

NOTES

James Guo

November 1, 2024

Contents

I Preliminaries	1
I.1 Introduction	1
I.2 Review: Vector Spaces and Subspaces	1
I.3 Normed Vector Space	8
I.4 Review: Eigenspace	14
I.5 Diagonalization and Singular Value Decomposition	18
I.6 Projection	25
I.7 QR Decomposition	33
II Applications with Computer Programming	40
II.1 MATLAB Preliminaries	40
II.2 Representation of Numbers	41
III Computational Methods	49
III.1 Least Square Approximations	49
III.2 Conditioning and Condition Number	53
III.3 Stability	56
III.4 Stability and Gaussian Elimination	61

I Preliminaries

I.1 Introduction

The aim of the course is to deal with matrices with *much more* entries, with computation of eigenvalues and eigenvectors. In particular, the cost of computing the *determinant* is large, *i.e.*, the characteristic polynomial has high degree. Therefore, we want *computability* and *closeness* of the eigenvalues.

Numerically, we want to evaluate the *deviation* of the value and the actual computation.

Remark I.1.1. Notations and Conventions.

This course aligns to the following notations and conventions:

- Vectors will be represented as bold cases, such as \mathbf{x} and \mathbf{y} .
- Matrices ($m \times n$, with m rows and n columns) will be represented with capitalized letters, such as A , where a_{ij} represents the the entry of row i and column j .
- The *field* of the R -module, unless otherwise specified, will be assumed to be \mathbb{C} .
- For $A \cdot \mathbf{x}$, let $f : \mathbb{C}^n \rightarrow \mathbb{C}^m$ as $\mathbf{y} = A \cdot \mathbf{x}$, which can be alternatively represented as $f : \mathbf{x} \mapsto A\mathbf{x}$.
- The class of m -by- n matrices with entries in field \mathbb{F} is denoted as $\mathbb{F}^{m \times n}$.

I.2 Review: Vector Spaces and Subspaces

The Numerical Linear Algebra has its foundations on linear algebra.

Note: Unless otherwise specified, the proofs of the theorems in this section is omitted, as they are assumed backgrounds in a typical linear algebra course.

Definition I.2.1. Linearity.

A function f is linear if:

$$\begin{cases} f(\mathbf{x} + \mathbf{y}) = f(\mathbf{x}) + f(\mathbf{y}) \text{ for all } \mathbf{x}, \mathbf{y} \in \mathbb{C}^n, \\ f(\alpha \mathbf{x}) = \alpha f(\mathbf{x}) \text{ for all } \alpha \in \mathbb{C} \text{ and } \mathbf{x} \in \mathbb{C}^n. \end{cases}$$

Theorem I.2.2. Linear Function as a Matrix.

If $f : \mathbb{C}^n \rightarrow \mathbb{C}^m$ is linear transformation, then there exists $A \in \mathbb{C}^{m \times n}$ such that $f(\mathbf{x}) = A \cdot \mathbf{x}$ for all $\mathbf{x} \in \mathbb{C}^n$.

Remark I.2.3. Linear Combinations.

For $A \in \mathbb{C}^{n \times m}$, we may write:

$$A = [a_{ij}] = [\mathbf{A}_1 \quad \mathbf{A}_2 \quad \cdots \quad \mathbf{A}_m],$$

and for $\mathbf{x} \in \mathbb{C}^n$:

$$A \cdot \mathbf{x} = [\mathbf{A}_1 \quad \mathbf{A}_2 \quad \cdots \quad \mathbf{A}_m] \cdot \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = x_1 \mathbf{A}_1 + x_2 \mathbf{A}_2 + \cdots + x_n \mathbf{A}_n.$$

Definition I.2.4. Matrix Multiplication.

Let $A \in \mathbb{C}^{l \times m}$ and $C \in \mathbb{C}^{m \times n}$, their matrix multiplication is defined as:

$$[b_{ij}] = B = AC = [a_{ik}][c_{kj}],$$

in which the entry b_{ij} in B is:

$$b_{ij} = \sum_{k=1}^m a_{ik}c_{kj}.$$

Remark I.2.5. Multiplication of Matrices as Columns.

For the above multiplication that $B = AC$, if we write B as vectors we have:

$$\begin{bmatrix} \mathbf{B}_1 & \mathbf{B}_2 & \cdots & \mathbf{B}_n \end{bmatrix} = A \begin{bmatrix} \mathbf{C}_1 & \mathbf{C}_2 & \cdots & \mathbf{C}_n \end{bmatrix} = \begin{bmatrix} A \cdot \mathbf{C}_1 & A \cdot \mathbf{C}_2 & \cdots & A \cdot \mathbf{C}_n \end{bmatrix},$$

hence:

$$\mathbf{B}_j = A \cdot \mathbf{C}_j = \begin{bmatrix} \mathbf{A}_1 & \mathbf{A}_2 & \cdots & \mathbf{A}_m \end{bmatrix} \cdot \begin{bmatrix} c_{1j} \\ c_{2j} \\ \vdots \\ c_{mj} \end{bmatrix},$$

hence each column of $B = AC$ is a linear combination of the columns of A .

Definition I.2.6. Inner and Outer Product.

Let $\mathbf{u}, \mathbf{v} \in \mathbb{C}^n$, their inner product (or dot product) is:

$$\mathbf{u} \cdot \mathbf{v} = \mathbf{u}^T \cdot \bar{\mathbf{v}} = u_1 \bar{v}_1 + u_2 \bar{v}_2 + \cdots + u_n \bar{v}_n \in \mathbb{C},$$

where as the outer product is:

$$\mathbf{u} \cdot \mathbf{v}^T = \begin{bmatrix} u_1 v_1 & u_1 v_2 & \cdots & u_1 v_n \\ u_2 v_1 & u_2 v_2 & \cdots & u_2 v_n \\ \vdots & \vdots & \ddots & \vdots \\ u_n v_1 & u_n v_2 & \cdots & u_n v_n \end{bmatrix} = \begin{bmatrix} v_1 \mathbf{u} & v_2 \mathbf{u} & \cdots & v_n \mathbf{u} \end{bmatrix}.$$

Definition I.2.7. Range and Null Space.

For $A \in \mathbb{C}^{m \times n}$, its range (or image) is:

$$\text{im } A = \{\mathbf{y} : \mathbf{y} = A \cdot \mathbf{x} \text{ for some } \mathbf{x} \in \mathbb{C}^n\} = \text{span}\{\mathbf{A}_1, \mathbf{A}_2, \cdots, \mathbf{A}_n\} \subset \mathbb{C}^m.$$

Its null space (or kernel) is:

$$\ker A = \{\mathbf{x} : A \cdot \mathbf{x} = \mathbf{0}\} \subset \mathbb{C}^n.$$

Theorem I.2.8. Rank-Nullity Theorem (Fundamental Theorem of Linear Maps).

For $A \in \mathbb{C}^{m \times n}$, the dimension of the range (or rank) and the dimension of the null space (or nullity) follows:

$$\dim(\text{im } A) + \dim(\ker A) = n.$$

Moreover, for the transpose of A , we have:

$$\dim(\operatorname{im} A^\top) = \dim(\operatorname{im} A).$$

Proposition I.2.9. Properties on the Rank of A .

Assume that $A \in \mathbb{C}^{m \times n}$, the following holds:

- (i) $\dim(\operatorname{im} A) \leq n$.
- (ii) $\dim(\operatorname{im} A) \leq m$.

The result (ii) is a direct result of (i) and Rank-Nullity for transpose of A . Moreover, this implies that $\dim(\operatorname{im} A) \leq \min\{m, n\}$.

Definition I.2.10. Full Rank.

For $A \in \mathbb{C}^{m \times n}$, A is full rank if $\dim(\operatorname{im} A) = \min\{m, n\}$.

An example of a full rank matrix is $\begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}$.

Definition I.2.11. Nonsingular Square Matrix.

A m -by- m square matrix A is nonsingular if $A \cdot \mathbf{x} = \mathbf{0}$ has unique solution $\mathbf{x} = \mathbf{0}$.

Proposition I.2.12. Equivalent Facts for Square Matrices.

The following are equivalent:

- (i) A is nonsingular,
- (ii) A is invertible, i.e., there exists A^{-1} such that $A^{-1}A = AA^{-1} = \operatorname{Id}$,
- (iii) For all $\mathbf{b} \in \mathbb{C}^m$, $A \cdot \mathbf{x} = \mathbf{b}$ has unique solution,
- (iv) Columns of A are linearly independent,
- (v) $\dim(\operatorname{im} A) = m$,
- (vi) $\operatorname{im} A = \mathbb{C}^m$,
- (vii) $\dim(\operatorname{ker} A) = 0$,
- (viii) $\operatorname{ker} A = \{\mathbf{0}\}$,
- (ix) $\det A \neq 0$,
- (x) A is an isomorphism.

In particular, since A being nonsingular is the equivalent of being injective for square matrices, it is an isomorphism, hence all (ii) to (x) are equivalent to (i).

Remark I.2.13. Coordinates in \mathbb{C}^m .

