

Maths for Signals and Systems

Session 3 [Week 4-5]

Pier Luigi Dragotti

EEE Department

Imperial College London

p.dragotti@imperial.ac.uk

Outline and Motivation

- The goal of this group of lectures is to solve $Ax = b$
 - Specifically, given $b \in \mathbb{C}^m$ and $A \in \mathbb{C}^{m \times n}$, we want to understand when a solution x exists and if yes when it is unique.
 - If there is no solution we want to pick the x that *best approximate* $Ax = b$
- We want to answer these questions by checking the properties of A (range and null space) and by using the geometry of vector spaces and linear mapping
- The more interesting cases in engineering are when $m \neq n$ and when there is no solution or the solution is not unique...

Motivation

- Estimating x from $Ax = b$ is an *inverse problem* which appears in many engineering applications
- Assume that the equation models a filtering operation, then the goal is to remove the distortion due to the filter. For example,
 - you take a picture and (motion) blur has distorted it and you want to remove the blur
 - b is some audio recorded and A models the reflections due to the room. Estimating x is like removing the echo experienced during the recording

$x \rightarrow$



$b = Ax \rightarrow$



$Ax=b$ – Geometrical Interpretation

- We can think of $Ax = b$ in pure maths and geometrical terms as a system of linear equations
 - if $m = n$, we have n equations in n unknowns
 - e.g. $\begin{pmatrix} 1 & 2 \\ 1 & -1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 4 \\ 1 \end{pmatrix}$ each equation is a line, two non-parallel lines intersect in one point which is the **unique** solution,
 - If the lines are parallel we either have no-solution or an infinite number of solutions (two overlapping lines). Clearly $rank(A)$ tells us when the solution is **unique**
 - if $m > n$, we have more equations than unknowns (A is 'tall'), more than two lines normally (but not necessarily) intersect in more than one point unless several lines overlap
 - if $m < n$, we have more unknowns than equation (A is 'fat'), e.g., when $m = 2, n = 3$, we have two planes in 3-D and unless they are parallel they intersect along a line (infinite number of solutions). So the solution, if it exists, is not unique.

$Ax=b$ – ‘Vector Space’ Interpretation

- We can think of the matrix A as describing a linear mapping from \mathbb{C}^n to \mathbb{C}^m
- We can then think of Ax as computing linear combinations of columns of A ,
- **Thus a solution to $Ax = b$ exists if b is in the range space of A**
- Any vector in \mathbb{C}^m can be reached by A when $\mathcal{R}(A) = \mathbb{C}^m$ that is when A is full row rank
- While the range space characterizes existence, the null space characterizes uniqueness.
- **When the null space of A is non-trivial, if a solution to $Ax = b$ exists, the solution is not unique.**

$Ax=b$ - Existence

- We can summarize this discussion with the following theorem:
- **Theorem (existence):** Let $A \in \mathbb{C}^{m \times n}$ then the following statements are equivalent:
 1. For each $b \in \mathbb{C}^m$ there exists at least one solution to $Ax = b$
 2. The range space of A is full: $\mathcal{R}(A) = \mathbb{C}^m$
 3. $\text{rank}(A) = m$
 4. The row of A are linearly independent: A is full row rank
- Note that each of the above conditions implies that $m \leq n$ ('fat/square' matrix)

Existence (examples)

1. $A = \begin{pmatrix} 4 & 8 & 6 \\ 0 & 4 & 7 \end{pmatrix}$. The rows are linearly independent so for every $\mathbf{b} \in \mathbb{C}^2$ there is at **least** one solution.

2. $A = \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 1 & 1 & 2 \end{pmatrix}$. The matrix has three rows but only rank 2. The existence theorem is not satisfied and \mathbf{b} exist for which $A\mathbf{x} = \mathbf{b}$ has no solution. For example, $\mathbf{b} = (0 \ 0 \ 1)^T$

3. $A = \begin{pmatrix} 6 & 5 & 6 & 3 \\ 2 & 3 & 0 & 1 \\ 4 & 4 & 3 & 2 \end{pmatrix}$. Although this matrix has more columns than rows **it is not full row rank**. Hence, existence theorem is not satisfied and a solution to $A\mathbf{x} = \mathbf{b}$ **does not always exist!** So remember 'more unknowns than equations' or 'underdetermined' does not imply existence. 'Full row rank' does.

$Ax=b$ – ‘Vector Space’ Interpretation

- Thus a solution to $Ax = b$ exists if b is in the range space of A
- While the range space characterizes existence, the null space characterizes uniqueness.
- When the null space of A is non-trivial, if a solution to $Ax = b$ exists, the solution is not unique.

$Ax=b$ - Uniqueness

- **Theorem (uniqueness):** Let $A \in \mathbb{C}^{m \times n}$ then the following statements are equivalent:
 1. If there exist a solution to $Ax = b$, it is unique
 2. The null space of A is trivial
 3. $\text{rank}(A) = n$
 4. The columns of A are linearly independent: A is full column rank
- Note that each of the above conditions implies that $m \geq n$ ('tall/square' matrix)

Uniqueness (examples)

- $\mathbf{A} = \begin{pmatrix} 0 & 1 \\ 1 & 1 \\ 1 & 0 \end{pmatrix}$. The columns are linearly independent, if a solution exists it is unique.
- For example, if $\mathbf{b} = [4, 6, 2]^T$, the unique solution is $\mathbf{x} = [2, 4]^T$.
- However, if $\mathbf{b} = [1, 0, 0]^T$ no solution exists.

Uniqueness (examples)

- $\mathbf{A} = \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 1 & 1 & 2 \end{pmatrix}$. The columns are linearly dependent. Null space = $\text{span}\{(1 \ 1 \ -1)^T\}$. So uniqueness not satisfied nor existence.

- $\mathbf{A} = \begin{pmatrix} 2 & 5 & 3 \\ 3 & 5 & 2 \\ 3 & 7 & 4 \\ 7 & 8 & 1 \end{pmatrix}$. The columns are linearly dependent so despite being overdetermined uniqueness cannot be guaranteed

$Ax=b$ – Existence and Uniqueness

- **Theorem (existence and uniqueness):** Let $A \in \mathbb{C}^{m \times n}$ then the following statements are equivalent:
 1. For each $b \in \mathbb{C}^m$ there exists a solution to $Ax = b$ and it is unique
 2. The null space of A is trivial and dimension of range space is m
 3. $\text{rank}(A) = n = m$
 4. A is square and nonsingular
- The case of A being square and nonsingular is of less interest, so we now focus on finding x given b and A for the other cases.

