in \mathcal{N}(\mathbf{A})$  is “equally good” from the point of view of the LLS cost function  $\|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2$ .

### Dealing with non-uniqueness

---

To pin down a unique solution in the under-determined case, we must introduce some additional criterion to select on  $\hat{\mathbf{x}}$  from the (infinitely) many candidates of the form  $\hat{\mathbf{x}} = \mathbf{A}^+ \mathbf{y} + \hat{\mathbf{x}}_{\mathcal{N}}$ .

A modern way to do this is to use a **sparsity** model and seek a sparse solution, *e.g.*:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \text{nonsparsity}(\mathbf{x}) \text{ s.t. } \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 \leq \epsilon.$$

But this is an advanced topic for later in the course!

Instead we focus for now on the “classical” approach of choosing the **minimum norm** solution, *i.e.*, the LLS minimizer where  $\|\hat{\mathbf{x}}\|_2^2$  is the smallest. This is a “double minimization” problem, because (conceptually) we first find all the minimizers of  $\|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2$ , and then among those minimizers we pick the one where  $\|\mathbf{x}\|_2^2$  is the smallest.

But first we give a geometric interpretation of the LLS solution.



## Orthogonality principle

To examine LLS geometrically, recall from the **normal equations** (4.5) that any LLS solution must satisfy

$$A'A\hat{x} = A'y \implies A'(y - A\hat{x}) = 0 \implies A'r = 0 \implies z'A'r = 0 \implies$$

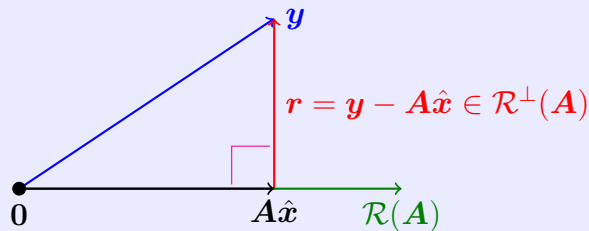
where the **residual** after fitting is denoted  $r = y - A\hat{x}$ .

(The normal equations are a **necessary condition** for  $\hat{x}$  to be a minimizer, regardless of the rank of  $A$ .)

The following **orthogonality principle** of LLS estimation is a direct consequence of the above:

$$\langle Az, y - A\hat{x} \rangle = (y - A\hat{x})'Az = 0, \text{ i.e., } \quad (4.18)$$

In words, the **residual** is **perpendicular** to  $\mathcal{R}(A)$ . There is an important geometric interpretation:



Terminology that (unfortunately) is often used inter-changeably:

- $\mathbf{y} - \mathbf{A}\hat{\mathbf{x}}$  is called the **residual** (from fitting)
- $\mathbf{y} - \mathbf{A}\mathbf{x}_{\text{true}}$  is called the **error** in the data
- $\hat{\mathbf{x}} - \mathbf{x}_{\text{true}}$  is called the **error** in the parameter estimate.

Often one must determine from context which meaning of “error” is meant.

If  $\mathbf{y} \in \mathcal{R}(\mathbf{A})$  (which rarely happens for noisy data when  $M > N$ ), then there is a solution with zero residual  $\mathbf{y} - \mathbf{A}\hat{\mathbf{x}}$ . However, the parameter error  $\hat{\mathbf{x}} - \mathbf{x}_{\text{true}}$  may still be nonzero!

**Alternate proof of optimality**

(Read)

Here is an alternate derivation that uses the orthogonality principle (4.18) to confirm that  $\hat{\mathbf{x}}$  in (4.17) is optimal for LLS estimation (for  $\mathbf{A}$  of any size or rank).

Suppose  $\mathbf{x}_? = \mathbf{A}^+ \mathbf{y} + \hat{\mathbf{x}}_{\mathcal{N}} + \mathbf{z} = \hat{\mathbf{x}} + \mathbf{z}$  for some arbitrary vector  $\mathbf{z} \in \mathbb{F}^N$ .

Could this  $\mathbf{x}_?$  be a better estimator?

Using the norm of a sum in (1.17) and the orthogonality principle (4.18), the squared error criterion for such an  $\mathbf{x}_?$  has the following lower bound:

$$\begin{aligned} \|\mathbf{A}\mathbf{x}_? - \mathbf{y}\|_2^2 &= \|\mathbf{A}(\hat{\mathbf{x}} + \mathbf{z}) - \mathbf{y}\|_2^2 = \|(\mathbf{A}\hat{\mathbf{x}} - \mathbf{y}) + \mathbf{A}\mathbf{z}\|_2^2 \\ &= \|\mathbf{A}\hat{\mathbf{x}} - \mathbf{y}\|_2^2 + \underbrace{2 \operatorname{real}\{(\mathbf{A}\hat{\mathbf{x}} - \mathbf{y})' \mathbf{A}\mathbf{z}\}}_0 + \underbrace{\|\mathbf{A}\mathbf{z}\|_2^2}_{\geq 0} \geq \|\mathbf{A}\hat{\mathbf{x}} - \mathbf{y}\|_2^2, \end{aligned}$$

where the lower bound is achieved by  $\mathbf{z} = \mathbf{0}$ .

In words, the LLS fit is best when  $\hat{\mathbf{x}} = \mathbf{A}^+ \mathbf{y} + \hat{\mathbf{x}}_{\mathcal{N}}$  for any vector  $\hat{\mathbf{x}}_{\mathcal{N}} \in \mathcal{N}(\mathbf{A})$ .

Mathematically, the set of LLS solutions is the sum of a vector and a subspace (the nullspace of  $\mathbf{A}$ ):

$$\begin{aligned} \{\tilde{\mathbf{x}} \in \mathbb{F}^N : \|\mathbf{A}\tilde{\mathbf{x}} - \mathbf{y}\|_2^2 \leq \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2, \forall \mathbf{x} \in \mathbb{F}^N\} &= \text{[yellow box]} \\ &\triangleq \{\hat{\mathbf{x}} = \mathbf{A}^+ \mathbf{y} + \hat{\mathbf{x}}_{\mathcal{N}} : \hat{\mathbf{x}}_{\mathcal{N}} \in \mathcal{N}(\mathbf{A})\}. \end{aligned}$$

The sum of a vector plus a subspace is called a **linear variety** or a **flat**.

## Minimum-norm LS solution via pseudo-inverse

L§8.1

We have seen that the set of optimal LLS estimates is  $\hat{\mathbf{x}} = \mathbf{A}^+ \mathbf{y} + \mathcal{N}(\mathbf{A})$ .

- If  $\mathbf{A}$  has full column rank, then  $\mathcal{N}(\mathbf{A}) = \mathbf{0}$  and we have a unique solution  $\hat{\mathbf{x}} = \mathbf{A}^+ \mathbf{y}$
- If  $\mathbf{A}$  does not have full column rank, then we want to pick one choice from the set of LLS estimates.

The classical way to pick one of the many possible estimates in the under-determined case is to choose the one with minimum norm.

Fact. The **minimum norm LLS solution** is:

$$\hat{\mathbf{x}} \triangleq \arg \min_{\mathbf{x} \in \{\mathbf{A}^+ \mathbf{y} + \mathcal{N}(\mathbf{A})\}} \|\mathbf{x}\|_2^2 = \text{ } \quad (4.19)$$

This is a kind of “double minimization” because first we found a set of candidate solutions by finding minimizers of  $\|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2$ , and then we solve a different minimization problem involving  $\|\mathbf{x}\|_2^2$  to select one final solution from that set of candidates.

Proof. If  $\mathbf{x} = \mathbf{A}^+ \mathbf{y} + \hat{\mathbf{x}}_{\mathcal{N}}$  where  $\hat{\mathbf{x}}_{\mathcal{N}} \in \mathcal{N}(\mathbf{A})$ , then  $\hat{\mathbf{x}}_{\mathcal{N}} \in \text{span}(\mathbf{V}_0)$ , where  $\mathbf{V} = [\mathbf{V}_r \mid \mathbf{V}_0]$ . (Read)

Using (4.15):  $\mathbf{A}^+ \mathbf{y} = \mathbf{V}_r \Sigma_r^{-1} \mathbf{U}_r' \mathbf{y} \in \mathcal{R}(\mathbf{V}_r)$ .

Because  $\mathbf{V}$  is unitary, the columns of  $\mathbf{V}_r$  and  $\mathbf{V}_0$  are orthogonal. Thus  $\mathbf{A}^+ \mathbf{y} \perp \hat{\mathbf{x}}_{\mathcal{N}}$ . Using (1.17):

$$\|\mathbf{x}\|_2^2 = \|\mathbf{A}^+ \mathbf{y} + \hat{\mathbf{x}}_{\mathcal{N}}\|_2^2 = \|\mathbf{A}^+ \mathbf{y}\|_2^2 + 2 \operatorname{real}\{\hat{\mathbf{x}}_{\mathcal{N}}' \mathbf{A}^+ \mathbf{y}\} + \|\hat{\mathbf{x}}_{\mathcal{N}}\|_2^2 = \|\mathbf{A}^+ \mathbf{y}\|_2^2 + \|\hat{\mathbf{x}}_{\mathcal{N}}\|_2^2 \geq \|\mathbf{A}^+ \mathbf{y}\|_2^2,$$

where the minimum is achieved when  $\hat{\mathbf{x}}_{\mathcal{N}} = \mathbf{0}$ . Thus the minimum norm solution is  $\hat{\mathbf{x}} = \mathbf{A}^+ \mathbf{y}$ .  $\square$

The set over which we minimize in (4.19) is  $\{\mathbf{A}^+\mathbf{y} + \mathcal{N}(\mathbf{A})\}$ .