There are many basis in \mathbb{C}^m :

(i) The Canonical Basis:

$$\alpha = \{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n\},$$

where:

$$\mathbf{e}_1 = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \mathbf{e}_2 = \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}, \dots, \mathbf{e}_n = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}.$$

Hence, for any $\mathbf{b} \in \mathbb{C}^m$, we have:

$$\mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix} = b_1 \mathbf{e}_1 + b_2 \mathbf{e}_2 + \dots + b_n \mathbf{e}_n.$$

(ii) Another Basis (Arbitrary):

$$\beta = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}.$$

Here, for any $\mathbf{b} \in \mathbb{C}^m$, we are looking forward to having:

$$\mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix} = c_1 \mathbf{v}_1 + c_2 \mathbf{v}_2 + \dots + c_m \mathbf{v}_m = \begin{bmatrix} \mathbf{v}_1 & \mathbf{v}_2 & \dots & \mathbf{v}_m \end{bmatrix} \cdot \mathbf{c},$$

and since the matrix is invertible, it leads to that:

$$\mathbf{c} = \begin{bmatrix} \mathbf{v}_1 & \mathbf{v}_2 & \dots & \mathbf{v}_m \end{bmatrix}^{-1} \cdot \mathbf{b},$$

in which by solving for \mathbf{c} here, we have completed the decomposition.

Remark I.2.14. Remarks on Complex Numbers.

For $z \in \mathbb{C} = \mathbb{R}(i)$ in which $i^2 = -1$, we may represent $z = a + ib$, where $a, b \in \mathbb{R}$. Thus, we can visualize the real and imaginary parts, respectively.

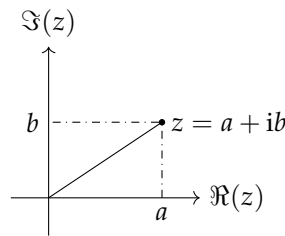


Figure I.1. Real and Complex Plane.

Here, we have the magnitude as:

$$|z| = \sqrt{a^2 + b^2},$$

and we have the complex conjugate as:

$$\bar{z} = a - ib.$$

Example I.2.15. Matrix Operations.

Let $A \in \mathbb{C}^{m \times n}$, that is:

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots \\ a_{21} & a_{22} & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix},$$

its transpose is:

$$A^T = \begin{bmatrix} a_{11} & a_{21} & \cdots \\ a_{12} & a_{22} & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix},$$

its complex conjugate is:

$$\bar{A} = \begin{bmatrix} \bar{a}_{11} & \bar{a}_{12} & \cdots \\ \bar{a}_{21} & \bar{a}_{22} & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix},$$

and the Hermitian conjugate is:

$$A^* = (\bar{A})^T.$$

In \mathbb{R}^n , the dot product for $\mathbf{x}, \mathbf{y} \in \mathbb{R}^m$, we have:

$$\mathbf{x} \cdot \mathbf{y} = \mathbf{x}^T \cdot \mathbf{y} = \sum_{k=1}^m x_k y_k.$$

In particular, we can define the norm as:

$$\|\mathbf{x}\| = \sqrt{\mathbf{x}^T \cdot \mathbf{x}} = \sqrt{\sum_{k=1}^m x_k^2}.$$

In \mathbb{C}^n , i.e., when $\mathbb{F} = \mathbb{C}$, and the inner product for $\mathbf{x}, \mathbf{y} \in \mathbb{C}^m$, we have:

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^* \cdot \mathbf{y} = \sum_{k=1}^m \bar{x}_k y_k.$$

Again, we have the norm as:

$$\|\mathbf{x}\| = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle} = \sqrt{\sum_{k=1}^m \bar{x}_k x_k} = \sqrt{\sum_{k=1}^m |x_k|^2},$$

since we have $\bar{x}_k x_k = |x_k|^2$, known as the *modulus*.

In particular $\mathbf{x} = \mathbf{0}$ if and only if $\|\mathbf{x}\| = 0$.

Proposition I.2.16. Distributivity of Matrix Operator.

For any $A \in \mathbb{C}^{m \times n}$ and $B \in \mathbb{C}^{n \times p}$, we have:

$$(AB)^T = B^T A^T \text{ and } (AB)^* = B^* A^*.$$

Definition I.2.17. Orthogonality.

For any $\mathbf{x}, \mathbf{y} \in \mathbb{C}^m$, \mathbf{x} and \mathbf{y} are orthogonal if $\mathbf{x}^* \cdot \mathbf{y} = 0$.

For sets $S = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\} \subset \mathbb{C}^m$ is orthogonal set if the vectors are pairwise orthogonal, i.e., $\mathbf{v}_i^* \cdot \mathbf{v}_j = 0$ for all $i \neq j$. J

Theorem I.2.18. Orthogonality \implies Linear Independence.

If $S = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\} \subset \mathbb{C}^m$ is a orthogonal set of *nonzero* vectors, then S is linearly independent.

Proof. Here, we let λ_k 's be set such that:

$$\sum_{k=1}^m \lambda_k \mathbf{v}_k = \mathbf{0}.$$

For all $1 \leq i \leq m$ the inner product with \mathbf{v}_i , giving us that:

$$\sum_{k=1}^m \lambda_k \mathbf{v}_k^* \mathbf{v}_i = \mathbf{0}^* \mathbf{v}_i = 0.$$

Thus, by orthogonality and the nonzero vectors, we have:

$$\lambda_i = 0 \text{ for all } 1 \leq i \leq m,$$

which enforces S to be linearly independent. □

Corollary I.2.19. Surjective + Orthogonal \implies Basis.

Suppose that $S = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\} \subset V$ is a basis in which $V = \text{span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$, then S is a orthogonal basis of V .

Definition I.2.20. Orthogonal Basis.

An orthogonal basis can be $\beta = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r\} \subset \mathbb{C}^m$ with property that $\mathbf{v}_i^* \cdot \mathbf{v}_j = 0$ for all $i \neq j$.

In particular, we have $V = \text{span}(\beta)$ as a subspace of \mathbb{C}^m . J

Theorem I.2.21. Orthogonal Projection.

Given an orthogonal basis:

$$\beta = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r\} \subset \mathbb{C}^m$$

of $V = \text{span}(\beta)$ and a vector $\mathbf{b} \in V$, we have the coordinates of \mathbf{b} with respect to β , i.e., the unique scalars c_1, \dots, c_r such that $\mathbf{b} = \sum_{k=1}^r c_k \mathbf{v}_k$, are:

$$c_k = \frac{\mathbf{v}_k^* \mathbf{b}}{\|\mathbf{v}_k\|^2} \text{ for all } 1 \leq k \leq r.$$

Proof. For all $1 \leq i \leq m$, we take the inner product of the linear combinations with \mathbf{v}_i^* , we have:

$$\mathbf{v}_i^* \mathbf{b} = \sum_{k=1}^r c_k \mathbf{v}_i^* \cdot \mathbf{v}_k = c_i \|\mathbf{v}_i\|^2,$$

as desired. \square

Definition I.2.22. Orthonormal Basis.

A basis $\beta = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r\} \subset \mathbb{C}^m$ is *orthonormal* if:

- (i) β is orthogonal, and
- (ii) each vector in β is unit, i.e., $\|\mathbf{v}_i\| = 1$.

In the orthonormal basis, the orthogonal projections will, in turn, be:

$$\mathbf{b} = \sum_{k=1}^r (\mathbf{v}_k^* \cdot \mathbf{b}) \mathbf{v}_k.$$

Definition I.2.23. Orthogonal Matrix with $\mathbb{F} = \mathbb{R}$.

A matrix $A \in \mathbb{R}^{m \times m}$ is called *orthogonal* if $A^\top = A^{-1}$, i.e., $A^\top A = \text{Id}_m$ and $AA^\top = \text{Id}_m$.

This can be similarly defined in matrices with complex entries.

Definition I.2.24. Unitary Matrix with $\mathbb{F} = \mathbb{C}$.

A matrix $Q \in \mathbb{C}^{m \times m}$ is called *unitary* if $Q^* = Q^{-1}$, i.e., $Q^* Q = \text{Id}_m$ and $QQ^* = \text{Id}_m$.

Theorem I.2.25. Equivalences with Unitary.

For a matrix $Q \in \mathbb{C}^{m \times m}$, the following conditions are the equivalent:

- (i) Q is unitary, i.e., $Q^* Q = \text{Id}_m$,
- (ii) The columns of $Q = \begin{bmatrix} \mathbf{Q}_1 & \mathbf{Q}_2 & \dots & \mathbf{Q}_m \end{bmatrix}$ are orthonormal, i.e.:
 $\|\mathbf{Q}_i\| = 1$ and $\mathbf{Q}_i^* \mathbf{Q}_j = 0$ for all $i \neq j$,
- (iii) For all $\mathbf{x} \in \mathbb{C}^m$, we have $\|Q \cdot \mathbf{x}\| = \|\mathbf{x}\|$, i.e., the action $f : \mathbf{x} \mapsto Q \cdot \mathbf{x}$ is an *isometry*.

Proof. (i) \implies (ii): Suppose that $Q^* Q = \text{Id}_m$, then we have:

$$\begin{bmatrix} \mathbf{Q}_1^* \\ \mathbf{Q}_2^* \\ \vdots \\ \mathbf{Q}_m^* \end{bmatrix} \cdot \begin{bmatrix} \mathbf{Q}_1 & \mathbf{Q}_2 & \dots & \mathbf{Q}_m \end{bmatrix} = \begin{bmatrix} \mathbf{Q}_1^* \mathbf{Q}_1 & \mathbf{Q}_1^* \mathbf{Q}_2 & \dots & \mathbf{Q}_1^* \mathbf{Q}_m \\ \mathbf{Q}_2^* \mathbf{Q}_1 & \mathbf{Q}_2^* \mathbf{Q}_2 & \dots & \mathbf{Q}_2^* \mathbf{Q}_m \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{Q}_m^* \mathbf{Q}_1 & \mathbf{Q}_m^* \mathbf{Q}_2 & \dots & \mathbf{Q}_m^* \mathbf{Q}_m \end{bmatrix} = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix},$$

which follows along with (ii).