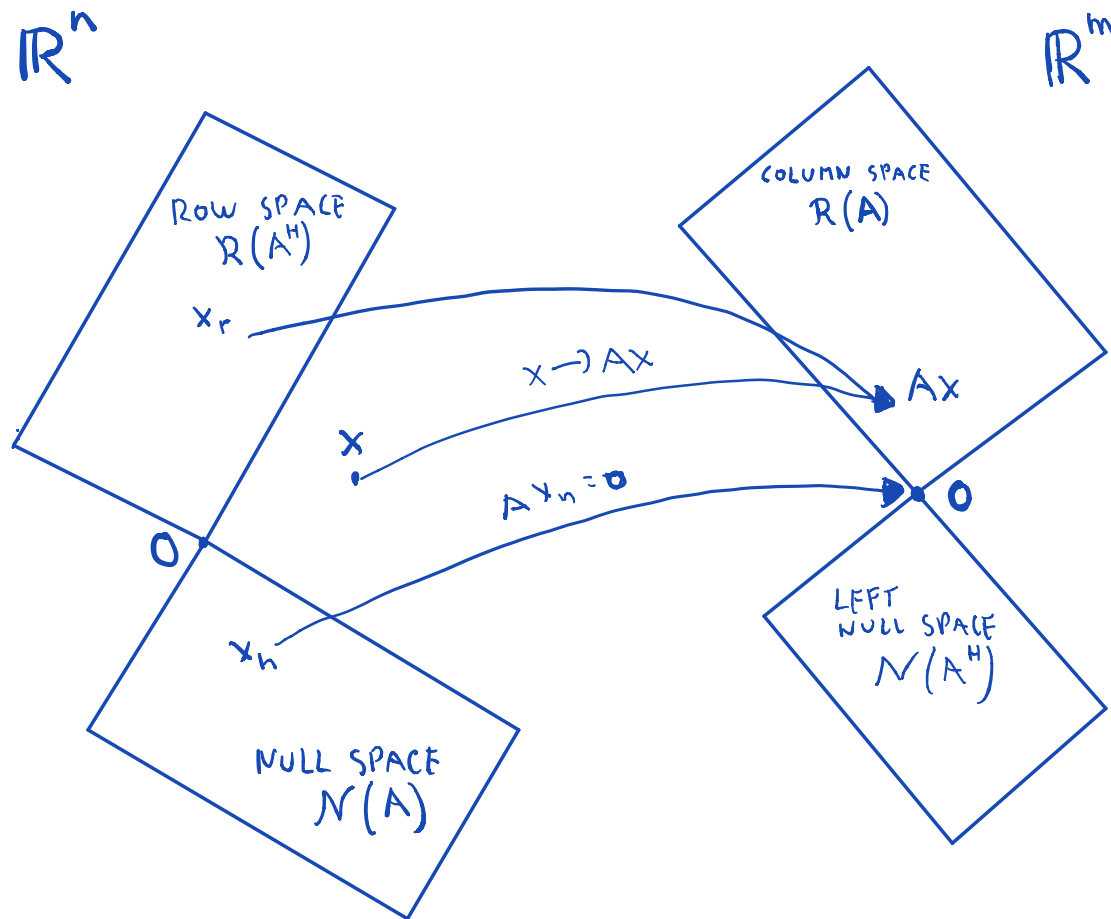
Solutions to $Ax=b$

- Recall that we are given A and b and we want to find x
- We begin by focusing on the full column rank case
- This means that A is tall (or square) and that if the solution exists it is unique
- So our goal is to develop a single method that either finds the x which is the solution to $Ax = b$ if it exists or finds the '*best approximation*' to $Ax = b$
- For *best approximation* we mean the x that minimizes $\|Ax - b\|_2$ (least-squares approximation)
- In order to achieve that we need to introduce a few new notions

Four Fundamental Subspaces

- Given a linear mapping described by $A \in \mathbb{C}^{m \times n}$, we have previously introduced the notion of range of A , $\mathcal{R}(A)$, and null space of A , $\mathcal{N}(A)$
- Consider now the linear mapping described by A^H , this maps vectors from \mathbb{C}^m to \mathbb{C}^n (however, remember that A^H is normally **not** the inverse of A)
- We then have $\mathcal{R}(A^H)$ and $\mathcal{N}(A^H)$
- These four subspaces are related as follows (without proof):
 - $\mathbb{C}^m = \mathcal{R}(A) \oplus \mathcal{N}(A^H)$, $\mathcal{R}^\perp(A) = \mathcal{N}(A^H)$, $\mathcal{N}^\perp(A^H) = \mathcal{R}(A)$,
 - $\mathbb{C}^n = \mathcal{R}(A^H) \oplus \mathcal{N}(A)$, $\mathcal{R}^\perp(A^H) = \mathcal{N}(A)$, $\mathcal{N}^\perp(A) = \mathcal{R}(A^H)$,

Four Fundamental Subspaces



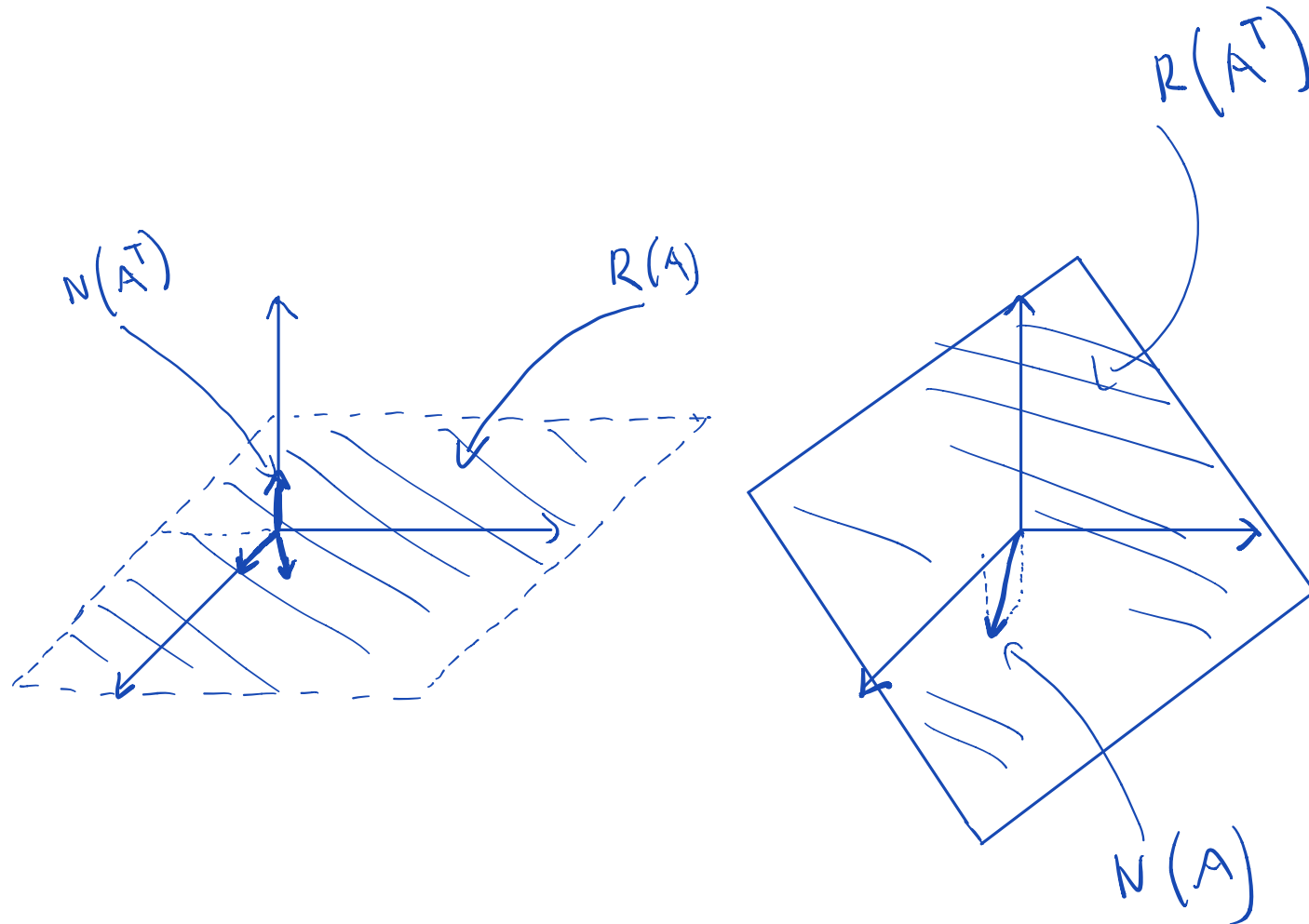
Four Fundamental Subspaces – Worked Example