- |    |   |
|----|---|
| 8. | What is the cardinality of this set when $\text{rank}(\mathbf{A}) = N$ ?  |
|    | <div style="display: flex; justify-content: space-between; padding: 0 10px;"> <span>A: 0</span> <span>B: 1</span> <span>C: <math>r</math></span> <span>D: <math>N</math></span> <span>E: <math>\infty</math></span> <span style="border: 1px solid black; padding: 2px 5px;">??</span> </div> |
| 9. | What is the cardinality of this set when $M < N$ ?  |
|    | <div style="display: flex; justify-content: space-between; padding: 0 10px;"> <span>A: 0</span> <span>B: 1</span> <span>C: <math>r</math></span> <span>D: <math>N</math></span> <span>E: <math>\infty</math></span> <span style="border: 1px solid black; padding: 2px 5px;">??</span> </div> |

In summary we have the following fortuitous situation.

- If  $\mathbf{A}$  has full column rank, then  $\mathcal{N}(\mathbf{A}) = \mathbf{0}$  and we have a unique solution  $\hat{\mathbf{x}} = \mathbf{A}^+\mathbf{y}$ .
- If  $\mathbf{A}$  does not have full column rank, then there are multiple LLS estimates and the one with smallest Euclidean norm is  $\hat{\mathbf{x}} = \mathbf{A}^+\mathbf{y}$ .

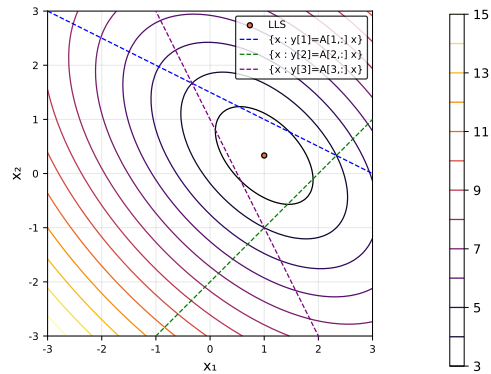
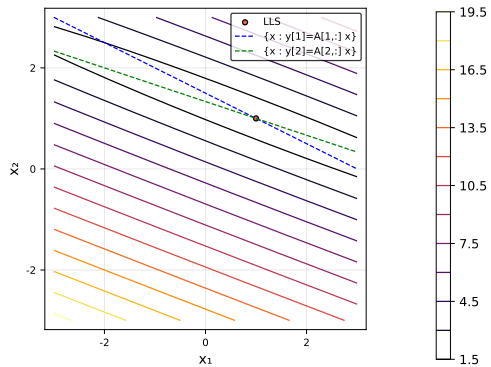
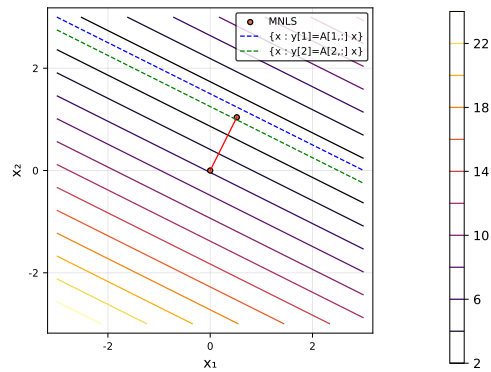
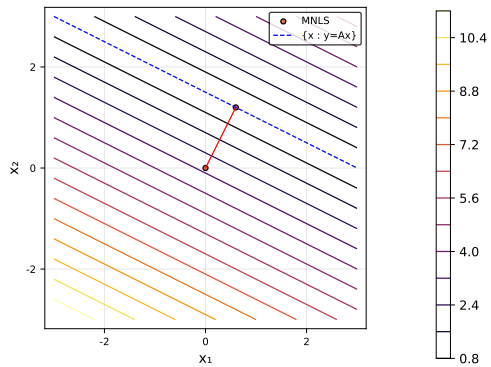
Hence the **pseudo-inverse** solution has been very popular historically, and remains important today, except in many highly under-determined problems of the kind known as **compressed sensing**.

There are **quantum-computing methods** for the pseudo-inverse solution.

---

Example. To visualize the cost function (4.3) and solution  $\hat{\mathbf{x}}$  see:

[https://web.eecs.umich.edu/~fessler/course/551/julia/demo/04\\_ls\\_cost1.html](https://web.eecs.umich.edu/~fessler/course/551/julia/demo/04_ls_cost1.html)  
[https://web.eecs.umich.edu/~fessler/course/551/julia/demo/04\\_ls\\_cost1.ipynb](https://web.eecs.umich.edu/~fessler/course/551/julia/demo/04_ls_cost1.ipynb)  
 and plots on next page.



**Interpreting the pseudo-inverse solution**

(Read)

Using (4.14), we can interpret the pseudo-inverse solution

$$\hat{\mathbf{x}} = \mathbf{A}^+ \mathbf{y} = \mathbf{V} \Sigma^+ \mathbf{U}' \mathbf{y}$$

as a cascade of three separate transformations:

$$\mathbf{y} \xrightarrow{\mathbf{U}'} \tilde{\mathbf{y}} \xrightarrow{\Sigma^+} \mathbf{z} \xrightarrow{\mathbf{V}} \hat{\mathbf{x}}.$$

The geometric interpretation is similar to what we did for the SVD except “in reverse” and with  $\Sigma^+$ .

To interpret the residual, note that

$$\begin{aligned} \mathbf{r} &= \mathbf{y} - \mathbf{A}\hat{\mathbf{x}} = \mathbf{y} - \mathbf{A}\mathbf{A}^+ \mathbf{y} = (\mathbf{I} - \mathbf{A}\mathbf{A}^+) \mathbf{y} = (\mathbf{I} - \mathbf{U}\Sigma\Sigma^+ \mathbf{U}') \mathbf{y} = \mathbf{U}(\mathbf{I} - \Sigma\Sigma^+) \mathbf{U}' \mathbf{y} \\ &= [\mathbf{U}_r \mid \mathbf{U}_0] \left( \mathbf{I} - \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \right) [\mathbf{U}_r \mid \mathbf{U}_0]' \mathbf{y} = [\mathbf{U}_r \mid \mathbf{U}_0] \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} [\mathbf{U}_r \mid \mathbf{U}_0]' \mathbf{y} = \mathbf{U}_0 \mathbf{U}_0' \mathbf{y} \\ &\implies \|\mathbf{r}\|_2^2 = \|\mathbf{U}_0 \mathbf{U}_0' \mathbf{y}\|_2^2 = \mathbf{y}' \mathbf{U}_0 \mathbf{U}_0' (\mathbf{U}_0 \mathbf{U}_0' \mathbf{y}) = \mathbf{y}' \mathbf{U}_0 \mathbf{U}_0' \mathbf{y} = \|\mathbf{U}_0' \mathbf{y}\|_2^2, \end{aligned}$$

where from p. 3.37,  $\mathcal{R}^\perp(\mathbf{A}) = \text{span}(\mathbf{U}_0)$ . Thus the residual norm squared comes from the portion of  $\mathbf{y}$  in  $\mathcal{R}^\perp(\mathbf{A})$ , consistent with our earlier picture about the orthogonality principle.

---

#### 4.4 Truncated SVD solution

We have seen that the minimum-norm LLS minimizer of  $\|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2$  is the pseudo-inverse solution:

$$\hat{\mathbf{x}} = \mathbf{A}^+ \mathbf{y} = \mathbf{V} \Sigma^+ \mathbf{U}' \mathbf{y} = \mathbf{V}_r \Sigma_r^{-1} \mathbf{U}_r' \mathbf{y} =$$

This solution is mathematically elegant, but in practice sometimes it works very poorly.

This section provides one remedy based on the **truncated SVD**.

---

Example. Consider an application where  $r = \text{rank}(\mathbf{A}) = 2$  with singular values  $\sigma_1 = 1$ ,  $\sigma_2 = 10^{-8}$ . Then the pseudo-inverse solution here is:

$$\hat{\mathbf{x}} = \mathbf{A}^+ \mathbf{y} = \sum_{k=1}^2 \frac{1}{\sigma_k} \mathbf{v}_k (\mathbf{u}_k' \mathbf{y}) = \frac{1}{\sigma_1} \mathbf{v}_1 (\mathbf{u}_1' \mathbf{y}) + \frac{1}{\sigma_2} \mathbf{v}_2 (\mathbf{u}_2' \mathbf{y}) = (1) \mathbf{v}_1 (\mathbf{u}_1' \mathbf{y}) + (10^8) \mathbf{v}_2 (\mathbf{u}_2' \mathbf{y}).$$

The pseudo-inverse solution “blows up” when any of the  $\sigma_k$  values are “too small” relative to the others, *i.e.*, if  $\sigma_k/\sigma_1$  is near the **floating-point precision** limits.