(i) \implies (iii): Still, suppose that $Q^* Q = \text{Id}_m$, then for all $\mathbf{x} \in \mathbb{C}^m$, we have:

$$\|Q \cdot \mathbf{x}\|^2 = (Q \cdot \mathbf{x})^* (Q \cdot \mathbf{x}) = \mathbf{x}^* Q^* Q \mathbf{x} = \mathbf{x}^* (\text{Id}_m) \mathbf{x} = \mathbf{x}^* \mathbf{x} = \|\mathbf{x}\|^2,$$

as desired. \square

Suppose that we want to solve that $Q.\mathbf{x} = \mathbf{b}$, where Q is unitary. If we denote:

$$Q = \begin{bmatrix} \mathbf{Q}_1 & \mathbf{Q}_2 & \cdots & \mathbf{Q}_m \end{bmatrix},$$

we have that:

$$x_1\mathbf{Q}_1 + x_2\mathbf{Q}_2 + \cdots + x_m\mathbf{Q}_m = \mathbf{b},$$

so we have:

$$\mathbf{x} = Q^*Q.\mathbf{x} = Q^*.\mathbf{b}.$$

I.3 Normed Vector Space

Definition I.3.1. Euclidean Length.

Let $\mathbf{x} \in \mathbb{C}^m$, the Euclidean length of \mathbf{x} is:

$$\|\mathbf{x}\| = \sqrt{\mathbf{x} * \mathbf{x}} = \sqrt{\sum_{k=1}^m |x_k|^2}$$

Definition I.3.2. Norm.

A norm is a function $\|\bullet\| : \mathbb{C}^n \rightarrow \mathbb{R}$ with the following properties for all $\mathbf{x}, \mathbf{y} \in \mathbb{C}^m$ and $\alpha \in \mathbb{C}$:

- (i) Positivity: $\|\mathbf{x}\| \geq 0$,
- (ii) Definiteness: $\mathbf{x} = 0 \iff \|\mathbf{x}\| = 0$,
- (iii) Triangle Inequality: $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$, and
- (iv) Homogeneity: $\|\alpha\mathbf{x}\| = |\alpha|\|\mathbf{x}\|$.

Example I.3.3. Examples of Norms.

The 1-norm is $\|\mathbf{x}\|_1 = \sum_{k=1}^m |x_k|$.

The p -norm (where $1 \leq p < \infty$) is $\|\mathbf{x}\|_p = (\sum_{k=1}^m |x_k|^p)^{1/p}$.

The ∞ -norm is $\|\mathbf{x}\|_\infty = \max_{1 \leq k \leq m} |x_k|$.

For a fixed norm, we have unit sphere:

$$S = \{\mathbf{x} : \|\mathbf{x}\| = 1\},$$

and the unit ball as:

$$B = \{\mathbf{x} : \|\mathbf{x}\| \leq 1\}.$$

Example I.3.4. 2-norm (Euclidean Norm) in \mathbb{R}^2 .

The unite sphere of 2-norm would satisfy that $\|\mathbf{x}\|_2 = 1$, or equivalently, $\|\mathbf{x}\|^2 = 1$, hence equivalent to $x_1^2 + x_2^2 = 1$ since we are in the real field. Thus it is the unit circle.

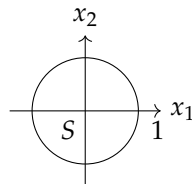


Figure I.2. Unit Sphere in 2-norm.

Example I.3.5. 1-norm in \mathbb{R}^2 .

The unite sphere of 1-norm would satisfy that $\|\mathbf{x}\|_1 = 1$, or equivalently, $\max\{|x_1|, |x_2|\} = 1$ since we are in the real field. Thus it is the diamonds.

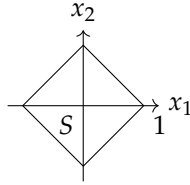
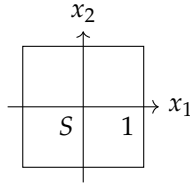


Figure I.3. Unit Sphere in 1-norm.

Example I.3.6. ∞ -norm in \mathbb{R}^2 .

The unite sphere of ∞ -norm would satisfy that $\|\mathbf{x}\|_\infty = 1$, or equivalently, $|x_1| + |x_2| = 1$ since we are in the real field. Thus it is the diamonds.

Figure I.4. Unit Sphere in ∞ -norm.**Remark I.3.7. Note on Invalid p -norms.**

When $0 < p < 1$, $\|\bullet\|_p$ is no longer a norm as it violates triangle inequality.

Example I.3.8. Calculation on Norms.

Let $\mathbf{x} = \begin{bmatrix} 2 \\ 1 \end{bmatrix} \in \mathbb{C}^2$, we have the norms as:

$$\|\mathbf{x}\|_1 = |2| + |1| = 3,$$

$$\|\mathbf{x}\|_2 = \sqrt{|2|^2 + |1|^2} = \sqrt{5},$$

$$\|\mathbf{x}\|_\infty = \max\{|2|, |1|\} = 2.$$

Proposition I.3.9. Monotonicity of Norms.

For any $\mathbf{x} \in \mathbb{C}^n$, $\|\mathbf{x}\|_1 \geq \|\mathbf{x}\|_2 \geq \|\mathbf{x}\|_\infty$.

Remark I.3.10. Stretch of a Norm.

Let $A \in \mathbb{C}^{m \times n}$, consider the function $f : \mathbb{C}^n \rightarrow \mathbb{C}^m$ such that $\mathbf{x} \mapsto A.\mathbf{x}$. The stretch of $\mathbf{x} \neq \mathbf{0}$ caused by multiplication by A is:

$$\frac{\|A.\mathbf{x}\|_*}{\|\mathbf{x}\|_*},$$

where the norms $(\|\bullet\|_*)$ are identical.

Definition I.3.11. Matrix Norms.

Let $A \in \mathbb{C}^{m \times n}$, the matrix norm induced by a vector norm is:

$$\|A\|_{M^{m \times n}} := \sup_{\substack{\mathbf{x} \in \mathbb{C}^n \\ \mathbf{x} \neq \mathbf{0}}} \frac{\|A\mathbf{x}\|_*}{\|\mathbf{x}\|_*}.$$

Since we can restrict the search on the unit circle, which is *compact*, we can use max instead of the sup.

Proof. Let $\mathbf{x} \in \mathbb{C}^n$ be nonzero and arbitrary, we want to show that the stretch of \mathbf{x} and $\mathbf{u} = \mathbf{x}/\|\mathbf{x}\|$, that is $\mathbf{x} = \|\mathbf{x}\|\mathbf{u}$, hence we have:

$$\frac{\|A\mathbf{x}\|}{\|\mathbf{x}\|} = \frac{\|A(\|\mathbf{x}\|\mathbf{u})\|}{\|\mathbf{x}\|} = \frac{\|\mathbf{x}A\mathbf{u}\|}{\|\mathbf{x}\|} = \|A\mathbf{u}\| = \frac{\|A\mathbf{u}\|}{\|\mathbf{u}\|},$$

which implies that they are equivalent. \square

Theorem I.3.12. Equivalence of Matrix Norm.

Let $A \in \mathbb{C}^{m \times n}$, the matrix norm can be computed as:

$$\|A\|_{M^{m \times n}} = \max_{\substack{\mathbf{x} \in \mathbb{C}^n \\ \|\mathbf{x}\|_* = 1}} \|A\mathbf{x}\|_*.$$

Example I.3.13. Computing Matrix Norm.

Let $A = \begin{bmatrix} 1 & 2 \\ 0 & 2 \end{bmatrix}$, we may compute its matrix norm induced by the vector norms for norm $\|\bullet\|_1$, as:

$$\|A\|_1 = \max_{\substack{\mathbf{x} \in \mathbb{C}^n \\ \|\mathbf{x}\|_1 = 1}} \|A\mathbf{u}\|_1,$$

where for all $\mathbf{u} = \begin{bmatrix} u_1 & u_2 \end{bmatrix}^\top$, we have $A\mathbf{u} = u_1\mathbf{A}_1 + u_2\mathbf{A}_2$. Since $\|\mathbf{u}\|_1 = 1$, we have $|u_1| + |u_2| = 1$. Consider $A\mathbf{e}_1 = (1, 0)$ and $A\mathbf{e}_2 = (2, 2)$, we may demonstrate the transformation as:

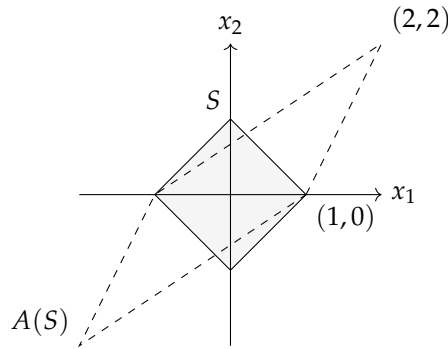


Figure I.5. Transformation of Unit Sphere in 1-norm with Matrix A .