- **Question:** Find the four fundamental subspaces of $A = \begin{pmatrix} 1 & 1 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix}$
- We first should row-reduce A to find the dimension of the range space and null space. A is already in a reduced form and clearly the range space has dimension 2 and a possible basis is given by the first two columns of A , so $\mathcal{R}(A) = \text{span}\{[1,0,0]^T, [1,1,0]^T\}$,
- $\dim(\mathcal{N}(A))=3-\dim(\mathcal{R}(A))=1$ and the basis is found by computing $An = \mathbf{0}$ which in this case is trivially given by $[1,0,-1]^T$, so $\mathcal{N}(A) = \text{span}\{[1,0,-1]^T\}$
- $\dim(\mathcal{N}(A^H))=3-\text{rank}(A)=1$; and we find a basis by imposing $\mathcal{N}^\perp(A^H) = \mathcal{R}(A)$, so $\mathcal{N}(A^H)=\text{span}\{[0,0,1]^T\}$
- Finally a basis for A^H is given by the first two rows of A therefore $\mathcal{R}(A^H) = \text{span}\{[1,1,1]^T, [0,1,0]^T\}$

Four Fundamental Subspaces – Worked Example

- **Question:** Find the four fundamental subspaces of $A = \begin{pmatrix} 1 & 1 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix}$
- **Answer:** In conclusion:
- $\mathcal{R}(A) = \text{span}\{[1,0,0]^T, [1,1,0]^T\}$, $\mathcal{N}(A) = \text{span}\{[1,0,-1]^T\}$
- $\mathcal{R}(A^H) = \text{span}\{[1,1,1]^T, [0,1,0]^T\}$, $\mathcal{N}(A^H) = \text{span}\{[0,0,1]^T\}$

Four Fundamental Subspaces - Example



Matrix Inverses

- Assume A is nonsingular and square, then we know that it has an inverse A^{-1}
- This means that $AA^{-1} = A^{-1}A = I$
- If for example the matrix is not square, it does not have an inverse but it may have a:
 - **Left Inverse** if a matrix B is such that $BA = I$ (but $AB \neq I$)
 - **Right Inverse** if there exists a matrix C such that $AC=I$
- The left inverse exists if A is full column rank (tall matrix)
- The right inverse exists if A is full row rank (fat matrix)
- Note that these matrices might not be unique!

Matrix Inverses: Examples

Assume $A = \begin{pmatrix} 2 & 0 & 0 \\ 0 & 7 & 0 \end{pmatrix}$, this matrix is full row-rank so it has a right inverse given

$$\text{by: } A_R = \begin{pmatrix} 1/2 & 0 \\ 0 & 1/7 \\ c_1 & c_2 \end{pmatrix}$$

Clearly $AA_R = I$ but $A_RA \neq I$

Exercise: Find a left inverse of A^T

Projections and Approximations

Assume that we have a vector $x \in H$ where H is a vector space and assume we want to find an approximation \hat{x} of x in a subspace V with $V \subset H$. The projection theorem tells us how to achieve that:

Projection Theorem: *Let V be a closed subspace of space H and let x be a vector in H .*

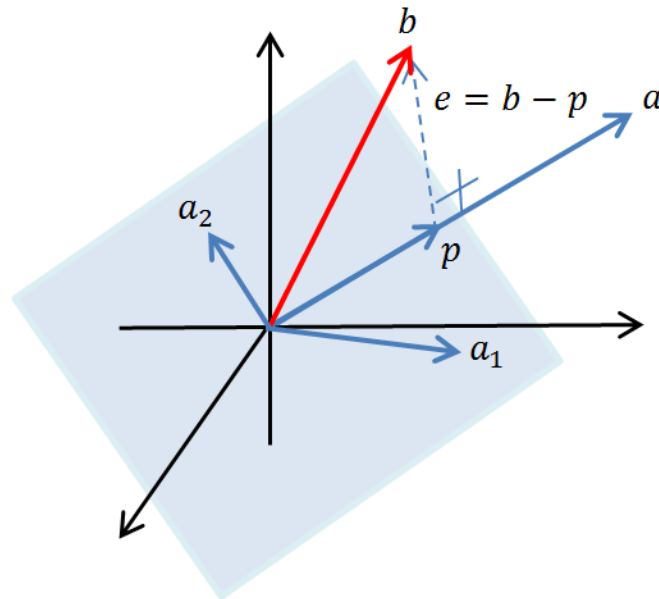
1. *Existence: There exists $\hat{x} \in V$ such that $\|x - \hat{x}\| \leq \|x - v\|$ for all $v \in V$.*
2. *Orthogonality: $x - \hat{x} \perp V$ is necessary and sufficient for determining \hat{x} .*
3. *Uniqueness: The vector \hat{x} is unique.*
4. *Linearity: $\hat{x} = Px$ where P is a linear operator that depends on V and not on x .*
5. *Idempotency: $P(Px) = Px$ for all $x \in H$.*
6. *Self-adjointness: $P = P^H$.*

Projections and Approximations

- Linearity of the projection means we can represent it using a matrix P
- Idempotent means that the matrix satisfies $P^2 = PP = P$
- Self-Adjoint means $P = P^H$
- Idempotency really means that once you have computed the projection, if you apply the projection operator again you get the same result

Projections and Approximations

- Assume we want to project a vector $x \in \mathbb{C}^m$ onto the subspace \mathbb{C}^n with $n < m$ (see the figure below for the case $m = 3$ and $n = 2$).
- Assume the vectors $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n$ span \mathbb{C}^n , we stack them together to form the tall matrix $\mathbf{A} \in \mathbb{C}^{m \times n}$. This matrix is full column rank.
- We now want to find the projection matrix \mathbf{P}



Projections and Approximations

Claim: the matrix $P = A(A^H A)^{-1} A^H$ is the orthogonal projection operator

Proof: First note that $A^H A$ is square and nonsingular because

$$A^H A x = 0 \Rightarrow x^H A^H A x = 0 \Rightarrow \|Ax\|^2 = 0 \Rightarrow Ax = 0 \Rightarrow x = 0$$

since the columns of A are linearly independent and so only the zero vector can give us zero. This means that P is well defined. Moreover, P is idempotent since

$$P^2 = A(A^H A)^{-1} A^H \cdot A(A^H A)^{-1} A^H = A(A^H A)^{-1} A^H = P$$

P is also self-adjoint so it is the projection operator.