## Condition number

Such problems are called **poorly conditioned** because the  $\mathbf{A}$  has an undesirably large **condition number**, defined as

$$\kappa(\mathbf{A}) \triangleq \begin{cases} \infty, & M < N \\ \text{[blank]}, & M \geq N \end{cases} = \begin{cases} \infty, & M < N \\ \text{[blank]}, & M \geq N \end{cases} = \text{[blank]} = \text{[blank]} \quad (4.20)$$

Caution: an alternate definition of condition number is  $\sigma_1/\sigma_r$ , e.g., [5, p. 69], but we will use (4.20).

Caution: Many sources consider only the case where  $\mathbf{A}$  is square; hence the expression  $\mathbf{A}'\mathbf{A}$  above.



This problem motivates the **truncated SVD** solution where we discard any singular values that are “too small” and write:

$$\hat{\mathbf{x}}_K = \text{[blank]} \quad (4.21)$$

where we choose fewer terms,  $K < r$ , such that  $\sigma_K > \delta > 0$  for some tolerance  $\delta$ .

## Practical implementation of truncated SVD solution (Read)

The tolerance  $\delta$  may depend on factors such as the noise level in the data, the size of  $\mathbf{A}$ , and whether one is using half, single, double or extended precision variables.

In JULIA, these variable types are `Float16`, `Float32`, `Float64`, and `BigFloat`.

JULIA's `pinv` has an optional second argument for specifying the tolerance.

The default value (see [HW](#)) is `eps(real(float(one(eltype(M)))))*minimum(size(A))`

The backslash function in JULIA `xh = A \ y` does not have any tolerance parameter, so one must use `pinv` instead of backslash to control the tolerance for poorly conditioned problems.

Alternatively, one can use some other method to improve the condition number, such as **Tikhonov regularization**, also called **ridge regression**, described later on p. [4.46](#).

## Low-rank approximation interpretation of truncated SVD solution

L§8.2

One way to interpret the truncated SVD solution (4.21) is as follows.

- First we form a **low-rank approximation**  $\mathbf{A}_K$  of  $\mathbf{A}$ , defined as

$$\mathbf{A}_K = \mathbf{U}_K \mathbf{\Sigma}_K \mathbf{V}_K' =$$

with  $K < r \leq \min(M, N)$ . (See Ch. 6 for more details about such approximations.)

- Then we express the pseudo-inverse LLS solution using that approximation:

$$\hat{\mathbf{x}} =$$

which is the same expression as (4.21). This approach is also called **principal component regression** [6].

We visualize a low-rank approximation as follows.

$$\underbrace{\mathbf{A}}_{M \times N} \approx \underbrace{\mathbf{A}_K}_{M \times N} = \underbrace{\mathbf{U}_K}_{M \times K} \underbrace{\mathbf{\Sigma}_K}_{K \times K} \underbrace{\mathbf{V}_K'}_{K \times N}$$

## Noise effects and perturbations

(Read)

There are two sources of perturbations in LLS problems that can degrade the estimate  $\hat{\mathbf{x}} = \mathbf{A}^+ \mathbf{y}$ .

- Additive noise:  $\mathbf{y} = \mathbf{A}\mathbf{x} + \boldsymbol{\varepsilon}$
- Errors in the model  $\mathbf{A}$ .

The  $\mathbf{A}$  we use for estimation might differ from the “true model”  $\mathbf{A}_{\text{true}}$ . Define  $\Delta\mathbf{A} \triangleq \mathbf{A}_{\text{true}} - \mathbf{A}$ .

[5, Theorem 6.12, p. 69]. (See [7].) Let  $\mathbf{A} \in \mathbb{F}^{M \times N}$  with  $M \geq N$  and  $\text{rank}(\mathbf{A}) = N$ .

Let  $\kappa \triangleq \sigma_1/\sigma_N = \|\mathbf{A}\|_2/\sigma_N$  denote the condition number of  $\mathbf{A}$ , using the matrix 2-norm defined in (2.5).

Assume that the model perturbations are not “too large” as follows:

$$\eta \triangleq \frac{\|\Delta\mathbf{A}\|_2}{\sigma_N} = \kappa\epsilon_A < 1, \quad \epsilon_A \triangleq \frac{\|\Delta\mathbf{A}\|_2}{\|\mathbf{A}\|_2}.$$

If the perturbed matrix  $\mathbf{A} + \Delta\mathbf{A}$  has full rank, then the solution perturbation  $\delta\hat{\mathbf{x}} = \hat{\mathbf{x}} - \mathbf{x}_{\text{true}}$  has the following bound in terms of the residual  $\mathbf{r} = \mathbf{A}\hat{\mathbf{x}} - \mathbf{y}$ :

$$\|\delta\hat{\mathbf{x}}\|_2 \leq \frac{\kappa}{1 - \eta} \left( \epsilon_A \|\mathbf{x}\|_2 + \frac{\|\boldsymbol{\varepsilon}\|_2}{\|\mathbf{A}\|_2} + \epsilon_A \kappa \frac{\|\mathbf{r}\|_2}{\|\mathbf{A}\|_2} \right).$$

See also [7–11].

- 
- In the (full-rank) square case where  $M = N$ , the residual is  $\mathbf{r} = \mathbf{0}$ , so the solution error  $\|\delta\hat{\mathbf{x}}\|_2$  depends on the condition number  $\kappa$ .
  - In the usual over-determined case where  $M \geq N$  and  $\mathbf{y} \notin \mathcal{R}(\mathbf{A})$ , then the solution error is proportional to  $\kappa^2$  so it is particularly important to try to keep  $\kappa$  small.

This bound is a motivation for using the truncated SVD where  $\sigma_1/\sigma_K$  is better (lower) than  $\sigma_1/\sigma_N$ .

Challenge. Find an error bound that depends on the truncated SVD where  $\text{rank}(\mathbf{A}) = K$ , instead of the full SVD where  $\text{rank}(\mathbf{A}) = N$  was assumed above.

### Tikhonov regularization aka ridge regression

A drawback of the **truncated SVD** solution to LLS problems is that it requires one to compute an SVD of  $\mathbf{A}$ , which can be impractical for large problems. An alternate approach to address ill-conditioned problems is to use **Tikhonov regularization**, also known as **ridge regression**.

We saw that the pseudo-inverse solution to LLS problems involves  $1/\sigma_k$  terms that can “blow up” for small singular values, leading to very large values of  $\hat{\mathbf{x}}$ . Instead of directly modifying the singular values, **Tikhonov regularization** modifies the LS cost function to include a term that discourages the estimate from having excessively high values:

$$\hat{\mathbf{x}}_\beta = \arg \min_{\mathbf{x} \in \mathbb{R}^N} \quad \text{(4.22)}$$

where  $\beta > 0$  is a **regularization parameter** that one must **tune** to trade-off between how well  $\hat{\mathbf{x}}_\beta$  fits the data and how high is the energy of  $\hat{\mathbf{x}}_\beta$ .

By combining terms, we can rewrite the Tikhonov estimate (4.22) as:

$$\hat{\mathbf{x}}_\beta = \arg \min_{\mathbf{x} \in \mathbb{R}^N} \left\| \begin{bmatrix} \mathbf{A} \\ \sqrt{\beta} \mathbf{I} \end{bmatrix} \mathbf{x} - \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix} \right\|_2^2 = \arg \min_{\mathbf{x} \in \mathbb{R}^N} \left\| \tilde{\mathbf{A}} \mathbf{x} - \tilde{\mathbf{y}} \right\|_2^2, \quad \tilde{\mathbf{A}} \triangleq \begin{bmatrix} \text{ } \\ \text{ } \end{bmatrix}, \quad \tilde{\mathbf{y}} \triangleq \begin{bmatrix} \text{ } \\ \text{ } \end{bmatrix}.$$

10.

If  $\mathbf{A}$  is  $M \times N$ , how many rows does  $\tilde{\mathbf{A}}$  have?

A:  $M$ B:  $N$ C:  $M + N$ D:  $2M$ E:  $2N$ 

??

11.

What is the rank of  $\tilde{\mathbf{A}}$ ?A:  $r$ B:  $M$ C:  $N$ D:  $M + N$ 

E: None of these.

☐

In this simplified form, we know that the (unique!) LLS solution to (4.22) is

$$\hat{\mathbf{x}}_{\beta} = \tilde{\mathbf{A}}^+ \tilde{\mathbf{y}} =$$

If  $N$  is large, this solution requires inverting a  $N \times N$  matrix (actually, solving a  $N \times N$  system of equations) which is impractical. Instead we usually apply an optimization approach (such as a **conjugate gradient method**) directly to (4.22). Nevertheless, the closed-form expression for the solution is useful for analysis.

In particular it is insightful to examine the solution in terms of an **SVD**  $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}'$ . Here,

$$\mathbf{A}'\mathbf{A} + \beta\mathbf{I} =$$

so

$$\hat{\mathbf{x}}_{\beta} = (\mathbf{A}'\mathbf{A} + \beta\mathbf{I})^{-1} \mathbf{A}'\mathbf{y} =$$

=

If  $\beta \rightarrow 0$  then the ratio  $\frac{\sigma_k}{\sigma_k^2 + \beta} \rightarrow \frac{1}{\sigma_k}$ , unless  $\sigma_k = 0$  (which never happens for  $k = 1, \dots, r$ ).