Notice that the maximum of the 1-norm after transformation A is, in fact, $\|A(0,1)\|_1 = \|(2,2)\|_1 = |2| + |2| = 4$. Hence, we have $\|A\| = 4$.

Note that the same computation can be done with 2-norms or other norms, however, the computation will be more lengthy and complicated. For example, $\|A\|_2 \approx 2.9308$ and $\|A\|_\infty = 3$.

Theorem I.3.14. 1-norm of Matrix is Maximum of Norm of Vectors.

For any $A \in \mathbb{C}^{m \times n}$, where $A = [\mathbf{A}_1 \ \mathbf{A}_2 \ \cdots \ \mathbf{A}_n]$. The 1-norm of A is given by:

$$\|A\|_1 = \max_{1 \leq j \leq n} \|\mathbf{A}_j\|_1 = \max_{1 \leq j \leq n} \left(\sum_{k=1}^m |a_{kj}| \right).$$

Observe that this conclusion aligns with the above example, where $\|\mathbf{A}_1\| = 1$ and $\|\mathbf{A}_2\| = 4$, so $\|A\|_1 = 4$.

Proof. Suppose $A \in \mathbb{C}^{m \times n}$, $A = [\mathbf{A}_1 \ \mathbf{A}_2 \ \cdots \ \mathbf{A}_n]$, and $\mathbf{u} = [u_1 \ u_2 \ \cdots \ u_n]^\top$ such that $\|\mathbf{u}\|_1 = 1$, hence we can represent:

$$A \cdot \mathbf{u} = u_1 \mathbf{A}_1 + u_2 \mathbf{A}_2 + \cdots + u_n \mathbf{A}_n.$$

When we take the norm for both sides, we have:

$$\|A \cdot \mathbf{u}\| = \|u_1 \mathbf{A}_1 + u_2 \mathbf{A}_2 + \cdots + u_n \mathbf{A}_n\|.$$

Then, by triangle inequality and properties of norm, we can have:

$$\begin{aligned} \|A \cdot \mathbf{u}\| &\leq \|u_1 \mathbf{A}_1\| + \|u_2 \mathbf{A}_2\| + \cdots + \|u_n \mathbf{A}_n\| \\ &= |u_1| \|\mathbf{A}_1\| + |u_2| \|\mathbf{A}_2\| + \cdots + |u_n| \|\mathbf{A}_n\| \\ &\leq \max_{1 \leq j \leq n} \|\mathbf{A}_j\|. \end{aligned}$$

Hence, for all \mathbf{u} , we have:

$$\|A \cdot \mathbf{u}\| \leq \max_{1 \leq j \leq n} \|\mathbf{A}_j\|_1.$$

Note that the above equality holds when $\mathbf{u} = \mathbf{e}_{j^*}$ where j^* is the largest where $\|A \cdot \mathbf{A}_{j^*}\|$ is the largest of $\|\mathbf{A}_1\|, \dots, \|\mathbf{A}_n\|$. Hence, the maximum over all the unit vectors of $A \cdot \mathbf{u}$ is equal to the maximum of the magnitude of the columns, hence the equality holds trivially. \square

Theorem I.3.15. ∞ -norm of Matrix.

For any $A \in \mathbb{C}^{m \times n}$, where $A = [\mathbf{R}_1 \ \mathbf{R}_2 \ \cdots \ \mathbf{R}_m]^\top$. The ∞ -norm of A is given by:

$$\|A\|_\infty = \max_{1 \leq k \leq m} \|\mathbf{R}_k^*\|_1.$$

Example I.3.16. Calculation of Matrix Norm.

Let $A = \begin{bmatrix} 1 & 2 \\ 0 & 2 \end{bmatrix}$, note that:

$$\|\mathbf{A}_1\| = |1| + |0| = 1,$$

$$\|\mathbf{A}_2\| = |0| + |2| = 2,$$

$$\|\mathbf{R}_1^*\| = |1| + |2| = 3,$$

$$\|\mathbf{R}_2^*\| = |0| + |2| = 2.$$

Hence, $\|A\|_1 = 4$ and $\|A\|_\infty = 3$.

Theorem I.3.17. Cauchy-Schwarz Inequality.

For any $\mathbf{x}, \mathbf{y} \in \mathbb{C}^n$, we have:

$$|\mathbf{x}^* \mathbf{y}| \leq \|\mathbf{x}\|_2 \cdot \|\mathbf{y}\|_2,$$

and the equality holds if and only if they are scalar multiples of each other, *i.e.*, parallel.

Recall that in \mathbb{R}^2 or \mathbb{R}^3 , and any \mathbf{x}, \mathbf{y} in the space, we have:

$$\mathbf{x} \cdot \mathbf{y} = \mathbf{x}^T \mathbf{y} = \|\mathbf{x}\| \|\mathbf{y}\| \cos \theta.$$

Hence, in such case:

$$-1 \leq \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|} \leq 1.$$

Theorem I.3.18. Hölder's Inequality.

For any p and q such that:

$$\frac{1}{p} + \frac{1}{q} = 1,$$

which is called *harmonic conjugates*, then:

$$|\mathbf{x}^* \mathbf{y}| \leq \|\mathbf{x}\|_p \|\mathbf{y}\|_q.$$

Proposition I.3.19. Matrix Norm for Outer Product.

For any $\mathbf{x} \in \mathbb{C}^m$, $\mathbf{y} \in \mathbb{C}^n$, and the matrix $A = \mathbf{x} \mathbf{y}^* \in \mathbb{C}^{m \times n}$. We have:

$$\|A\|_2 = \|\mathbf{x}\| \|\mathbf{y}\|.$$

Proof. Note that by the equivalent definition of norms:

$$\|A\| = \max_{\substack{\|\mathbf{u}\|=1 \\ \mathbf{u} \in \mathbb{C}^n}} \|A \cdot \mathbf{u}\|.$$

Note that by associativity:

$$\|A \cdot \mathbf{u}\| = \|(\mathbf{x} \mathbf{y}^*) \mathbf{u}\| = \|\mathbf{x} (\mathbf{y}^* \mathbf{u})\| = \|(\mathbf{y}^* \mathbf{u}) \mathbf{x}\| = \|\mathbf{y}^* \mathbf{u}\| \|\mathbf{x}\|,$$

and by Cauchy-Schwarz, we have:

$$\|A \cdot \mathbf{u}\| = \|\mathbf{y}^* \mathbf{u}\| \|\mathbf{x}\| \leq \|\mathbf{y}\| \|\mathbf{u}\| \|\mathbf{x}\|,$$

with equality when \mathbf{u} is parallel to \mathbf{y} , thus:

$$\|A\| = \max_{\substack{\|\mathbf{u}\|=1 \\ \mathbf{u} \in \mathbb{C}^n}} \|A \cdot \mathbf{u}\| = \|\mathbf{x}\| \|\mathbf{y}\| \|\mathbf{u}\| = \|\mathbf{x}\| \|\mathbf{y}\|.$$

□

Theorem I.3.20. Inequality for Matrix and Vector Norm.

For any vector norm and induced matrix norm, and for all $A \in \mathbb{C}^{m \times n}$ and $\mathbf{x} \in \mathbb{C}^n$, we have:

$$\|A \cdot \mathbf{x}\| \leq \|A\|_M \|\mathbf{x}\|.$$

Proof. By definition:

$$\|A\| = \max_{\substack{\mathbf{x} \in \mathbb{C}^m \\ \|\mathbf{x}\| \neq 0}} \frac{\|A \cdot \mathbf{x}\|}{\|\mathbf{x}\|},$$

hence for eneric $\mathbf{x} \in \mathbb{C}^m$, we have:

$$\|A\| \geq \frac{\|A \cdot \mathbf{x}\|}{\|\mathbf{x}\|},$$

hence implying that $\|A \cdot \mathbf{x}\| \leq \|A\| \|\mathbf{x}\|$. □

Corollary I.3.21. Inequality for Matrix Multiplication.

For any matrix norm induced by vector norm, and for all $A \in \mathbb{C}^{m \times n}$ and $B \in \mathbb{C}^{n \times p}$, then:

$$\|AB\| \leq \|A\| \|B\|$$

Proof. Note that:

$$\|AB\| = \max_{\substack{\|\mathbf{u}\|=1 \\ \mathbf{u} \in \mathbb{C}^p}} \|AB \cdot \mathbf{u}\|.$$

For the inequality:

$$\|AB \cdot \mathbf{u}\| \leq \|A\| \|B \cdot \mathbf{u}\| \leq \|A\| \|B\| \|\mathbf{u}\|$$

Taking when the equality holds (by Cauchy-Schwarz), the equality holds. □

Remark I.3.22. Recalling Axioms of Normed Vector Space.

For all $A, B \in \mathbb{C}^{m \times n}$, for all $\alpha \in \mathbb{C}$, the following must hold:

- (i) Positivity and definiteness: $\|A\| \geq 0$ and $\|A\| = 0$ if and only if $A = 0$,
- (ii) Triangular inequality: $\|A + B\| \leq \|A\| + \|B\|$, and
- (iii) Homogeneity: $\|\alpha A\| = |\alpha| \|A\|$.

Definition I.3.23. Forbenius norm.

Let $A \in \mathbb{C}^{m \times n}$, its Forbenius norm is defined to be:

$$\|A\|_F = \left(\sum_{i=1}^m \sum_{j=1}^n |a_{i,j}^2| \right)^{1/2}.$$

Theorem I.3.24. Equivalent Definition of Forbenius Norm.