It is also easy to verify that orthogonality is satisfied: $A^H(x - Px) = 0$

An Application of the Projection Theorem - The Gram-Schmidt orthogonalization process

- We introduced in the first week the notion of orthogonal matrices:
- Consider a set of vector v_i $i = 1, 2, \dots, n$ which forms an **orthonormal** basis and create a matrix A from stacking them one after the other:

$$A = \begin{pmatrix} v_{1,1} & \cdots & v_{1,n} \\ \vdots & \ddots & \vdots \\ v_{n,1} & \cdots & v_{n,n} \end{pmatrix}$$

- Clearly $A^H A = A A^H = I$, where A^H is the Hermitian transpose of A and I is the identity matrix. This means that $A^H = A^{-1}$.
- This type of matrices are called orthogonal matrices and **they preserve the induced norm**:
- If the set of vector v_i $i = 1, 2, \dots, n$ is not orthonormal, the process to make them orthogonal is known as **Gram-Schmidt process**.

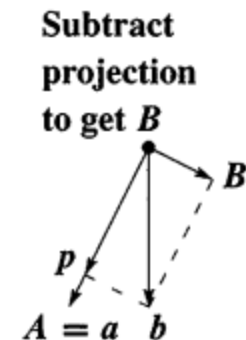
The Gram-Schmidt orthogonalization process

- The goal here is to start with three independent vectors $\mathbf{a}, \mathbf{b}, \mathbf{c}$ and construct three orthogonal vectors $\mathbf{A}, \mathbf{B}, \mathbf{C}$ and finally three orthonormal vectors.

$$\mathbf{q}_1 = \mathbf{A}/\|\mathbf{A}\|, \mathbf{q}_2 = \mathbf{B}/\|\mathbf{B}\|, \mathbf{q}_3 = \mathbf{C}/\|\mathbf{C}\|$$
- We begin by choosing $\mathbf{A} = \mathbf{a}$. This first direction is accepted.
- The next direction \mathbf{B} must be perpendicular to \mathbf{A} . We start with \mathbf{b} and subtract its projection along \mathbf{A} . This leaves the part of \mathbf{b} which we call vector \mathbf{B} (what we knew before as the error of projection), defined as (remember that the projection operator is $\mathbf{P} = \mathbf{A}(\mathbf{A}^H \mathbf{A})^{-1} \mathbf{A}^H$):

$$\mathbf{B} = \mathbf{b} - \frac{\mathbf{A}\mathbf{A}^T}{\mathbf{A}^T \mathbf{A}} \mathbf{b}$$

- By the projection theorem we know that the projection error is orthogonal to \mathbf{a}

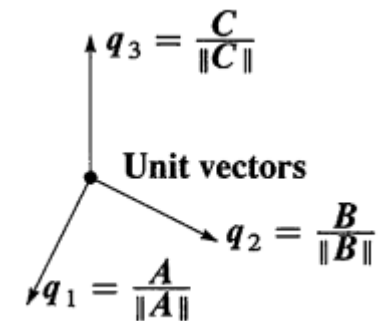


The Gram-Schmidt process

- The third direction starts with \mathbf{c} . This is not a combination of \mathbf{A} and \mathbf{B} .
- Most likely \mathbf{c} is not already perpendicular to \mathbf{A} and \mathbf{B} .
- Therefore, we subtract the projections of \mathbf{c} along \mathbf{A} and \mathbf{B} to get \mathbf{C} :

$$\mathbf{C} = \mathbf{c} - \frac{\mathbf{A}\mathbf{A}^T}{\mathbf{A}^T\mathbf{A}}\mathbf{c} - \frac{\mathbf{B}\mathbf{B}^T}{\mathbf{B}^T\mathbf{B}}\mathbf{c}$$

- **In general we subtract from every new vector its projections in the directions already set.**
- If we had a fourth vector \mathbf{d} , we would subtract three projections onto $\mathbf{A}, \mathbf{B}, \mathbf{C}$ to get \mathbf{D} .
- We make the resulting vectors orthonormal.
- This is done by dividing the vectors with their magnitudes.



The Gram-Schmidt process and QR Factorization

- Gram-Schmidt can be used to factorize a matrix \mathbf{A} into the product of an orthogonal and an upper triangular matrix. We focus only on the case \mathbf{A} is squared and invertible.
- Consider a squared matrix \mathbf{A} made of n linearly independent vectors $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n$:

$$\mathbf{A} = \begin{pmatrix} | & | & \cdots & | \\ \mathbf{a}_1 & \mathbf{a}_2 & \cdots & \mathbf{a}_n \\ | & | & \cdots & | \end{pmatrix}$$

- Apply Gram-Schmidt:

$$\begin{aligned} \bullet \quad \mathbf{u}_1 &= \mathbf{a}_1; & \mathbf{e}_1 &= \frac{\mathbf{u}_1}{\|\mathbf{u}_1\|} \\ \bullet \quad \mathbf{u}_2 &= \mathbf{a}_2 - \frac{\langle \mathbf{u}_1, \mathbf{a}_2 \rangle}{\|\mathbf{u}_1\|^2} \mathbf{u}_1; & \mathbf{e}_2 &= \frac{\mathbf{u}_2}{\|\mathbf{u}_2\|} \\ \bullet \quad \mathbf{u}_k &= \mathbf{a}_k - \sum_{j=1}^{k-1} \frac{\langle \mathbf{u}_j, \mathbf{a}_k \rangle}{\|\mathbf{u}_j\|^2} \mathbf{u}_j & \mathbf{e}_k &= \frac{\mathbf{u}_k}{\|\mathbf{u}_k\|} \end{aligned}$$

The Gram-Schmidt process and QR Factorization

- We now have (note that: $\langle \mathbf{e}_i, \mathbf{a}_i \rangle = ||\mathbf{u}_i||$):
- $\mathbf{a}_1 = \langle \mathbf{e}_1, \mathbf{a}_1 \rangle \mathbf{e}_1$
- $\mathbf{a}_2 = \langle \mathbf{e}_1, \mathbf{a}_2 \rangle \mathbf{e}_1 + \langle \mathbf{e}_2, \mathbf{a}_2 \rangle \mathbf{e}_2$
- $\mathbf{a}_3 = \langle \mathbf{e}_1, \mathbf{a}_3 \rangle \mathbf{e}_1 + \langle \mathbf{e}_2, \mathbf{a}_3 \rangle \mathbf{e}_2 + \langle \mathbf{e}_3, \mathbf{a}_3 \rangle \mathbf{e}_3$
- \vdots
- This means we can write $\mathbf{A} = \mathbf{QR}$ where \mathbf{Q} is orthogonal and \mathbf{R} is upper triangular with:

$$\bullet \quad \mathbf{Q} = (\mathbf{e}_1 \quad \cdots \quad \mathbf{e}_n) \text{ and } \mathbf{R} = \begin{pmatrix} \langle \mathbf{e}_1, \mathbf{a}_1 \rangle & \langle \mathbf{e}_1, \mathbf{a}_2 \rangle & \cdots & \langle \mathbf{e}_1, \mathbf{a}_n \rangle \\ 0 & \langle \mathbf{e}_2, \mathbf{a}_2 \rangle & \cdots & \langle \mathbf{e}_2, \mathbf{a}_n \rangle \\ \vdots & 0 & \ddots & \vdots \\ 0 & 0 & \cdots & \langle \mathbf{e}_n, \mathbf{a}_n \rangle \end{pmatrix}$$