So  $\hat{\mathbf{x}}_{\beta} \rightarrow \mathbf{A}^+ \mathbf{y}$  as  $\beta \rightarrow 0$ .

## Regularization parameter selection (Read)

An important practical question is how does one choose the regularization parameter  $\beta$ . Unsupervised methods include **cross validation** and **Stein's unbiased risk estimate (SURE)**. One can also pursue supervised approaches when training data is available. The details are beyond the scope of this course, but the foundations from this course are essential for understanding such methods. ♦♦



### 4.5 Summary of LLS solution methods in terms of SVD

This chapter has discussed three different ways of solving the linear least-squares (LLS) problem (4.2), each of which have SVD expressions as follows.

- “Optimal” solution (minimum-norm LLS):

$$\hat{\mathbf{x}} = \sum_{k=1}^r \frac{1}{\sigma_k} \mathbf{v}_k (\mathbf{u}'_k \mathbf{y})$$

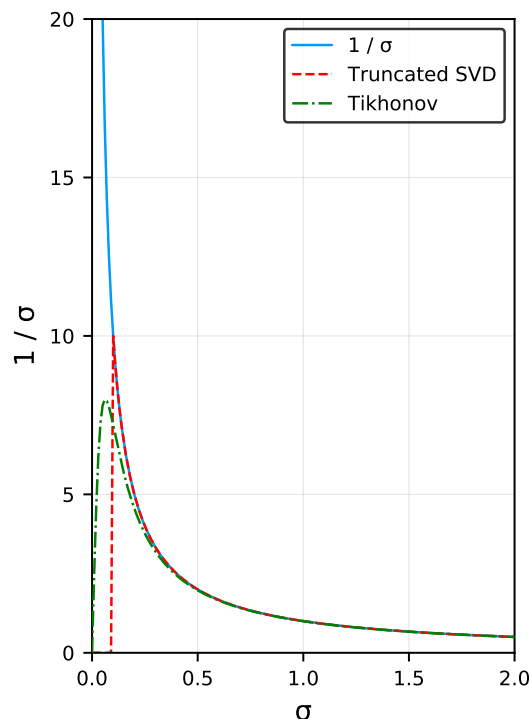
- Truncated SVD:

$$\hat{\mathbf{x}} = \sum_{k=1}^{K < r} \frac{1}{\sigma_k} \mathbf{v}_k (\mathbf{u}'_k \mathbf{y})$$

- Tikhonov regularized:

$$\hat{\mathbf{x}} = \sum_{k=1}^r \left( \frac{\sigma_k}{\sigma_k^2 + \beta} \right) \mathbf{v}_k (\mathbf{u}'_k \mathbf{y})$$

The SVD is a key tool for understanding LLS problems.



### 4.6 Frames and tight frames

A LS minimization problem  $\arg \min_{\mathbf{x}} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2$  is the easiest to solve when  $\mathbf{A}$  is unitary because in that case  $\hat{\mathbf{x}} = \mathbf{A}'\mathbf{y}$ . We now discuss a generalization of unitary matrices that leads to equally easy solutions.

For simplicity we focus on the case of a finite number of vectors in a finite-dimensional vector space [13]. The concepts generalize to a countable collection of vectors in a general Hilbert space.

**Define.** A collection of  $M$  vectors  $\{\phi_1, \dots, \phi_M\}$  in  $\mathbb{F}^N$  is called a **frame** in  $\mathbb{F}^N$  iff there exist real numbers  $0 < \alpha \leq \beta < \infty$ , called the **frame bounds**, such that

$$\alpha \|\mathbf{x}\|_2^2 \leq \sum_{m=1}^M |\langle \phi_m, \mathbf{x} \rangle|^2 \leq \beta \|\mathbf{x}\|_2^2, \quad \forall \mathbf{x} \in \mathbb{F}^N. \quad (4.23)$$

In other words, if we arrange those vectors into a  $N \times M$  matrix  $\Phi \triangleq [\phi_1 \ \dots \ \phi_M]$ , then the collection of vectors is a **frame** iff there exist real numbers  $0 < \alpha \leq \beta < \infty$  such that

$$\alpha \|\mathbf{x}\|_2^2 \leq \|\Phi' \mathbf{x}\|_2^2 \leq \beta \|\mathbf{x}\|_2^2, \quad \forall \mathbf{x} \in \mathbb{F}^N. \quad (4.24)$$

For brevity, we call such a  $\Phi$  a **frame**.

The upper bound is important in infinite-dimensional inner product spaces, but is not very informative in  $\mathbb{F}^N$ .

If  $\alpha > 0$ , then what is  $\beta$  here in terms of the **singular values** of  $\Phi$ ?

A:  $\sigma_1$

B:  $\sigma_1^2$

C:  $\sigma_r^2$

D:  $\sigma_M$

E: None of these

??

When is the upper frame bound  $\beta$  positive? *Always unless  $\Phi = 0$ .*

So in  $\mathbb{F}^N$ , the key to whether  $\Phi$  is a frame or not depends on existence of  $\alpha > 0$ .

### Properties of a frame

- Fact.  $\alpha > 0$  in (4.24) iff  $\Phi$  has full **row rank** (i.e.,  $\text{rank}(\Phi) = N$ ).
- Thus  $\Phi$  must be wide (usually) or square, i.e.,  $M \geq N$ .
- Thus  $\alpha > 0 \implies \Phi = U \Sigma_N V_N'$  (**compact SVD** is same as **economy SVD**) where  $\sigma_N > 0$ , because  $U_r = U_N = U$  and  $V_r = V_N$ .
- So  $\Phi^+ = V_N \Sigma_N^{-1} U' = \Phi'(\Phi\Phi')^{-1}$  and  $\Phi\Phi^+ = I_N$ .

Now suppose we have a vector  $x \in \mathbb{F}^N$  that we would like to express as a linear combination of the frame vectors  $\{\phi_1, \dots, \phi_M\}$ , i.e., we want to write

$$x = \Phi c$$

for some coefficient vector  $c \in \mathbb{F}^M$ . In the usual case where  $\Phi$  is wide, there will be numerous possible coefficient vectors  $c$  for which  $\|x - \Phi c\|_2 = 0$ . The unique such  $c$  having minimum norm is given by the Moore-Penrose pseudo-inverse:

$$c = \Phi^+ x = \Phi'(\Phi\Phi')^{-1} x.$$

This is not too exciting yet because of the expensive matrix inverse. So we impose more conditions on  $\Phi$ .

## Tight frame

Define. A matrix  $\Phi$  is a **tight frame** [14–16] iff  $\Phi$  is a **frame** with  $\alpha = \beta$ , i.e.,

$$\alpha \|\mathbf{x}\|_2^2 = \|\Phi' \mathbf{x}\|_2^2 = \sum_{m=1}^M |\langle \phi_m, \mathbf{x} \rangle|^2, \quad \forall \mathbf{x} \in \mathbb{F}^N. \quad (4.25)$$

If  $\Phi$  is a **tight frame**, then

- $\Phi \Phi' = \alpha \mathbf{I}_N$ , so  $\alpha = \sigma_1^2 = \dots = \sigma_N^2$ , where  $\{\sigma_k\}$  denotes the singular values of  $\Phi$ ,
- its pseudo-inverse is simply  $\Phi^+ = \frac{1}{\alpha} \Phi' = \frac{1}{\sigma_1^2} \Phi'$ .

Proof. Using (4.25):

$$\alpha \|\mathbf{x}\|_2^2 = \|\Phi' \mathbf{x}\|_2^2 \implies \mathbf{x}'(\alpha \mathbf{I} - \Phi \Phi') \mathbf{x} = 0, \quad \forall \mathbf{x} \in \mathbb{F}^N \implies \alpha \mathbf{I} - \Phi \Phi' = \mathbf{0} \implies \alpha = \text{eig}\{\Phi \Phi'\} = \{\sigma_k^2\},$$

because the singular values are the square roots of those eigenvalues. □

Normally, finding  $\sigma_1$  requires an SVD, which is expensive for large problems. But for a tight frame,

$$\|\Phi' \mathbf{e}_1\|_2 = \sqrt{\alpha} = \sigma_1$$

so here we can find  $\sigma_1$  by a simple matrix-vector multiplication and a norm.

Furthermore, one can show that  $\Phi' \Phi \preceq \alpha \mathbf{I}$ , i.e.,  $\mathbf{I} - \Phi' \Phi \succeq \mathbf{0}$ , using the SVD  $\Phi = U [\sqrt{\alpha} \mathbf{I} \quad \mathbf{0}] V'$ . (HW)

**Parseval tight frame**

Define. A matrix  $\Phi$  is a **Parseval tight frame** [14, 15] iff  $\Phi$  is a **tight frame** with  $\alpha = \beta = 1$ , i.e.,

$$\|x\|_2^2 = \|\Phi'x\|_2^2 = \sum_{m=1}^M |\langle \phi_m, x \rangle|^2, \quad \forall x \in \mathbb{F}^N. \quad (4.26)$$

What are the singular values  $\sigma_1$  and  $\sigma_N$  of  $\Phi$  in this case?  $\sigma_1 = \cdots = \sigma_N = 1$

**Properties of Parseval tight frames**

If  $\Phi$  is a **Parseval tight frame**, then

- $\Phi\Phi' = I_N$ ,
- its pseudo-inverse is simply  $\Phi^+ = \Phi'$ .