Let $A \in \mathbb{C}^{m \times n}$, its Forbenius norm is:

$$\|A\|_F = \sqrt{\text{Tr}(A^* A)}.$$

Proof. It is easy to find that the diagonals of $A^* A$ are $\mathbf{A}_i^* \mathbf{A}_i$, and since the trace is the sum of the squares, then:

$$\text{Tr}(A^* A) = \sum_{i=1}^n \|\mathbf{A}_i\|^2 = \sum_{i=1}^n |a_{i,1}|^2 + \cdots + \sum_{i=1}^n |a_{i,m}|^2 = \sum_{j=1}^m \sum_{i=1}^n |a_{i,j}|^2,$$

and the sums can be switched since it is finite sum. □

Example I.3.25. Isometry.

The following matrices are isometric:

- (i) Rotation matrices,
- (ii) Reflection matrices, and
- (iii) Permutation of coordinates.

Proposition I.3.26. Invariant of Unitary Matrix.

If $Q \in \mathbb{C}^{m \times n}$ is unitary, and $P \in \mathbb{C}^{n \times m}$, then for all $A \in \mathbb{C}^{m \times n}$, we have:

- (i) Left unit: $\|QA\|_2 = \|A\|_2$ and $\|QA\|_F = \|A\|_F$.
- (ii) Right unit: $\|AP\|_2 = \|A\|_2$ and $\|AP\|_F = \|A\|_F$.

Proof. (i) By definition:

$$\|QA\|_2 = \max_{\substack{\mathbf{u} \in \mathbb{C}^n \\ \|\mathbf{u}\|_2=1}} \|QA \cdot \mathbf{u}\|_2 = \max_{\substack{\mathbf{u} \in \mathbb{C}^n \\ \|\mathbf{u}\|_2=1}} \|Q \cdot (A \cdot \mathbf{u})\|_2 = \max_{\substack{\mathbf{u} \in \mathbb{C}^n \\ \|\mathbf{u}\|_2=1}} \|A \cdot \mathbf{u}\|_2 = \|A\|_2.$$

For Forbenius norm, we have:

$$\|QA\|_F = \sqrt{\text{Tr}((QA)^*(QA))} = \sqrt{\text{Tr}(A^*Q^*QA)} = \sqrt{\text{Tr}(A^*A)} = \|A\|_F.$$

(ii) For the right unitary:

$$\|AP\|_2 = \max_{\substack{\mathbf{u} \in \mathbb{C}^n \\ \|\mathbf{u}\|_2=1}} \|AP \cdot \mathbf{u}\|_2 = \max_{\substack{\mathbf{u} \in \mathbb{C}^n \\ \|\mathbf{u}\|_2=1}} \|A \cdot \mathbf{v}\|_2,$$

where $\mathbf{v} \in \mathbb{C}^m$ and $\|\mathbf{v}\|_2 = \|P\mathbf{u}\|_2 = 1$. Note that the map is valid since P is isometry, hence it is invertible, thus $S \rightarrow S$ is invertible, hence, we have the above to be $\|A\|_2$.

For Forbenius norm, we have:

$$\|AP\|_F = \sqrt{\text{Tr}((AP)^*(AP))} = \sqrt{\text{Tr}(P^*A^*AP)} = \sqrt{\text{Tr}(PP^*A^*A)} = \sqrt{\text{Tr}(A^*A)} = \|A\|_F. \quad \square$$

I.4 Review: Eigenspace**Definition I.4.1. Eigenvalue.**

Let $A \in \mathbb{C}^{m \times m}$ be a square matrix, \mathbf{x} is an eigenvector of A with associate eigenvalue $\lambda \in \mathbb{C}$ such that:

$$\mathbf{x} \neq \mathbf{0} \text{ and } A \cdot \mathbf{x} = \lambda \mathbf{x}.$$

Example I.4.2. Finding Eigenvalue and Eigenvector.

Let $A = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$, we may notice that for $\mathbf{x} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$, we have:

$$A \cdot \mathbf{x} = 2\mathbf{x}.$$

Else, for the rotation matrix $A = \begin{bmatrix} \cos \pi/4 & -\sin \pi/4 \\ \sin \pi/4 & \cos \pi/4 \end{bmatrix}$, and note that there is no real eigenvalues, the complex eigenvalues are:

$$\lambda = e^{\pm i\pi/4}.$$

Theorem I.4.3. Singular \iff 0 is Eigenvalue.

Let A be a m -by- m matrix, A is singular if and only if 0 is an eigenvalue.

Proof. A is singular \iff There exists $\mathbf{x} \neq \mathbf{0}$ such that $A\mathbf{x} = \mathbf{0} \iff A\mathbf{x} = 0\mathbf{x} \iff \lambda = 0$ is eigenvalue. \square

Definition I.4.4. Characteristic Polynomial.

The characteristic polynomial of a matrix $A \in \mathbb{C}^{m \times m}$ is:

$$\det(A - \lambda \text{Id}) = 0.$$

Hence, we may define the polynomial that:

$$\det(A) = \sum_{\sigma \in S_n} \text{sgn}(\sigma) \prod_{i=1}^n a_{i, \sigma(i)},$$

where σ is the permutation, S_m is the n -th cyclic group, and sgn is the sign function of the permutation, i.e., even or odd.

The degree of the characteristic polynomial is the dimension of the square matrix, and by the fundamental theorem of algebra, since \mathbb{C} is algebraically closed, we are guaranteed with full set of complex eigenvalues.

Proposition I.4.5. Roots of Characteristic Polynomial are Eigenvalues.

The roots of the characteristic polynomial are the eigenvalues with respective multiplicity. If an eigenvalue has multiplicity 1, then it is a simple eigenvalue.

Remark I.4.6. Diagonalized Matrices.

When A is diagonal, the eigenvalues are the entries on the diagonal.

Theorem I.4.7. Determinant and Eigenvalues.

$\det A$ is the product of the eigenvalues counted with multiplicity.

Theorem I.4.8. Complex Conjugates in Real Field.

Suppose $A \in \mathbb{R}_{m,m}$ has only real entries, then the coefficients in p_A are also real. If $\lambda = a + ib$ is an eigenvalues, then:

- (i) $\bar{\lambda} = a - ib$ is also an eigenvalue, and
- (ii) $\bar{\lambda}$ has the same multiplicity with λ .

Proof. Note that:

$$\overline{p_A(\lambda)} = \overline{\sum_{k=1}^n p_k \lambda^k} = 0,$$

$$\sum_{k=1}^n p_k \overline{\lambda^k} = 0.$$

Hence $\overline{\lambda}$ is also a root. □

If A is real and symmetric, then we have:

$$A^T = A \text{ and } A^* = A,$$

then we have only real eigenvalues.

Definition I.4.9. Eigenspace.

The eigenspace associated with eigenvalue λ_i is:

$$\ker(A - \lambda_i \text{Id}) = \{x : (A - \lambda_i \text{Id}).x = 0\}.$$

Every vector in this subspace is an eigenvector of A with eigenvalue λ_i . J

Definition I.4.10. Geometric Multiplicity.

The geometric multiplicity of an eigenvalue λ_i is defined as:

$$\dim(\text{eigenspace of } \lambda_i) = \dim(\ker(A - \lambda_i \text{Id})) = \# \text{ of Linearly Independent Eigenvectors of } \lambda_i. \quad J$$

Proposition I.4.11. Monotonicity of Multiplicity.

For each matrix $A \in \mathbb{C}^{n \times n}$ with an eigenvalue λ , we have:

$$1 \leq \text{Geometric Multiplicity}(\lambda) \leq \text{Algebraic Multiplicity}(\lambda).$$

The second equality holds when A is diagonal.

Example I.4.12. Example of Strict Inequality.

Let $A = \begin{bmatrix} 5 & 1 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 5 \end{bmatrix}$, note that 5 is an eigenvalue, we may observe that its algebraic multiplicity is 3.

Consider $A - 5\text{Id} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} := B$, and note that:

$$\ker B = \{(x_1, 0, x_3) : x_1, x_3 \in \mathbb{C}\}.$$

Hence, the geometric multiplicity, we have 2, which is strictly less than the algebraic multiplicity. J

Definition I.4.13. Defective Eigenvalue and Non-Defective Matrix.

If $\text{Geometric Multiplicity}(\lambda) < \text{Algebraic Multiplicity}(\lambda)$, then λ_i is a *defective* eigenvalue.

$A \in \mathbb{C}^{m \times m}$ is *non-defective* if it has m independent eigenvectors, which is equivalently $\text{Geometric Multiplicity}(\lambda) = \text{Algebraic Multiplicity}(\lambda)$ for all eigenvalues of A . J

Proposition I.4.14. Distinct Eigenvalues \implies LI Eigenvectors.

Suppose $A \in \mathbb{C}^{m \times m}$ is a matrix with eigenvectors $\{x_1, \dots, x_k\}$ whose corresponding eigenvalues are $\{\lambda_1, \dots, \lambda_k\}$. If all eigenvalues are distinct, then the eigenvectors are linearly independent.

If $k = m$, then all eigenvectors of A are linearly independent.

The above case does not account for invertible matrix, as zero can be an eigenvalue.

Definition I.4.15. Similar Matrices.