Solution to $Ax=b$

- Let's go back to the problem of finding x such that $Ax = b$
- We are focusing on the full column rank case
- This means that A is tall and that if the solution exists it is unique
- We essentially need to find the left inverse of A
- We will achieve that through a somewhat convoluted but geometrical approach
- We look for x_{LS} such that $\|Ax_{LS} - b\| \leq \|Ax - b\|$ for any x
- So that if a solution to $Ax = b$ exists x_{LS} is that solution

Solution to $Ax=b$

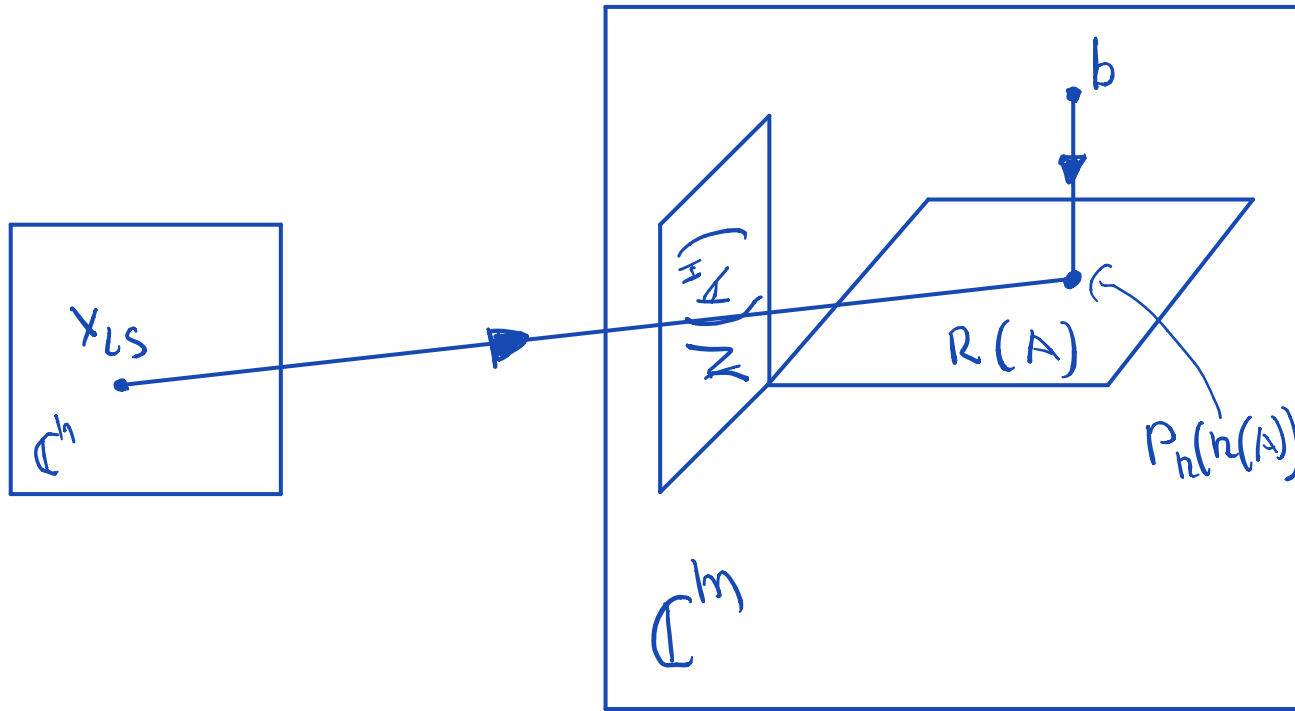
Claim: the vector \hat{x} minimizes $\|Ax - b\|$ if and only if $A^H A \hat{x} = A^H b$ (normal equation)

Proof:

- Minimizing $\|Ax - b\|$ is equivalent to minimizing $\|\hat{b} - b\|$ where $\hat{b} \in \mathfrak{R}(A)$.
- By the projection theorem: $b - \hat{b} \in \mathfrak{R}(A)^\perp \Rightarrow b - \hat{b} \in \mathcal{N}(A^H)$
- Therefore, $A^H(b - \hat{b}) = 0 \Rightarrow A^H b = A^H \hat{b} \Rightarrow A^H A \hat{x} = A^H b$
- Conversely $A^H A \hat{x} = A^H b \Rightarrow A^H(b - A\hat{x}) = 0$