Proof. Using (4.26):

$$\|x\|_2^2 = \|\Phi'x\|_2^2 \implies x'(I - \Phi\Phi')x = 0, \quad \forall x \in \mathbb{F}^N \implies I - \Phi\Phi' = 0.$$

13. The converse also holds, i.e., if  $\Phi\Phi' = I_N$ , then  $\Phi$  is a **Parseval tight frame**. (?)

A: True

B: False

??

14. Every **unitary matrix** is a **tight frame**. (?)

A: True

B: False

??

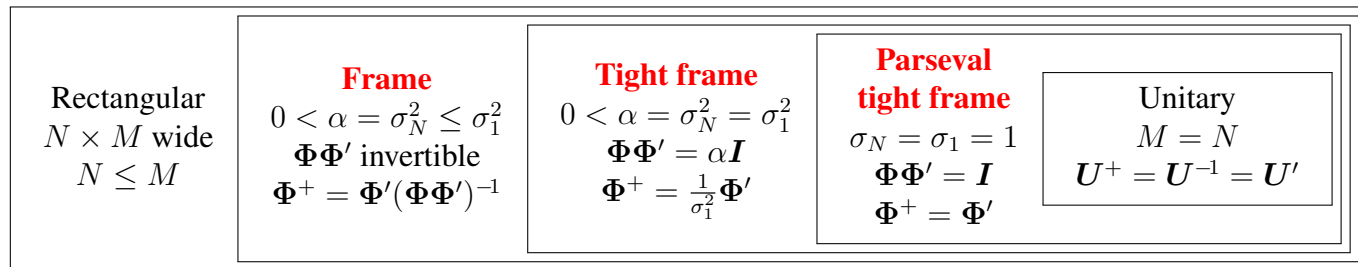
Example. A simple **Parseval tight frame** is the “Mercedes Benz” frame:  $\Phi = \sqrt{\frac{2}{3}} \begin{bmatrix} 0 & -\sqrt{3}/2 & \sqrt{3}/2 \\ 1 & -1/2 & -1/2 \end{bmatrix}$ .

One can verify that  $\Phi\Phi' = \mathbf{I}_2$ .

15. Every **Parseval tight frame** is a **unitary matrix**. (?)  
 A: True B: False

??

### Venn diagram of frames



Another way to construct a **tight frame** is to combine multiple unitary matrices:

$$\Phi = [U_1 \ \dots \ U_K].$$

This construction is fairly common in signal processing, *e.g.*, combining orthogonal wavelet transforms with other transforms.

16.

If  $U_1$  and  $U_2$  are  $N \times N$  unitary matrices, what is the frame bound of  $\Phi \triangleq [U_1 \ U_2]$  ?

A: 1

B: 2

C:  $N$ D:  $2N$ 

E: None

??

- Frames have numerous uses in signal processing [14–17].
- In the usual case where  $\Phi$  is wide, they are considered “robust redundant signal representations” [18].
- Frame theory underpins recent advances in image denoising and attempts to provide insights into CNNs by relating them to perfect reconstruction filter banks [19] [20] [21].
- Solving LLS problems for a **Parseval tight frame**  $\Phi$  is as easy as for a unitary matrix, because  $\Phi^+ = \Phi'$ . So given a vector  $x$  in  $\mathbb{F}^N$ , the representation of  $x$  using the (typically linearly dependent!) set of “atoms”  $\{\phi_1, \dots, \phi_M\}$  for which the coefficient vector has minimum 2-norm is

$$x = \Phi c, \quad c = \Phi' x.$$

In this sense, Parseval tight frames are generalizations of orthonormal bases and unitary matrices.

- Other minimum-norm representations are also of interest, such as [22]

$$\arg \min_{c: x = \Phi c} \|c\|_p.$$

## 4.7 Projection and orthogonal projection

**Idempotent matrix**

L§7.1

Define. A (square) matrix  $\mathbf{P}$  is called a **projection matrix** iff  $\mathbf{P}^2 = \mathbf{P}\mathbf{P} = \mathbf{P}$ .

Such a (square) matrix is also called an **idempotent matrix**.

Example.  $\mathbf{P} = \begin{bmatrix} 1 & a \\ 0 & 0 \end{bmatrix} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^{-1}$ ,  $\mathbf{V} = \begin{bmatrix} 1 & -a/\sqrt{1+a^2} \\ 0 & 1/\sqrt{1+a^2} \end{bmatrix}$ ,  $\mathbf{\Lambda} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$  is idempotent for any  $a$ .

All eigenvalues of an idempotent matrix are either 0 or 1.

Proof. Suppose  $\mathbf{x}$  is an eigenvector of an idempotent matrix  $\mathbf{P}$ , i.e.,  $\mathbf{P}\mathbf{x} = \lambda\mathbf{x}$ . Multiplying both sides by  $\mathbf{P}$  yields  $\mathbf{P}(\mathbf{P}\mathbf{x}) = \lambda(\mathbf{P}\mathbf{x}) \implies \mathbf{P}\mathbf{x} = \lambda^2\mathbf{x}$ . Combining we have  $\lambda^2 = \lambda$ , so  $\lambda = 0$  or  $\lambda = 1$ . □

Fact. Every projection matrix is **diagonalizable** [wiki].

17. The converse of that property also holds, i.e., if  $\mathbf{A}$  is a (square) diagonalizable matrix with eigenvalues that are all either 0 or 1, then  $\mathbf{A}$  is always idempotent. (?)

A: True

B: False

??

??

Example. For any matrix  $\mathbf{A}$ , the matrix  $\mathbf{P} = \mathbf{A}\mathbf{A}^+$  is **idempotent**, because  $\mathbf{P}^2 = (\mathbf{A}\mathbf{A}^+)(\mathbf{A}\mathbf{A}^+) = \mathbf{A}(\mathbf{A}^+\mathbf{A}\mathbf{A}^+) = \mathbf{A}\mathbf{A}^+ = \mathbf{P}$ .



## Orthogonal projection matrix

Our primary interest will be the subset of projection matrices that are (Hermitian) symmetric matrices.

Define. A (square) matrix  $P$  is called an **orthogonal projector** or **orthogonal projection matrix** iff  $P$  is **idempotent** and  $P$  is **Hermitian**.

Caution: An **orthogonal projection matrix** typically is not an **orthogonal matrix**!

If  $P = P' = P^2$  is an **orthogonal projection matrix** and if  $P$  is also an **orthogonal matrix** then  $I = P'P = P^2 = P$ . So the only matrix that is both is the identity matrix  $I$ .



Example.  $P = \frac{1}{2} \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} = \left( \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix} \right) \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \left( \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} \right)$ .

18. Every **orthogonal projection matrix** has a **unitary eigendecomposition**. (?)

A: True

B: False

??

19. If  $x$  has unit norm, then  $xx'$  is an orthogonal projection matrix. (?)

A: True

B: False

??

20. Every orthogonal projection matrix  $P$  is positive semidefinite. (?)

A: True

B: False

??

Alternative explanation: if  $P$  is an orthogonal projection matrix, *i.e.*, both idempotent and Hermitian, then:

$$P =$$

More generally, if  $Q$  is any (possibly non-square) matrix with orthonormal columns, then  $P = QQ'$  is an orthogonal projection matrix, because

- $P = QQ'$  is Hermitian
- $P^2 = (QQ')(QQ') = QQ' = P$ , because  $Q'Q = I$ .

21. Is the converse true? Can a  $N \times N$  **orthogonal projection matrix**  $P$  be written as  $QQ'$  for some matrix  $Q$  with orthonormal columns?

A: Yes, always.

B: Yes, if  $\text{rank}(P) \geq 1$ .

C: Yes, if  $\text{rank}(P) \leq N - 1$ .

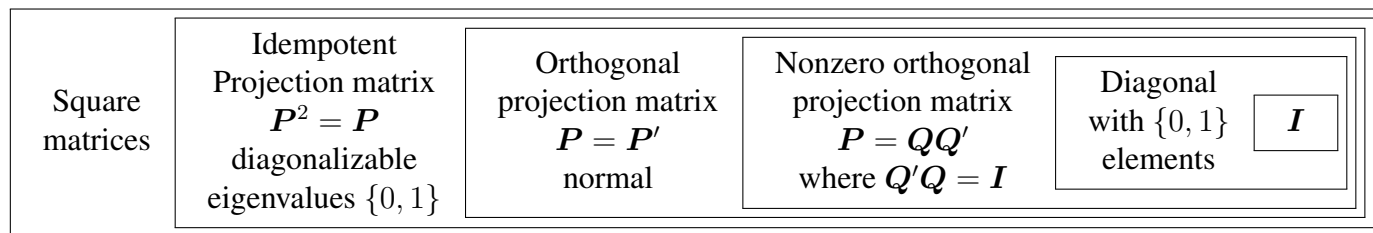
D: Only if  $P$  is nonsingular.

E: No.

??