Suppose $A, B \in \mathbb{C}^{m \times m}$, they are similar if $A = SBS^{-1}$ for some invertible matrix S .

Proposition I.4.16. Consequences of Similar Matrices.

Suppose $A, B \in \mathbb{C}^{m \times m}$ are similar, then:

- (i) $\det A = \det B$,
- (ii) $p_A(\lambda) = p_B(\lambda)$, and
- (iii) The set of all eigenvalues of A and B are identical.

Definition I.4.17. Diagonalizable.

$A \in \mathbb{C}^{m \times m}$ is diagonal if it is similar to a diagonal matrix $D \in \mathbb{C}^{m \times m}$, i.e., there exists invertible matrix $S \in \mathbb{C}^{m \times m}$ such that $A = SDS^{-1}$.

Theorem I.4.18. Diagonalizable \iff Non-defective.

Suppose $A \in \mathbb{C}^{m \times m}$. A is diagonalizable if and only if A has m linearly independent eigenvectors, i.e., A is non-defective.

Proof. Note that $A = SDS^{-1}$ is equivalent to $AS = SD$, hence equivalent to:

$$\begin{bmatrix} AS_1 & AS_2 & \cdots & AS_m \end{bmatrix} = \begin{bmatrix} d_{1,1}S_1 & d_{2,2}S_2 & \cdots & d_{m,m}S_m \end{bmatrix}.$$

Hence, we equivalently have $A \cdot S_i = d_{i,i}S_i$, so S_i is a set of eigenvector of A with eigenvalues $d_{i,i}$. Since S is invertible, all columns of S are linearly independent. \square

Corollary I.4.19. Distinct Eigenvalues \implies Diagonalizable.

If $A \in \mathbb{C}^{m \times m}$ has m distinct eigenvalues, then A is diagonalizable.

This corollary is an immediate consequence of Diagonalizable \iff Non-defective.

The Complex Spectral Theorem can be generalized to the real Spectral Theorem.

Proposition I.4.20. Kernel of A and A^*A are Same.

For any $A \in \mathbb{C}^{m \times n}$, we have:

$$\ker A = \ker(A^*A).$$

Proof. We first show that $\ker A \subset \ker(A^*A)$. Suppose $\mathbf{x} \in \ker A$, then $A\mathbf{x} = \mathbf{0}$, so $A^*A\mathbf{x} = \mathbf{0}$, hence $\mathbf{x} \in \ker(A^*A)$.

For the other inclusion, we suppose $\mathbf{x} \in \ker(A^*A) = 0$, then $\mathbf{x}^*A^*A\mathbf{x} = \mathbf{x}^*\mathbf{0}$. By collecting the terms, we have $(A\mathbf{x})^*(A\mathbf{x}) = 0$, which implies that $\|A\mathbf{x}\|^2 = 0$, and by the axiom of normed vector space, we have $A\mathbf{x} = \mathbf{0}$, which implies that $\mathbf{x} \in \ker A$. \square

I.5 Diagonalization and Singular Value Decomposition**Remark I.5.1. Matrices as Stretches.**

Let $A \in \mathbb{C}^{n \times m}$ be a matrix with complex entries. The map $T : \mathbb{C}^n \rightarrow \mathbb{C}^m$ which for any $\mathbf{x} \in \mathbb{C}^n$, having $\mathbf{x} \mapsto A\mathbf{x}$, deforms a unit circle $S = \{\mathbf{x} : \|\mathbf{x}\|_2 = 1\}$ into an ellipse.

Definition I.5.2. Singular Values of Matrix.

Let $A \in \mathbb{C}^{n \times m}$ be a matrix with complex entries. We let \mathbf{v}_1 be the vector such that:

$$\|A\mathbf{v}_1\| = \max_{\substack{\mathbf{u} \in \mathbb{C}^n \\ \|\mathbf{u}\|_2=1}} \|A\mathbf{u}\|_2,$$

and we let \mathbf{u}_1 be the unit vector of v_1 , and the singular value σ_1 satisfies that:

$$\sigma_1 = \|A\mathbf{v}_1\|.$$

There also exists $\mathbf{v}_2 \in \mathbb{C}^2$ such that $\mathbf{v}_2 \perp \mathbf{v}_1$ and $A\mathbf{v}_2 \perp A\mathbf{v}_1$, and we define the other singular value σ_2 :

$$\sigma_2 = \|A\mathbf{v}_2\|.$$

Proposition I.5.3. Properties of Full Rank 2-by-2 Matrices.

Let A be a 2-by-2 matrix with full rank, there exists $\mathbf{v}_1, \mathbf{v}_2 \in \mathbb{R}^2$ with the following properties:

- (i) $\beta = \{\mathbf{v}_1, \mathbf{v}_2\}$ is an orthonormal basis of \mathbb{R}^2 .
- (ii) If we define $A\mathbf{v}_1 = \sigma_1\mathbf{u}_1$ and $A\mathbf{v}_2 = \sigma_2\mathbf{u}_2$, we have $\sigma_1 \geq \sigma_2 > 0$ and $A\mathbf{v}_1 \perp A\mathbf{v}_2$.

In particular, we can rewrite it as:

$$\begin{bmatrix} A\mathbf{v}_1 & A\mathbf{v}_2 \end{bmatrix} = \begin{bmatrix} \sigma_1\mathbf{u}_1 & \sigma_2\mathbf{u}_2 \end{bmatrix},$$

which results in:

$$A \underbrace{\begin{bmatrix} \mathbf{v}_1 & \mathbf{v}_2 \end{bmatrix}}_V = \underbrace{\begin{bmatrix} \mathbf{u}_1 & \mathbf{u}_2 \end{bmatrix}}_U \underbrace{\begin{bmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{bmatrix}}_\Sigma.$$

Since U and V are unitary, we have:

$$A = U\Sigma V^*.$$

In general, for matrix $A \in \mathbb{C}^{m \times n}$ with $\dim(\text{im } A) = r$, we can write:

$$\begin{array}{ccccccc} A & = & U & \circ & \Sigma & \circ & V. \\ \mathfrak{m} & & \mathfrak{m} & & \mathfrak{m} & & \mathfrak{m} \\ \mathbb{C}^{m \times n} & & \mathbb{C}^{m \times m} & & \mathbb{C}^{m \times n} & & \mathbb{C}^{n \times n} \end{array}$$

Proposition I.5.4. Properties of $\bullet^* \bullet$ for Matrices.

For any matrix $A \in \mathbb{C}^{m \times n}$, with $\dim(\text{im } A) = r$, we have:

- (i) A^*A is n -by- n matrix and Hermitian (self-adjoint), i.e., $(A^*A)^* = A^*A$,
- (ii) $\dim(\text{im}(A^*A)) = r$,
- (iii) A^*A has r nonzero eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0$, while $\lambda_{r+1} = \lambda_{r+2} = \dots = \lambda_n = 0$.
- (iv) A^*A and AA^* have the same nonzero eigenvalues.

Proof. (i) Trivial.

(ii) Concerning the rank-nullity theorem, we have:

$$\dim(\text{im}(A^*A)) = n - \dim(\ker(A^*A)) = n - \dim(\ker A) = \dim(\text{im } A).$$

(iii) Consider the eigenspace for $\lambda = 0$, we have:

$$E_0 = \ker(A - 0 \cdot \text{Id}) = \ker A.$$

Then, we have the dimension of E_0 as $n - r$, so A has r nonzero eigenvalues. Suppose $\lambda \neq 0$ is an eigenvalue of A^*A , so there exists $\mathbf{x} \neq \mathbf{0}$ such that $A^*A\mathbf{x} = \lambda\mathbf{x}$, then we have $\mathbf{x}^*A^*A\mathbf{x} = \lambda\mathbf{x}^*\mathbf{x}$, so we have $\|A\mathbf{x}\|^2 = \lambda\|\mathbf{x}\|^2$, so:

$$\lambda = \frac{\|A\mathbf{x}\|^2}{\|\mathbf{x}\|^2} > 0.$$

(iv) Suppose that $\lambda \neq 0$ is an eigenvalue of A^*A , then there exists $\mathbf{x} \neq \mathbf{0}$ such that $A^*A\mathbf{x} = \lambda\mathbf{x}$, so:

$$AA^*A\mathbf{x} = \lambda A\mathbf{x},$$

hence the eigenvalue of AA^* is λ and the eigenvector is $A\mathbf{x}$ (since \mathbf{x} is eigenvector and $\lambda \neq 0$, $A\mathbf{x} \neq \mathbf{0}$, otherwise it is a contradiction). \square

Proposition I.5.5. Single Value Decomposition.

Suppose $A \in \mathbb{C}^{m \times n}$ such that $\dim(\text{im } A) = r \leq \min\{m, n\}$ in which $A : \mathbb{C}^n \rightarrow \mathbb{C}^m$, then there exists:

- (i) an orthonormal basis $\{\mathbf{v}_1, \dots, \mathbf{v}_r, \underbrace{\mathbf{v}_{r+1}, \dots, \mathbf{v}_n}_{\text{basis for } \ker A}\}$ of \mathbb{C}^n ,
- (ii) another orthonormal basis $\{\underbrace{\mathbf{u}_1, \dots, \mathbf{u}_r}_{\text{basis for } \text{im } A}, \mathbf{u}_{r+1}, \dots, \mathbf{u}_m\}$ of \mathbb{C}^m , and
- (iii) r singular values, with $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$,

such that:

(i) $U^*U = UU^* = \text{Id},$

(ii) $V^*V = VV^* = \text{Id},$ and

(iii) $\Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 & 0 \\ 0 & \ddots & 0 & 0 \\ 0 & 0 & \sigma_r & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}.$

Remark I.5.6. Single Value for Full Rank.