Least-Squares Solution

- Geometric Interpretation to LS Solution



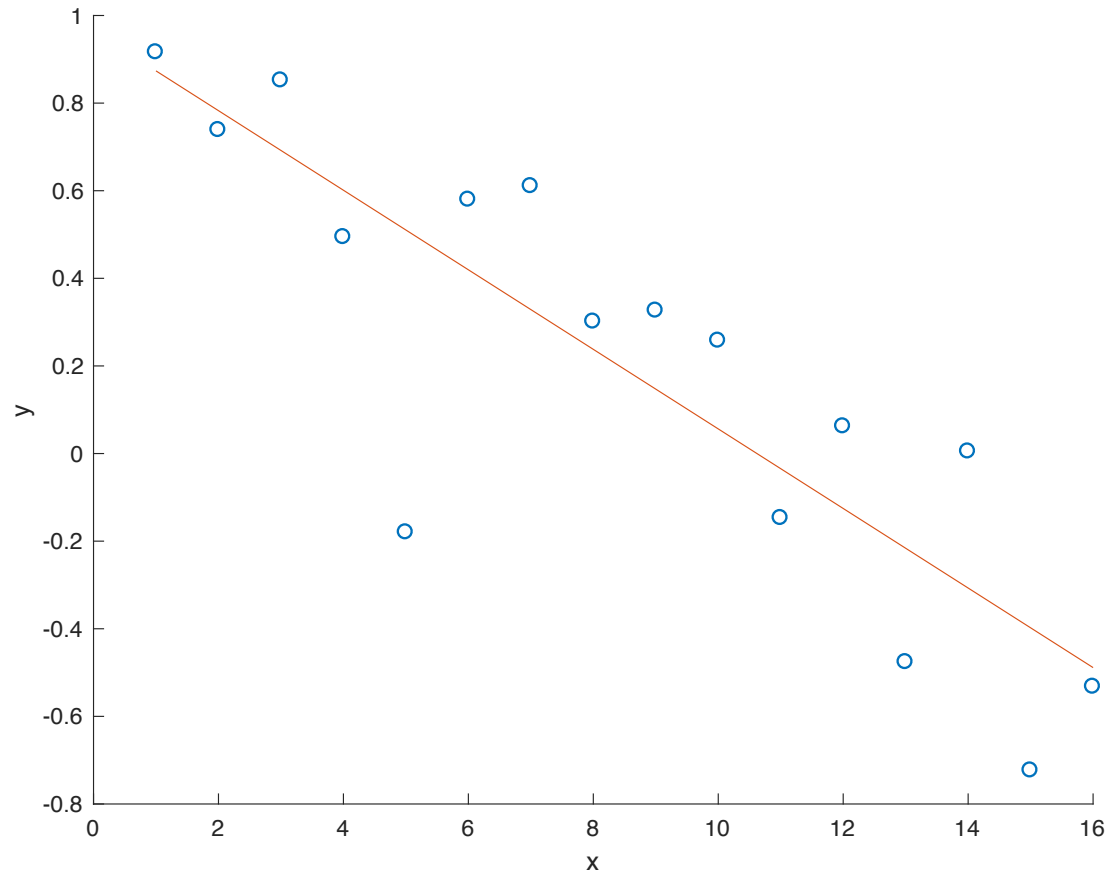
Example: Linear Regression

- Assume you observe a function $f(x)$ at instants x_1, x_2, \dots, x_n : $y_i = f(x_i)$ $i = 1, 2, \dots, n$
- You want to approximate $f(x)$ with a line, that is, $y_i \approx ax_i + b$
- In matrix/vector form: $\mathbf{y} = \mathbf{A}\mathbf{c} + \mathbf{e}$, where $\mathbf{y} = [y_1 \ y_2 \ \cdots \ y_n]^T$, $\mathbf{c} = \begin{bmatrix} a \\ b \end{bmatrix}$, \mathbf{e} is

the approximation error and $\mathbf{A} = \begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{bmatrix}$.

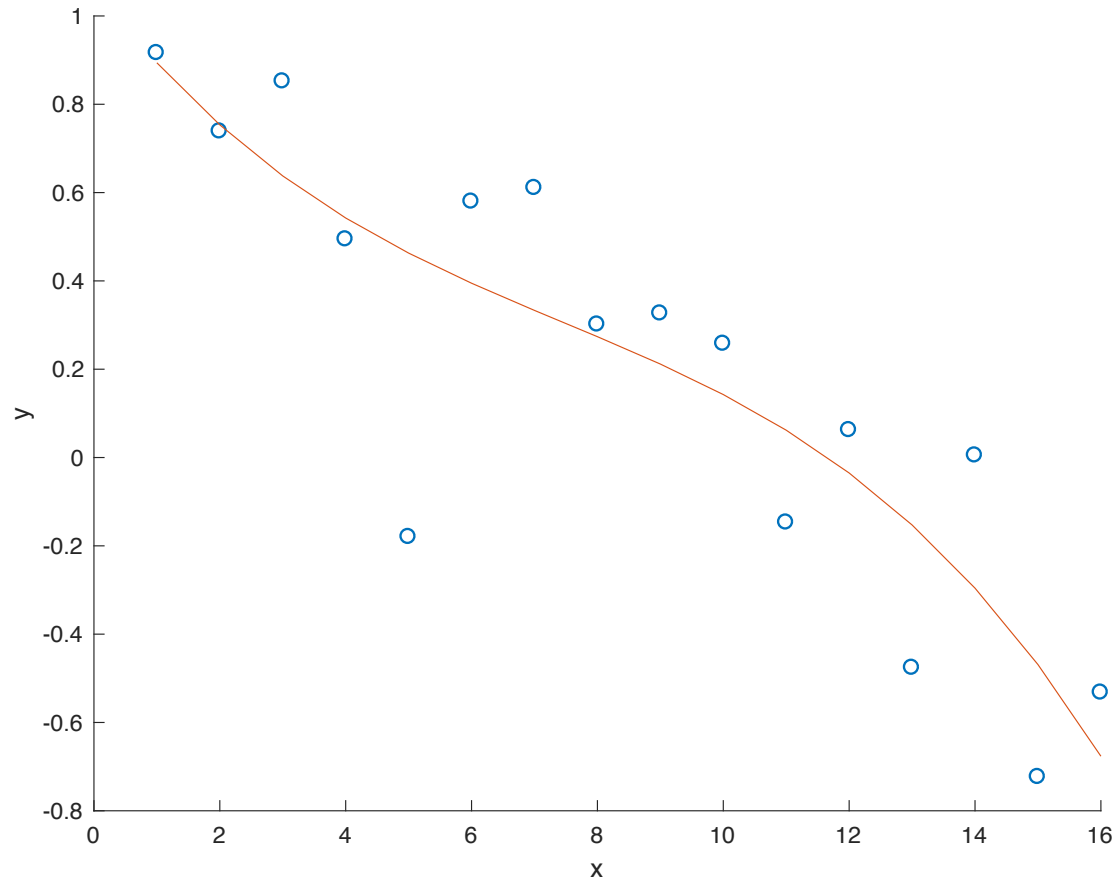
- You want to minimize the norm of the error. Note that \mathbf{A} is tall and full column rank, so you can use the least-square solution discussed before: $\mathbf{A}^H \mathbf{A} \hat{\mathbf{c}} = \mathbf{A}^H \mathbf{y}$

Example: Linear Regression



- Can you find A in the case we want to fit a higher order polynomial to the data?

Example: High-order Regression



- Can you find A in the case we want to fit a higher order polynomial to the data?

Example: least-square deconvolution

Assume the FIR deconvolution problem: the output of the filter is $\mathbf{y} \in \mathbb{C}^{n+K-1}$ and the unknown input signal is $\mathbf{x} \in \mathbb{C}^n$ and the two are related by $\mathbf{y} = \mathbf{A}\mathbf{x}$ where $\mathbf{A} \in \mathbb{C}^{(n+K-1) \times n}$ is the convolution (Toeplitz) matrix associated with the known K -tap filter \mathbf{h} .

Given noisy measurements $\mathbf{b} = \mathbf{y} + \mathbf{n}$ we want to find the \mathbf{x} that is closest to the correct one and when there is no noise we want the solution to be correct.

\mathbf{A} is tall and we can use the previous result, thus, $\mathbf{A}^H \mathbf{A} \hat{\mathbf{x}} = \mathbf{A}^H \mathbf{b}$

Example: least-square deconvolution

A is tall and we can use the previous result, thus, $A^H A \hat{x} = A^H b$

$x \rightarrow$



$y = Ax \rightarrow$



$b = y + n \rightarrow$



$x_{LS} \rightarrow$



$x_{LS} \rightarrow$



Least-Square Solution with Tikhonov Regularization

- In some cases you need to “regularize” the reconstruction process. This is an active research area. The most traditional regularization is known as “Tikhonov” regularization and tries to constrain the norm of \mathbf{x} :

$$\min_{\mathbf{x}} (\|\mathbf{Ax} - \mathbf{b}\|^2 + \lambda \|\mathbf{x}\|^2)$$

- This is equivalent to the modified least-square problem:

$$\min_{\mathbf{x}} (\|\tilde{\mathbf{A}}\mathbf{x} - \tilde{\mathbf{b}}\|^2) \text{ with } \tilde{\mathbf{A}} = \begin{bmatrix} \mathbf{A} \\ \sqrt{\lambda}\mathbf{I} \end{bmatrix} \text{ and } \tilde{\mathbf{b}} = \begin{bmatrix} \mathbf{b} \\ \mathbf{0} \end{bmatrix}$$

- This regularization essentially improves the condition number of $\mathbf{A}^H \mathbf{A}$.