22. When writing a nonzero **orthogonal projection matrix**  $P$  as  $QQ'$ , the  $Q$  is unique.  
 A: True B: False ??
23. The  $N \times N$  **Gram matrix**  $G = \Phi'\Phi$  of a  $N \times M$  **Parseval tight frame**  $\Phi$  is an **orthogonal projection matrix**. (?)  
 A: True B: False ??

The following **Venn diagram** summarizes relationships related to **projection** matrices.



## Convex sets

---

Define. A nonempty set  $\mathcal{C}$  in a vector space  $\mathcal{V}$  is a **convex set** iff

$$\mathbf{x}, \mathbf{z} \in \mathcal{C} \implies \alpha \mathbf{x} + (1 - \alpha) \mathbf{z} \in \mathcal{C}, \quad \forall 0 \leq \alpha \leq 1.$$

A **linear combination**  $\alpha \mathbf{x} + (1 - \alpha) \mathbf{z}$  for any  $0 \leq \alpha \leq 1$  is a **convex combination**.

It follows by induction that if  $\mathbf{x}_1, \dots, \mathbf{x}_K \in \mathcal{C}$  and  $\alpha_1, \dots, \alpha_K \geq 0$  and  $\sum_{k=1}^K \alpha_k = 1$ , then  $\sum_{k=1}^K \alpha_k \mathbf{x}_k \in \mathcal{C}$ . Such a sum is also called a **convex combination**.

Example. Any **subspace**  $\mathcal{S}$  in a vector space  $\mathcal{V}$  is a **convex set**.

If  $\mathcal{V}$  is a finite-dimensional vector space, then any subspace is a **closed** convex set.

Example. The **nonnegative orthant**  $\mathbb{R}_+^N \triangleq \{\mathbf{x} \in \mathbb{R}^N : \mathbf{x} \geq \mathbf{0}\}$  is a convex set.

Example. The ball of radius  $r > 0$  in  $\mathbb{F}^N$ , defined as  $\mathcal{B}_r \triangleq \{\mathbf{x} \in \mathbb{F}^N : \|\mathbf{x}\|_2 \leq r\}$ , is a convex set.

(HW)

## Projection onto convex sets

---

Finding the nearest point in a set is an important operation.

Ch. 3 defined the **projection** of a point  $\mathbf{v}$  onto a set  $\mathcal{S}$  as

$$\mathcal{P}_{\mathcal{S}}(\mathbf{v}) \triangleq \arg \min_{\mathbf{s} \in \mathcal{S}} \|\mathbf{v} - \mathbf{s}\|.$$

For general sets  $\mathcal{S}$ , the “arg min” in this definition need not be unique.

However, if  $\mathcal{C}$  is a **closed convex set** in a finite-dimensional vector space  $\mathcal{V}$ ,

then for any  $\mathbf{v} \in \mathcal{V}$  the nearest point in  $\mathcal{C}$  exists and is unique.

(See EECS 600 or IOE 611.) Thus the “arg min” in the definition of  $\mathcal{P}_{\mathcal{C}}$  above is well-defined.

---

We have used the term **projection** in two ways here.

- One use is in defining a  $N \times N$  **projection matrix**  $\mathbf{P}$  that has the property that  $\mathbf{P}^2 = \mathbf{P}$ .  
In other words, for any  $\mathbf{x} \in \mathbb{F}^N$ :  $\mathbf{P}(\mathbf{P}\mathbf{x}) = \mathbf{P}^2\mathbf{x} = \mathbf{P}\mathbf{x}$ .
- Another use is in defining the **projection** onto a (typically convex) set  $\mathcal{P}_{\mathcal{C}}(\cdot)$ .

One connection between the two uses is that if we project a point  $\mathbf{v}$  onto a set  $\mathcal{C}$  to get  $\hat{\mathbf{v}} = \mathcal{P}_{\mathcal{C}}(\mathbf{v})$  and then project  $\hat{\mathbf{v}}$  onto  $\mathcal{C}$  we get the same point:  $\hat{\mathbf{v}} = \mathcal{P}_{\mathcal{C}}(\hat{\mathbf{v}})$ .

In other words, for any  $\mathbf{v} \in \mathcal{V}$ :  $\mathcal{P}_{\mathcal{C}}(\mathcal{P}_{\mathcal{C}}(\mathbf{v})) = \mathcal{P}_{\mathcal{C}}(\mathbf{v})$ .

Put concisely, we have an expression very similar to  $\mathbf{P}^2 = \mathbf{P}$ , hence the similarity in terminology:

$$\mathcal{P}_{\mathcal{C}} \circ \mathcal{P}_{\mathcal{C}} = \mathcal{P}_{\mathcal{C}}.$$

---

## Projection onto a subspace

Now we move towards using projection as a tool for solving signal processing problems.

Let  $\mathbf{a}_1, \dots, \mathbf{a}_N$  denote any collection of  $N$  vectors, each in  $\mathbb{F}^M$ . The span of this set defines a subspace:

$$\mathcal{S} \triangleq \text{span}(\{\mathbf{a}_1, \dots, \mathbf{a}_N\}).$$

Equivalently (because we are working in  $\mathbb{F}^M$  here) we can group the vectors into a  $M \times N$  matrix and then describe the subspace as the **range** of this matrix:

$$\mathcal{S} = \mathcal{R}(\mathbf{A}) \triangleq \{\mathbf{A}\mathbf{x} : \mathbf{x} \in \mathbb{F}^N\}, \quad \mathbf{A} \triangleq [\mathbf{a}_1 \quad \dots \quad \mathbf{a}_N].$$

As introduced in Ch. 3, in many applications (including handwritten digit recognition via nearest subspace) we are given some test vector  $\mathbf{y}$  and we want to find the vector in a subspace  $\mathcal{S}$  that is closest to  $\mathbf{y}$ :

$$\hat{\mathbf{y}} = \arg \min_{\mathbf{s} \in \mathcal{S}=\mathcal{R}(\mathbf{A})} \|\mathbf{y} - \mathbf{s}\|_2.$$

The resulting vector  $\hat{\mathbf{y}}$  is called the **projection** of the point  $\mathbf{y}$  onto the subspace  $\mathcal{S}$ , and the process (operation) of finding that closest point is called **projecting** the point  $\mathbf{y}$  onto the subspace  $\mathcal{S}$ , as introduced in Ch. 3.

The key to solving this problem is to use the fact that every point in the subspace  $\mathcal{S}$  has the form  $\mathbf{A}\mathbf{x}$  for some  $\mathbf{x} \in \mathbb{F}^N$ . Thus  $\hat{\mathbf{y}} = \mathbf{A}\hat{\mathbf{x}}$  for some  $\hat{\mathbf{x}} \in \mathbb{F}^N$ , and we just have to find that best  $\hat{\mathbf{x}}$  to get  $\hat{\mathbf{y}}$ . In other words, the closest point or projection problem is equivalent to:

$$\hat{\mathbf{y}} =$$

because every  $\mathbf{s} \in \mathcal{S}$  is  $\mathbf{s} = \mathbf{A}\mathbf{x}$  for some  $\mathbf{x} \in \mathbb{F}^N$ .

This form involves solving a LS problem, and a solution to that part is  $\hat{\mathbf{x}} = \mathbf{A}^+\mathbf{y}$ , so

$$\hat{\mathbf{y}} = \arg \min_{\mathbf{s} \in \mathcal{S}=\mathcal{R}(\mathbf{A})} \|\mathbf{y} - \mathbf{s}\|_2 =$$

In other words:

$$\mathcal{P}_{\mathcal{R}(\mathbf{A})}(\mathbf{y}) = \mathbf{A}\mathbf{A}^+\mathbf{y}.$$

24. If  $\mathbf{A}$  does not have full rank, then there will be other LS solutions, *i.e.*,  $\hat{\mathbf{x}}$  is not unique. Thus,  $\hat{\mathbf{y}}$  is not unique. (?)

A: True

B: False


??

## Practical implementation

---

Reiterating, the point on the subspace  $\mathcal{R}(\mathbf{A})$  that is closest to the point  $\mathbf{y}$  is  $\hat{\mathbf{y}} = \mathbf{A}\mathbf{A}^+\mathbf{y}$

If we need to compute this just once, for one  $\mathbf{A}$  and one  $\mathbf{y}$ , then using code like `A * (A \ y)` is fine if  $\mathbf{A}$  is small enough for backslash to work. If  $\mathbf{A}$  is large, then we use an iterative method to compute the LS coefficients  $\hat{\mathbf{x}}$  first, then compute the **projection**  $\hat{\mathbf{y}} = \mathbf{A}\hat{\mathbf{x}}$ .

The parentheses are essential in `A * (A \ y)` because `*` and `\` have the same **operator precedence** in JULIA. 