Suppose A is square and *full rank*, it is $r \times r$ and $\ker A = \{\mathbf{0}\}$, $\text{im } A = \mathbb{C}^r$. Moreover, they have the bases $\beta = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r\}$ and $\gamma = \{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r\}$, while the Σ matrix is:

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_r \end{bmatrix}$$

is square, diagonal, and full rank.

Theorem I.5.7. Spectral Theorem.

Suppose $B \in \mathbb{C}^{n \times n}$ is Hermitian (i.e., $B^* = B$), then B has n real eigenvalues and it is orthogonally diagonalizable, i.e., B has n orthonormal eigenvectors $\beta = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$, such that, if we define:

$$S := \begin{bmatrix} \mathbf{v}_1 & \cdots & \mathbf{v}_n \end{bmatrix}$$

we have:

$$S^*AS = D = \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_n \end{bmatrix}.$$

Hence $A = SDS^*$.

Remark I.5.8. Remark on Real Spectral Theorem.

For the case above, we want to construct unitary $U \in \mathbb{C}^{m \times m}$, $V \in \mathbb{C}^{m \times m}$, and $\Sigma \in \mathbb{C}^{m \times n}$ that is diagonalizable such that $A = U\Sigma V^*$. The sketch of the proof would be solving for U, Σ, V , where we first compute:

$$A^*A = (U\Sigma V^*)^*(U\Sigma V^*) = V\Sigma^*U^*U\Sigma V^*,$$

hence implying that:

$$A^*A = V(\Sigma^*\Sigma)V^*,$$

which implies that:

$$V = \begin{bmatrix} \mathbf{v}_1 & \cdots & \mathbf{v}_n \end{bmatrix}$$

which gives:

$$\Sigma^* \Sigma = \begin{bmatrix} \sigma_1^2 & 0 & 0 & 0 \\ 0 & \ddots & 0 & 0 \\ 0 & 0 & \sigma_r^2 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

Hence, the singular values are $\sigma_i = \sqrt{\lambda_i}$ with $1 \leq i \leq r$.

For the process, we diagonalize A^*A with respect to orthonormal basis $\{v_i\}$ and take v_i 's to form columns of V and compute $\sigma_i \sqrt{\lambda_i}$ to form Σ .

Then, we choose u_i 's via:

$$u_i = \frac{A \cdot v_i}{\sigma_i} \text{ for } i = 1, 2, \dots, r.$$

For the remaining u_{r+1}, \dots, u_m are chosen to be orthonormal to u_1, \dots, u_m , using Gram Schmidt process.

Example I.5.9. Finding SVD.

Let $A = \begin{bmatrix} 1 & -1 \\ -1 & 1 \\ -4 & 4 \end{bmatrix}$, we note that $m = 3$, $n = 2$, and $r = 1$ since there is only one linearly independent column.

(i) First, we find A^*A , that is:

$$A^*A = \begin{bmatrix} 18 & -18 \\ -18 & 18 \end{bmatrix},$$

and we note that eigenvalues as:

$$A^*A \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \text{ and } A^*A \cdot \begin{bmatrix} 1 \\ -1 \end{bmatrix} = 36 \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

(ii) Now, since we have the eigenvalues as 36 and 0, whose eigenvectors are $v_1 = 1/\sqrt{2} \begin{bmatrix} 1 \\ -1 \end{bmatrix}$ and

$v_2 = 1/\sqrt{2} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$, then the first singular value is $\sigma_1 = 6$, then we have:

$$u_1 = \frac{A v_1}{\sigma_1} = \frac{1}{3\sqrt{2}} \begin{bmatrix} 1 \\ -1 \\ -4 \end{bmatrix}$$

(iii) Then, we look for $u_2 \perp u_3$, in which we want:

$$\frac{1}{3\sqrt{2}} \begin{bmatrix} y_1 & y_2 & y_3 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \\ -4 \end{bmatrix} = 0,$$

which returns to:

$$y_1 = y_2 + 4y_3$$

As of right now, we have:

$$y = \begin{bmatrix} y_2 + 4y_3 \\ y_2 \\ y_3 \end{bmatrix} = y_2 \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} + y_3 \begin{bmatrix} 4 \\ 0 \\ 1 \end{bmatrix}.$$

Note that we want them to be orthogonal, that is:

$$\begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \quad \text{and} \quad \begin{bmatrix} 2 \\ -2 \\ 1 \end{bmatrix}.$$

Technically, we should use Graham Schmidt to find an orthonormal basis, but this can be observed easily.

Proposition I.5.10. Block Matrix Multiplication.

Suppose $A \in \mathbb{C}^{m \times n}$ and $B \in \mathbb{C}^{n \times p}$, we can break:

$$\begin{aligned} m &= \sum_{i=1}^{\alpha} m_i, \\ n &= \sum_{i=1}^{\beta} n_i, \\ p &= \sum_{i=1}^{\gamma} p_i. \end{aligned}$$

In particular, we can write:

$$AB = \begin{bmatrix} A_{1,1} & \cdots & A_{1,\beta} \\ \vdots & \ddots & \vdots \\ A_{\alpha,1} & \cdots & A_{\alpha,\beta} \end{bmatrix} \begin{bmatrix} B_{1,1} & \cdots & B_{1,\gamma} \\ \vdots & \ddots & \vdots \\ B_{\beta,1} & \cdots & B_{\beta,\gamma} \end{bmatrix},$$

where computation is distributive as if they are scalar entries.

There is an alternative way to compute $\gamma = \{\mathbf{u}_1, \dots, \mathbf{u}_r, \mathbf{u}_{r+1}, \dots, \mathbf{u}_m\}$, hence:

$$A = U\Sigma V^* \iff AV = U\Sigma V^*V = U\Sigma.$$

The right hand side can be simplified into:

$$U\Sigma = \begin{bmatrix} \mathbf{u}_1 & \mathbf{u}_2 & \cdots & \mathbf{u}_m \end{bmatrix} \begin{bmatrix} \sigma_1 & 0 & 0 & 0 \\ 0 & \ddots & 0 & 0 \\ 0 & 0 & \sigma_r & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} \sigma_1 \mathbf{u}_1 & \sigma_2 \mathbf{u}_2 & \cdots & \sigma_m \mathbf{u}_m \end{bmatrix}.$$

For the left hand side, we have:

$$AV = \begin{bmatrix} A\mathbf{v}_1 & A\mathbf{v}_2 & \cdots & A\mathbf{v}_r & A\mathbf{v}_{r+1} & \cdots & A\mathbf{v}_n \end{bmatrix}.$$

Hence, for the first r columns, we have:

$$\mathbf{u}_i = \frac{A\mathbf{v}_i}{\sigma_i} \text{ for } i = 1, 2, \dots, r,$$

where as for the last vectors, they are the kernel of the map, and must be set to be an orthonormal set that is also orthogonal to prior entries.

Theorem I.5.11. Properties about Matrix Norms.

For any $m \times n$ matrix A , we have:

- (i) $\|A\|_2 = \sigma_1$, i.e., the largest singular value (as σ_i 's are positive and ordered from large to small), and
- (ii) $\|A\|_F = \sqrt{\sigma_1^2 + \sigma_2^2 + \cdots + \sigma_r^2}$.

Recall that $\|\bullet\|_2$ and $\|\bullet\|_F$ are invariant by multiplication by unitary matrices.

Proof. (i) Note that:

$$\begin{aligned} \|A\|_2 &= \|U\Sigma V^*\|_2 = \|U\Sigma\|_2 = \|\Sigma\|_2 \\ &= \max_{\substack{\mathbf{u} \in \mathbb{C}^n \\ \|\mathbf{u}\|_2=1}} \|\Sigma \cdot \mathbf{u}\|_2 = \max_{\substack{\mathbf{u} \in \mathbb{C}^n \\ \|\mathbf{u}\|_2=1}} \sqrt{|u_1|^2 \sigma_1^2 + |u_2|^2 \sigma_2^2 + \cdots + |u_r|^2 \sigma_r^2} = \sqrt{|1|^2 \sigma_1^2} = \sigma_1. \end{aligned}$$

(ii) For Forbenius norm, we have:

$$\begin{aligned} \|A\|_F &= \|U\Sigma V^*\|_F = \|U\Sigma\|_F = \|\Sigma\|_F \\ &= \sqrt{|\sigma_1|^2 + |\sigma_2|^2 + \cdots + |\sigma_r|^2} = \sqrt{\sigma_1^2 + \sigma_2^2 + \cdots + \sigma_r^2}. \end{aligned} \quad \square$$

Proposition I.5.12. Hermitian \implies Orthogonally Diagonalizable.

Suppose A is Hermitian, then A is orthogonally diagonalizable.

Moreover, the singular values of A are $\sigma_i = |\lambda_i|$, where λ_i 's are the *ordered* (by absolute value) of eigenvalues of A .