Least-Square Solution with Tikhonov Regularization



Original Image



Blurred Image with noise



LS with Tikhonov

Binary Classification Using Least-Squares

- Assume you have a set of “feature” vectors representing your training data: \mathbf{v}_i $i = 1, 2, \dots, N$. These vectors belong to two different classes which you label with “1” or “-1”. So we associate $y_i = \pm 1$ to each \mathbf{v}_i .
- We want to find a vector \mathbf{x} such that $\langle \mathbf{x}, \mathbf{v}_i \rangle > 0$ if $y_i = 1$ and $\langle \mathbf{x}, \mathbf{v}_i \rangle < 0$ if $y_i = -1$
- We create a matrix A where each row is one feature vector and then solve:
- $\min_{\mathbf{x}} (\|A\mathbf{x} - \mathbf{y}\|^2)$ where \mathbf{y} is the vector with entries y_i .
- Once you have learned \mathbf{x} , then given a new “test” feature vector you expect the inner product to help classify this new vector.
- This is a simple classifier, but it is very simple to implement and is fast.

Solution to $Ax=b$ (full row rank case)

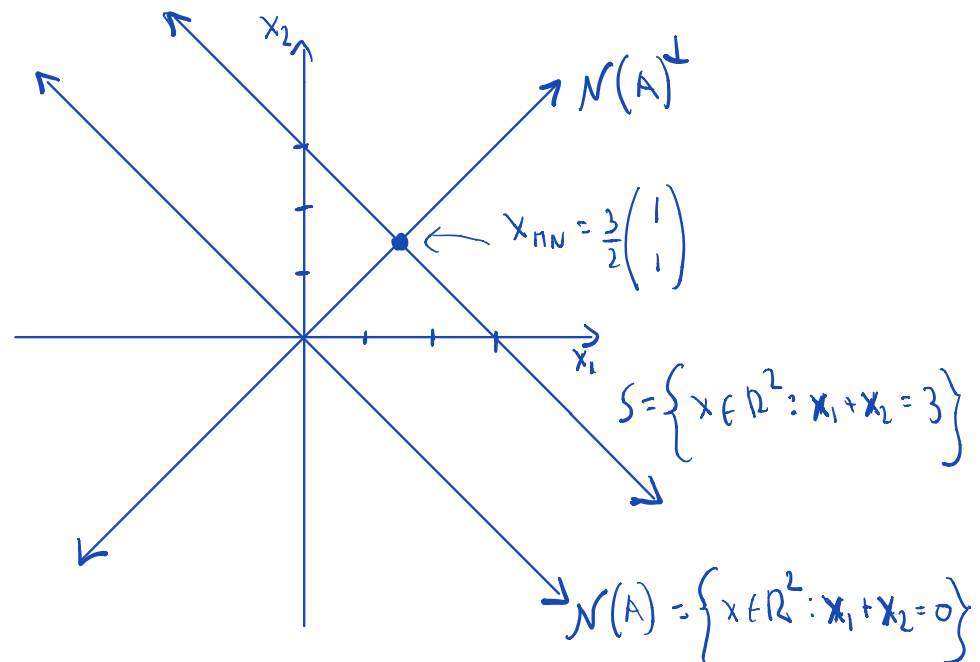
- We now understand how to find an approximate solution $Ax = b$ in the full column rank case where a solution may not exist.
- We now consider the full row rank case (A is 'fat') and we know that at least one solution exists.
- When the solution exists, but is not unique, the set of all solutions has a special structure.
- Let x_p be one solution to $Ax = b$. Then any other solution is $x = x_p + z$ with $z \in \mathcal{N}(A)$
- Then the solution to $Ax = b$ is: $S = \{x_p\} + \mathcal{N}(A)$
- Since $\mathcal{N}(A)$ is a subspace S is a *translated subspace* also called **affine subspace**

Solution to $Ax=b$ (full row rank case)

- We now consider the full row rank case (A is 'fat') and we know that at least one solution exists.
- We want to pick one of the (infinitely) many solutions
- We look for the *minimum norm solution*:
 - minimize $\|x\|$ subject to $Ax = b$

Minimum Norm Solution - Example

- Let $A = [1, 1]$ and suppose $b=3$.
- The set of all possible (real) solutions to $Ax=b$ is
- $T = \{x: x_1 + x_2 = 3\}$
- So $x = \begin{pmatrix} \alpha \\ 3 - \alpha \end{pmatrix}$ and $\|x\|^2 = \alpha^2 + (3 - \alpha)^2$
- By taking the derivative of the norm and setting it to zero we find the minimum norm solution which is: $x = \begin{pmatrix} 3/2 \\ 3/2 \end{pmatrix}$



Minimum Norm Solution

Dual Projection Theorem (without proof): Let $T = \{\mathbf{x}_p\} + \mathcal{S}$ be an affine subspace of \mathbb{C}^n . Then the element \mathbf{t}_{MN} of T with minimum norm exists, is unique, and satisfies $\mathbf{t}_{MN} \in T \cap \mathcal{S}^\perp$. Moreover, $\mathbf{t}_{MN} = P_{\mathcal{S}^\perp} \mathbf{t}$.

Minimum Norm Theorem: Let $\mathbf{Ax}=\mathbf{b}$ and let \mathbf{A} be a full-row rank matrix. Amongst all the $\mathbf{x} \in \mathbb{C}^n$ satisfying $\mathbf{Ax}=\mathbf{b}$, there exists a unique solution with minimum norm. It lies in $\mathfrak{R}(\mathbf{A}^H)$ and is given by $\mathbf{x}_{MN} = \mathbf{A}_{MN} \mathbf{b}$ with $\mathbf{A}_{MN} = \mathbf{A}^H (\mathbf{A} \mathbf{A}^H)^{-1}$.

Proof: If \mathbf{x}_p is a solution to $\mathbf{Ax}=\mathbf{b}$ then all the solutions are in the affine subspace $T = \{\mathbf{x}_p\} + \mathcal{N}(\mathbf{A})$. Using the Dual Projection Theorem we conclude that

$$\mathbf{x}_{MN} \in T \cap \mathcal{N}(\mathbf{A})^\perp = T \cap \mathfrak{R}(\mathbf{A}^H)$$

Since $\mathbf{x}_{MN} \in \mathfrak{R}(\mathbf{A}^H) \Rightarrow \mathbf{x}_{MN} = \mathbf{A}^H \mathbf{z} \Rightarrow \mathbf{b} = \mathbf{Ax}_{MN} = (\mathbf{A} \mathbf{A}^H) \mathbf{z}$.