But in most applications  $\mathbf{A}$  is fixed (after some training or modeling process) and we will need to perform projection for many different “test”  $\mathbf{y}$  vectors. In such cases, it is more efficient to use an SVD at the beginning (as part of the training process) to save computation at the test stage. Specifically, when  $r > 0$ :

$$\mathbf{P}_{\mathcal{R}(\mathbf{A})} \triangleq \mathbf{A}\mathbf{A}^+ = \mathbf{U}_r \boldsymbol{\Sigma}_r \mathbf{V}_r' \mathbf{V}_r \boldsymbol{\Sigma}_r^{-1} \mathbf{U}_r' = \mathbf{U}_r \mathbf{U}_r'.$$

So a practical implementation is something like the following JULIA code for  $r > 0$ :

```
(U, s, V) = svd(A)
r = sum(s .> threshold)
Ur = U[:, 1:r]
projector = (y) -> Ur * (Ur' * y)
```

We do the SVD once, then after that we need only use simple matrix-vector multiplies to perform projection.

Note that  $\mathbf{U}_r$  is an orthonormal basis for  $\mathcal{R}(\mathbf{A})$  and  $\mathbf{P}_{\mathcal{R}(\mathbf{A})} = \mathbf{U}_r \mathbf{U}_r'$ . This is not a coincidence!



**Orthonormal vs non-orthonormal bases for a subspace** \_\_\_\_\_ (Read)

More generally, if  $Q$  is a matrix whose columns are orthonormal, then those columns form an orthonormal basis for a subspace, namely  $\mathcal{R}(Q)$ , and the projection of a vector  $y$  onto that subspace is simply

$$\mathcal{P}_{\mathcal{R}(Q)}(y) = P_{\mathcal{R}(Q)}y = Q(Q'y).$$

So **orthonormal bases** for a subspace are extremely convenient for computation, because subspace projection using such bases requires mere matrix-vector multiplication.

If a basis  $B$  is not orthonormal, then to compute the projection we would need

$$P_{\mathcal{R}(B)}y = B(B^+y),$$

as derived above, which is much more expensive in general because  $B^+$  usually requires an SVD.

**Orthogonality principle revisited**

(Read)

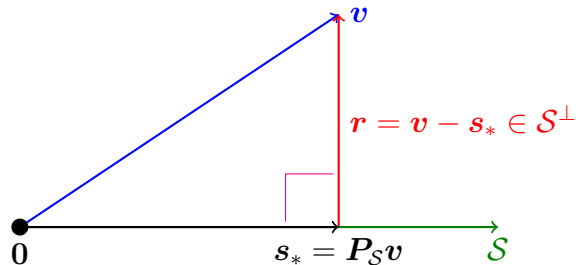
Another version of the **orthogonality principle** is the following. Let  $\mathcal{S}$  denote any subspace of a vector space  $\mathcal{V}$ , and  $\mathbf{v} \in \mathcal{V}$  be any point in that vector space. If the closest point in  $\mathcal{S}$  to  $\mathbf{v}$  is

$$\mathbf{s}_* = \arg \min_{\mathbf{s} \in \mathcal{S}} \|\mathbf{v} - \mathbf{s}\|_2 = \mathbf{P}_{\mathcal{S}} \mathbf{v},$$

then the **orthogonality principle** says that the residual vector  $\mathbf{r} = \mathbf{v} - \mathbf{s}_*$  is perpendicular to the entire subspace  $\mathcal{S}$ , i.e.,

$$\langle \mathbf{v} - \mathbf{s}_*, \mathbf{s} \rangle = 0, \quad \forall \mathbf{s} \in \mathcal{S}.$$

The following figure illustrates this principle.



Note that  $\mathbf{r} = \mathbf{v} - \mathbf{s}_* = \mathbf{v} - \mathbf{P}_{\mathcal{S}} \mathbf{v} = (\mathbf{I} - \mathbf{P}_{\mathcal{S}}) \mathbf{v} = \mathbf{P}_{\mathcal{S}}^{\perp} \mathbf{v}$ ,  $\mathbf{P}_{\mathcal{S}}^{\perp} \triangleq \mathbf{I} - \mathbf{P}_{\mathcal{S}}$ .

Proof sketch from [wiki]: For  $\mathbf{s} \in \mathcal{S}$  and any  $\alpha \in \mathbb{F}$ :  $0 \leq \|(\mathbf{s}_* + \alpha \mathbf{s}) - \mathbf{v}\|_2^2 - \|\mathbf{s}_* - \mathbf{v}\|_2^2$   
 $= 2 \operatorname{real}\{\alpha \langle \mathbf{s}_* - \mathbf{v}, \mathbf{s} \rangle\} + |\alpha|^2 \|\mathbf{s}\|_2^2 \implies \langle \mathbf{s}_* - \mathbf{v}, \mathbf{s} \rangle = 0.$

## Projection onto a subspace's orthogonal complement

---

Recall that in a finite-dimensional vector space  $\mathcal{V}$ , if  $\mathcal{S}$  is a subspace of  $\mathcal{V}$ , then

$$\mathcal{V} =$$

Clearly  $P_{\mathcal{V}} = I$ , so it follows that

$$I =$$

Thus the projection onto the subspace's **orthogonal complement**  $\mathcal{S}^\perp$  is simply:

$$P_{\mathcal{S}^\perp} =$$

If  $P$  is any projection matrix, then as a notation convention we define

$$P^\perp \triangleq$$

It then follows that

$$P_{\mathcal{S}^\perp} =$$

In words,  $P_{\mathcal{S}^\perp}$  is the orthogonal projector onto  $\mathcal{S}^\perp$ , the orthogonal complement of the subspace  $\mathcal{S}$ .

Example. If  $\mathcal{S} = \text{span}(\{\mathbf{1}_n\})$  then  $P_{\mathcal{S}}^\perp = I - \mathbf{1}_n \mathbf{1}_n^+ = I - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n'$ .

The matrix-vector product  $\mathbf{y} = P_{\mathcal{S}}^\perp \mathbf{x}$  returns a vector  $\mathbf{y}$  having zero mean, so  $\mathbf{y}$  is orthogonal to  $\mathbf{1}_n$ .

## Projectors and the four fundamental subspaces

Recall the SVD anatomy of a  $M \times N$  matrix with rank  $0 < r < \min(M, N)$ :

L§7.1.1

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}' = \sum_{k=1}^{\min(M,N)} \sigma_k \mathbf{u}_k \mathbf{v}_k' = \sum_{k=1}^r \sigma_k \mathbf{u}_k \mathbf{v}_k' = \left[ \begin{array}{c|c} \mathbf{U}_r & \mathbf{U}_0 \end{array} \right] \left[ \begin{array}{c|c} \mathbf{\Sigma}_r & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{0} \end{array} \right] \left[ \begin{array}{c|c} \mathbf{V}_r & \mathbf{V}_0 \end{array} \right]'$$

The four fundamental subspaces have the following **orthonormal bases** and **orthogonal projectors**.

space in terms of $\mathbf{A}$	subspace in terms of $\mathbf{A}'$	equivalent subspace matrix	orthonormal basis	orthogonal projector
$\mathbb{F}^N$	$\mathcal{N}(\mathbf{A})$	$\mathcal{R}^\perp(\mathbf{A}')$	$\mathbf{V}_0$	$\mathbf{V}_0 \mathbf{V}_0' = \mathbf{I} - \mathbf{V}_r \mathbf{V}_r'$
$\mathbb{F}^N$	$\mathcal{N}^\perp(\mathbf{A})$	$\mathcal{R}(\mathbf{A}')$	$\mathbf{V}_r$	$\mathbf{V}_r \mathbf{V}_r'$
$\mathbb{F}^M$	$\mathcal{R}(\mathbf{A})$	$\mathcal{N}^\perp(\mathbf{A}')$	$\mathbf{U}_r$	$\mathbf{U}_r \mathbf{U}_r'$
$\mathbb{F}^M$	$\mathcal{R}^\perp(\mathbf{A})$	$\mathcal{N}(\mathbf{A}')$	$\mathbf{U}_0$	$\mathbf{U}_0 \mathbf{U}_0' = \mathbf{I} - \mathbf{U}_r \mathbf{U}_r'$

I list two forms for the projectors involving  $\mathbf{U}_0$  and  $\mathbf{V}_0$  because if we have stored “only” the **compact SVD**, then we will need to use  $\mathbf{U}_r$  and  $\mathbf{V}_r$ .

25.

$$\mathcal{N}^\perp(\mathbf{A}') \oplus \mathcal{R}^\perp(\mathbf{A}) = \mathbb{F}^M. (?)$$

A: True

B: False

??

### Warm-up questions: projections

26. If  $Q$  is a matrix with orthonormal columns and  $D$  is diagonal matrix whose elements are all 0 or 1, then  $P = QDQ'$  is an **orthogonal projection matrix**. (?)  
 A: True B: False ??

27. Let rank  $r$  matrix  $A$  have **SVD**  $A = U\Sigma V'$  and **compact SVD**  $A = U_r \Sigma_r V_r'$ , where we partition unitary matrix  $V$  as usual as  $V = [V_r \ V_0]$ . Which of the following matrices is  $P_{\mathcal{N}(A'A)}$ ?  
 A:  $V_0 V_0'$  B:  $V_r V_r'$  C:  $V_0 V_r'$  D:  $V_r V_0'$  E: None of these. ??

28. Continuing the previous problem, and assuming the rank  $r$  is much smaller than the matrix dimensions, which of the following JULIA snippets is the most efficient way to implement  $P_{\mathcal{R}^\perp(A^+)}x$  after the line `U,s,V = svd(A)` ?  
 A: `V[:,1:r] * (V[:,1:r]' * x)` B: `V[:,1:r] * V[:,1:r]' * x`  
 C: `x - V[:,1:r] * V[:,1:r]' * x` D: `x - V[:,1:r] * (V[:,1:r]' * x)`  
 E: `V[:,(r+1):end] * (V[:,(r+1):end]' * x)`

??

This last question is a “mini review” that uses elements from Ch. 1 (multiplication order), Ch. 2 (SVD), Ch. 3 (subspaces), Ch. 4 (pseudo-inverse, orthogonal projection), and JULIA!

## Binary classifier design using least-squares

(Read)

Linear LS can be used as a simple tool for designing a binary classifier via supervised learning.

Given  $M$  feature vectors  $\mathbf{v}_i \in \mathbb{F}^N$  and corresponding binary class labels  $y_i = \pm 1$ . We want to find a weight vector  $\mathbf{x}$  such that  $\langle \mathbf{x}, \mathbf{v}_i \rangle$  has the same sign as  $y_i$ . Equivalently, we want  $\text{sign}(\langle \mathbf{x}, y_i \mathbf{v}_i \rangle) = 1$ .

Form a  $M \times N$  data matrix from the training feature vectors and a corresponding label vector:

$$\tilde{\mathbf{A}} = \begin{bmatrix} \mathbf{v}'_1 \\ \vdots \\ \mathbf{v}'_M \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix}.$$

Then solve a least-squares problem (or regularized variant thereof):

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x} \in \mathbb{R}^N} \left\| \tilde{\mathbf{A}} \mathbf{x} - \mathbf{y} \right\|_2.$$

Because each element of  $\mathbf{y}$  is  $\pm 1$ , the diagonal matrix  $\mathbf{D} = \text{Diag}\{\mathbf{y}\}$  is unitary and  $\mathbf{D}^{-1} = \mathbf{D}' = \mathbf{D}$ . Using unitary invariance of the 2-norm, and the fact that  $\mathbf{y} = \mathbf{D}\mathbf{1}$ , an equivalent formulation is

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x} \in \mathbb{R}^N} \left\| \mathbf{D} \tilde{\mathbf{A}} \mathbf{x} - \mathbf{1}_M \right\|_2 = \arg \min_{\mathbf{x} \in \mathbb{R}^N} \left\| \mathbf{A} \mathbf{x} - \mathbf{1}_M \right\|_2, \quad \mathbf{A} \triangleq \begin{bmatrix} y_1 \mathbf{v}'_1 \\ \vdots \\ y_M \mathbf{v}'_M \end{bmatrix}.$$

After “learning” the regression weights  $\hat{\mathbf{x}}$ , the classifier for a subsequent test data point  $\mathbf{v} \in \mathbb{R}^N$  is simply

$$\text{sign}(\mathbf{v}'\hat{\mathbf{x}}).$$

It is somewhat unconventional to use LS for a data vector  $\mathbf{y}$  that is binary ( $\pm 1$  values), whereas **logistic regression** (see Ch. 8) is truly designed for such data. Nevertheless, the LS approach is simple and fast and can work adequately in some cases. This approach will be explored in Discussion.

### 4.8 Summary

Key points about solving LS problems  $\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2$

- Every minimizer  $\hat{\mathbf{x}}$  satisfies the **normal equations**  $\mathbf{A}'\mathbf{A}\hat{\mathbf{x}} = \mathbf{A}'\mathbf{y}$
- If  $\mathbf{A}$  has full column rank then  $\mathbf{A}'\mathbf{A}$  is invertible and the unique solution is  $\hat{\mathbf{x}} = (\mathbf{A}'\mathbf{A})^{-1}\mathbf{A}'\mathbf{y}$
- The **minimum-norm LS solution** is always  $\hat{\mathbf{x}} = \mathbf{A}^+\mathbf{y}$
- If  $\mathbf{A}$  is **unitary** or a **Parseval tight frame**, then  $\mathbf{A}^+ = \mathbf{A}'$
- The projection of a point  $\mathbf{y}$  onto  $\mathcal{R}(\mathbf{A})$  is  $\mathbf{A}\mathbf{A}^+\mathbf{y} = \mathbf{U}_r\mathbf{U}_r'\mathbf{y}$  when  $r > 0$  and is  $\mathbf{0}$  when  $r = 0$ .

### Bibliography

- [1] A. J. Laub. *Matrix analysis for scientists and engineers*. Soc. Indust. Appl. Math., 2005 (cit. on pp. 4.2, 4.3).
- [2] T. Adali, P. J. Schreier, and L. L. Scharf. “Complex-valued signal processing: the proper way to deal with impropriety”. In: *IEEE Trans. Sig. Proc.* 59.11 (Nov. 2011), 5101–25 (cit. on p. 4.14).
- [3] L. Sorber, M. V. Barel, and L. D. Lathauwer. “Unconstrained optimization of real functions in complex variables”. In: *SIAM J. Optim.* 22.3 (2012), 879–98 (cit. on p. 4.14).
- [4] J. Liu, C. Garcia-Cardona, B. Wohlberg, and W. Yin. “First- and second-order methods for online convolutional dictionary learning”. In: *SIAM J. Imaging Sci.* 11.2 (Jan. 2018), 1589–628 (cit. on p. 4.14).
- [5] L. Eldén. *Matrix methods in data mining and pattern recognition*. Errata: <http://users.mai.liu.se/larel04/matrix-methods/index.html> <http://users.mai.liu.se/larel04/matrix-methods/>. Soc. Indust. Appl. Math., 2007 (cit. on pp. 4.41, 4.44).
- [6] L. Mor-Yosef and H. Avron. “Sketching for principal component regression”. In: *SIAM J. Matrix. Anal. Appl.* 40.2 (Jan. 2019), 454–85 (cit. on p. 4.43).
- [7] Per-AAke Wedin. “Perturbation theory for pseudo-inverses”. In: *BIT Numerical Mathematics* 13.2 (June 1973), 217–32 (cit. on p. 4.44).
- [8] G. W. Stewart. “On the perturbation of pseudo-inverses, projections and linear least squares problems”. In: *SIAM Review* 19.4 (1977), 634–62 (cit. on p. 4.44).



- [9] G. W. Stewart. “Stochastic perturbation theory”. In: *SIAM Review* 32.4 (1990), 579–610 (cit. on p. 4.44).
- [10] L. Meng and B. Zheng. “The optimal perturbation bounds of the Moore-Penrose inverse under the Frobenius norm”. In: *Linear Algebra and its Applications* 432.4 (Feb. 2010), 956–63 (cit. on p. 4.44).
- [11] L. Meng and B. Zheng. “New multiplicative perturbation bounds of the Moore-Penrose inverse”. In: *Linear and Multilinear Algebra* 63.5 (2015), 1037–48 (cit. on p. 4.44).
- [12] Per-AAke Wedin. “On the almost rank deficient case of the least squares problem”. In: *BIT Numerical Mathematics* 13.3 (Sept. 1973), 344–54.
- [13] S. F. D. Waldron. *An introduction to finite tight frames*. Springer, 2018 (cit. on p. 4.50).
- [14] J. Kovacevic and A. Chebira. “Life beyond bases: the advent of frames (Part I)”. In: *IEEE Sig. Proc. Mag.* 24.4 (July 2007), 86–104 (cit. on pp. 4.52, 4.53, 4.55).
- [15] J. Kovacevic and A. Chebira. “Life beyond bases: the advent of frames (Part II)”. In: *IEEE Sig. Proc. Mag.* 24.5 (Sept. 2007), 115–25 (cit. on pp. 4.52, 4.53, 4.55).
- [16] B. Adcock and D. Huybrechs. “Frames and numerical approximation”. In: *SIAM Review* 61.3 (2019), 443–73 (cit. on pp. 4.52, 4.55).
- [17] J. Kovacevic and A. Chebira. “An introduction to frames”. In: *Found. & Trends in Sig. Pro.* 2.1 (2008), 1–94 (cit. on p. 4.55).
- [18] A. Rahimi, G. Zandi, and B. Daraby. “Redundancy of fusion frames in Hilbert spaces”. In: *Complex Analysis and Operator Theory* 10.3 (Mar. 2016), 545–65 (cit. on p. 4.55).
- [19] R. Yin, T. Gao, Y. M. Lu, and I. Daubechies. “A tale of two bases: local-nonlocal regularization on image patches with convolution framelets”. In: *SIAM J. Imaging Sci.* 10.2 (Jan. 2017), 711–50 (cit. on p. 4.55).
- [20] E. Kang and J. C. Ye. “Framelet denoising for low-dose CT using deep learning”. In: *Proc. IEEE Intl. Symp. Biomed. Imag.* 2018, 311–4 (cit. on p. 4.55).
- [21] J. C. Ye, Y. Han, and E. Cha. “Deep convolutional framelets: A general deep learning framework for inverse problems”. In: *SIAM J. Imaging Sci.* 11.2 (Jan. 2018), 991–1048 (cit. on p. 4.55).
- [22] M. Akcakaya and V. Tarokh. “A frame construction and a universal distortion bound for sparse representations”. In: *IEEE Trans. Sig. Proc.* 56.6 (June 2008), 2443–50 (cit. on p. 4.55).