Therefore, $\mathbf{z} = (\mathbf{A} \mathbf{A}^H)^{-1} \mathbf{b}$ which implies that $\mathbf{x}_{MN} = \mathbf{A}^H (\mathbf{A} \mathbf{A}^H)^{-1} \mathbf{b}$

Minimum Norm Solution: worked example

Question: find the minimum norm solution to $\mathbf{Ax}=\mathbf{b}$ with

$$\mathbf{A} = \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix} \text{ and } \mathbf{b} = [1, 2]^T$$

Approach 1 (recommended): Compute directly $\mathbf{x}_{MN} = \mathbf{A}^H(\mathbf{A}\mathbf{A}^H)^{-1}\mathbf{b}$ which is in this case

$$\mathbf{x}_{MN} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}^{-1} \begin{bmatrix} 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 1/2 & 0 \\ 0 & 1 \\ 1/2 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 1/2 \\ 2 \\ 1/2 \end{bmatrix}$$

Approach 2: the constraints are $x_1 + x_3 = 1$; $x_2 = 2$, therefore $\mathbf{x}_p = [\alpha, 2, 1 - \alpha]^T$
 $\|\mathbf{x}_p\|^2 = \alpha^2 + 4 + (1 - \alpha)^2$ which implies $\alpha = \frac{1}{2}$ and $\mathbf{x}_{MN} = \left[\frac{1}{2}, 2, \frac{1}{2}\right]^T$

Minimum Norm Solution: worked example

Question: find the minimum norm solution to $\mathbf{Ax}=\mathbf{b}$ with

$$\mathbf{A} = \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix} \text{ and } \mathbf{b} = [1, 2]^T$$

Approach 3: $\mathcal{N}(\mathbf{A}) = \text{span}\{ [1, 0, -1]^T \}$ since the minimum norm solution $\mathbf{x} \in \mathcal{N}(\mathbf{A})^\perp$ we need that if $\mathbf{n} \in \mathcal{N}(\mathbf{A})$ then $\mathbf{n}^T \mathbf{x} = 0$, we know that $x_1 + x_3 = 1$; $x_2 = 2$, therefore $\mathbf{n}^T \mathbf{x} = 0$ implies $x_1 - x_3 = 0$ and $\mathbf{x}_{MN} = \left[\frac{1}{2}, 2, \frac{1}{2} \right]^T$

Approach 4: Given a possible solution \mathbf{x}_p project it onto $\mathcal{R}(\mathbf{A}^H)$ so in our case

$$\mathbf{x}_{MN} = \mathbf{A}^H (\mathbf{A} \mathbf{A}^H)^{-1} \mathbf{A} \mathbf{x}_p = \frac{1}{2} \begin{bmatrix} 1 & 0 & 1 \\ 0 & 2 & 0 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 0 \end{bmatrix} = \begin{bmatrix} 1/2 \\ 2 \\ 1/2 \end{bmatrix}$$

Computing left and right inverses

Claim: Assume \mathbf{A} is 'tall' and full column rank then \mathbf{A} has a left inverse \mathbf{A}_{left}

Sketch of proof: $\mathbf{A}^T \mathbf{A}$ is full rank so $(\mathbf{A}^T \mathbf{A})^{-1}$ exists. Then $\mathbf{A}_{left} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$

Likewise: Assume \mathbf{A} is 'fat' and full row rank then \mathbf{A} has a right inverse \mathbf{A}_{rig} and

$$\mathbf{A}_{rig} = \mathbf{A}^T (\mathbf{A} \mathbf{A}^T)^{-1}$$

Least-Square problems with linear constraints

- Let's revisit the minimum norm solution approach:
 - Remember that we are back to the full row rank case (\mathbf{A} is 'fat') and we know that at least one solution exists.
 - Since there are (infinitely) many solutions
 - We look for a specific one:
 - minimize $\|\mathbf{x}\|^2$ subject to $\mathbf{Ax} = \mathbf{b}$
- The idea now is to interpret the above equation differently: as a least square minimization with some (linear constraints)
- A more general form is: minimize $\|\mathbf{y} - \mathbf{Cx}\|^2$ subject to $\mathbf{Ax} = \mathbf{b}$

Least-Square problems with linear constraints

- Problem: minimize $\|\mathbf{y} - \mathbf{C}\mathbf{x}\|^2$ subject to $\mathbf{A}\mathbf{x} = \mathbf{b}$
 - We use Lagrangian multipliers to solve the problem:
 - Minimize $L(\mathbf{x}, \boldsymbol{\lambda}) = \|\mathbf{y} - \mathbf{C}\mathbf{x}\|^2 + \lambda_1(\mathbf{a}_1^T \mathbf{x} - b_1) + \lambda_2(\mathbf{a}_2^T \mathbf{x} - b_2) + \dots + \lambda_m(\mathbf{a}_m^T \mathbf{x} - b_m)$, where \mathbf{a}_i^T are the rows of \mathbf{A}
 - Compute $\frac{\partial L(\mathbf{x}, \boldsymbol{\lambda})}{\partial x_i} = 0 \Rightarrow 2\mathbf{C}^T \mathbf{C}\mathbf{x} - 2\mathbf{C}^T \mathbf{y} + \mathbf{A}^T \boldsymbol{\lambda} = \mathbf{0}$
 - Combining this set of linear equations with the feasibility conditions $\mathbf{A}\mathbf{x} = \mathbf{b}$, we can write the optimality conditions as one set of $m+n$ linear equations in the variables $(\mathbf{x}; \boldsymbol{\lambda})$:

$$\begin{pmatrix} 2\mathbf{C}^T \mathbf{C} & \mathbf{A}^T \\ \mathbf{A} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{x} \\ \boldsymbol{\lambda} \end{pmatrix} = \begin{pmatrix} 2\mathbf{C}^T \mathbf{y} \\ \mathbf{b} \end{pmatrix}$$

Least-Square problems with linear constraints

- Note that the square matrix $\begin{pmatrix} 2\mathbf{C}^T\mathbf{C} & \mathbf{A}^T \\ \mathbf{A} & \mathbf{0} \end{pmatrix}$ is invertible if and only if \mathbf{A} is full row rank and $\begin{pmatrix} \mathbf{C} \\ \mathbf{A} \end{pmatrix}$ is full column rank

Minimum Norm Least Square Solution to $Ax=b$

- We have so far dealt with the full column and full row rank cases
- When A is neither full column or row rank, a solution to $Ax = b$ may not exist and if it does it is not unique
- Let $A = \begin{pmatrix} 1 & 2 \\ 2 & 4 \end{pmatrix}$ and $b = \begin{pmatrix} 2 \\ 3 \end{pmatrix}$, we note $b \notin \mathcal{R}(A)$, furthermore $\mathcal{N}(A) = \text{span} \begin{pmatrix} 2 \\ -1 \end{pmatrix}$
- We may try to find x_{LS} by solving $A^H A x_{LS} = A^H b$ which yields

$$5 \begin{pmatrix} 1 & 2 \\ 2 & 4 \end{pmatrix} x_{LS} = 8 \begin{pmatrix} 1 \\ 2 \end{pmatrix}$$

- There are an infinity of solutions all in the affine subspace $T = \begin{pmatrix} 1 \\ 3/10 \end{pmatrix} + \text{span} \begin{pmatrix} 2 \\ -1 \end{pmatrix}$
- We may then pick the one with minimum norm

Minimum Norm Least Square Solution to $Ax=b$

- It turns out that the minimum norm least square solution is unique and is given by the Moore-Penrose pseudoinverse of A which we denote with A^+
- Given A its Moore-Penrose pseudoinverse A^+ is defined as the matrix satisfying:
 1. $AA^+A = A$
 2. $A^+AA^+ = A^+$
 3. $(AA^+)^H = AA^+$
 4. $(A^+A)^H = A^+A$
- The pseudo inverse is computed using the **Singular Value Decomposition** (SVD)