Proof. Let $A = SDS^*$, in particular:

$$\begin{aligned} A &= \underbrace{\begin{bmatrix} \mathbf{v}_1 & \mathbf{v}_2 & \cdots & \mathbf{v}_n \end{bmatrix}}_{\text{orthonormal set of eigenvectors of } A} \underbrace{\begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_n \end{bmatrix}}_{\text{Not necessarily SVD, as eigenvalues can be non-positive}} \begin{bmatrix} \mathbf{v}_1^* \\ \mathbf{v}_2^* \\ \vdots \\ \mathbf{v}_n^* \end{bmatrix} \\ &= \underbrace{\begin{bmatrix} \text{sgn}(\lambda_1)\mathbf{v}_1 & \text{sgn}(\lambda_2)\mathbf{v}_2 & \cdots & \text{sgn}(\lambda_n)\mathbf{v}_n \end{bmatrix}}_U \underbrace{\begin{bmatrix} |\lambda_1| & 0 & \cdots & 0 \\ 0 & |\lambda_2| & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & |\lambda_n| \end{bmatrix}}_{\Sigma} \underbrace{\begin{bmatrix} \mathbf{v}_1^* \\ \mathbf{v}_2^* \\ \vdots \\ \mathbf{v}_n^* \end{bmatrix}}_V, \end{aligned}$$

where the sign function is defined to be:

$$\text{sgn}(x) = \begin{cases} 1 & \text{if } x \geq 0, \\ -1 & \text{if } x < 0. \end{cases}$$

□

Now, as we consider the SVD for $A \in \mathbb{C}^{m \times n}$, where $\dim(\text{im } A) = r$, we have:

$$A = \begin{bmatrix} \overbrace{\mathbf{u}_1 \mathbf{u}_2 \cdots \mathbf{u}_r}^{r \text{ columns}} & \overbrace{\mathbf{u}_{r+1} \cdots \mathbf{u}_m}^{m-r \text{ columns}} \end{bmatrix} \begin{bmatrix} \overbrace{\begin{matrix} \sigma_1 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_r & 0 & \cdots & 0 \\ 0 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & 0 & \cdots & 0 \end{matrix}}^{\tilde{\Sigma}} \end{bmatrix} \begin{bmatrix} \left. \begin{matrix} \mathbf{v}_1^* \\ \vdots \\ \mathbf{v}_r^* \end{matrix} \right\} r \text{ rows} \\ \left. \begin{matrix} \mathbf{v}_{r+1}^* \\ \vdots \\ \mathbf{v}_n^* \end{matrix} \right\} n-r \text{ rows} \end{bmatrix}$$

Computationally, we find that:

- (i) $\mathbf{v}_1, \dots, \mathbf{v}_r, \mathbf{v}_{r+1}, \dots, \mathbf{v}_n$ are orthonormal eigenvectors of A^*A with respective eigenvalues, and
- (ii) $\mathbf{u}_1, \dots, \mathbf{u}_r, \mathbf{u}_{r+1}, \dots, \mathbf{u}_m$ are orthonormal eigenvectors of AA^* with respective eigenvalues.

Recall the block computation, we have:

$$A = \begin{bmatrix} \tilde{U} & U_2 \end{bmatrix} \begin{bmatrix} \tilde{\Sigma} & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \tilde{V}^* \\ V_2^* \end{bmatrix} = \begin{bmatrix} \tilde{U}\tilde{\Sigma} + U_2 0 & \tilde{U} 0 + U_2 0 \end{bmatrix} \begin{bmatrix} \tilde{V}^* \\ V_2^* \end{bmatrix} = \tilde{U}\tilde{\Sigma}\tilde{V}^*.$$

Definition I.5.13. Reduced Singular Value Decomposition.

With the construction above, with $A \in \mathbb{C}^{m \times n}$ being Hermitian, we have reduced Singular Value Decomposition, that is:

$$\begin{array}{ccccccc} A & = & \tilde{U} & \circ & \tilde{\Sigma} & \circ & \tilde{V} \\ \mathfrak{M} & & \mathfrak{M} & & \mathfrak{M} & & \mathfrak{M} \\ \mathbb{C}^{m \times n} & & \mathbb{C}^{m \times r} & & \mathbb{C}^{r \times r} & & \mathbb{C}^{r \times n} \end{array}$$

Note that with reduced SVD, we have:

$$A = \begin{bmatrix} \sigma_1 \mathbf{u}_1 & \sigma_2 \mathbf{u}_2 & \cdots & \sigma_r \mathbf{u}_r \end{bmatrix} = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^* + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^* + \cdots + \sigma_r \mathbf{u}_r \mathbf{v}_r^*,$$

which is called the *rank-one* decomposition of A . Hence A is the sum of rank-one matrices with weights $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_r > 0$.

Remark I.5.14. Approximating a Matrix.

By choosing a $p \leq r$, we have the rank- p approximation of A such that

$$\widehat{A}_p = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^* + \cdots + \sigma_p \mathbf{u}_p \mathbf{v}_p^*,$$

where we have $\dim(\text{im } \widehat{A}_p) = p$.

Example I.5.15. Using SVD to Compress Image(s).

Suppose we use a 400-by-400 matrix to represent a gray-scale image, i.e., $A \in [0, 255]^{400 \times 400}$, that is:

$$A = \begin{bmatrix} A_{1,1} & A_{1,2} & \cdots & A_{1,400} \\ A_{2,1} & A_{2,2} & \cdots & A_{2,400} \\ \vdots & \vdots & \ddots & \vdots \\ A_{400,1} & A_{400,2} & \cdots & A_{400,400} \end{bmatrix},$$

where each entry represents the gray-scale in that pixel, where 0 represents black and 255 represents white, and the grays are within $(0, 255)$, while becoming lighter as the number increases.

Consider the reduced SVD, we may write A into:

$$A = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^* + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^* + \cdots + \sigma_r \mathbf{u}_{400} \mathbf{v}_{400}^*,$$

in which we only consider the first p terms, we are able to reduce the rank in order to store less data for the image. In this way, for rank p , we only need to store $(1 + 400 + 400) \times p$ values rather than everything.

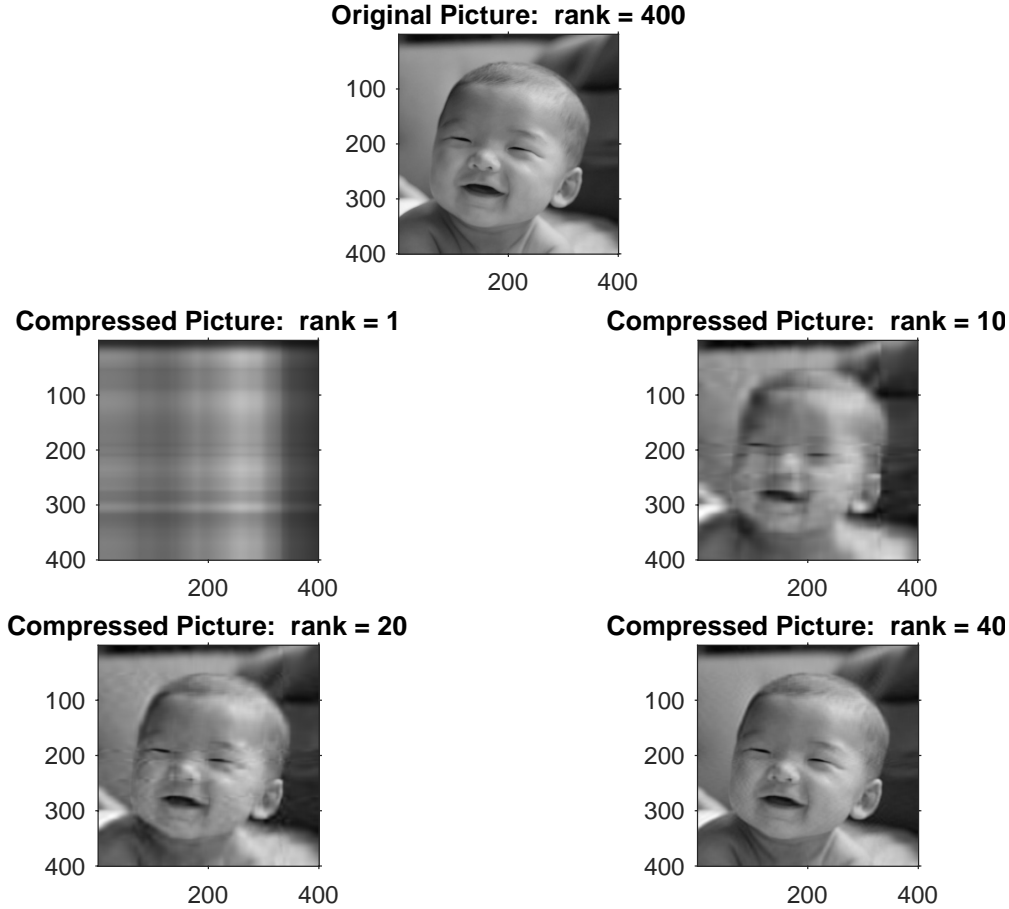


Figure I.6. Compressing an image of full rank 400 into lower ranks.

With the example image above, we note that when rank is 10, the image is roughly recognizable, where as at rank 40, the image is clear enough.

I.6 Projection

Definition I.6.1. Idempotent.

For matrix $P \in \mathbb{C}^{m \times m}$ is *idempotent* or *projector* if $P^2 = P$.

Hence, for all $\mathbf{x} \in \mathbb{C}^m$, we have $P.\mathbf{x} = P^2.\mathbf{x}$, i.e., $P.\mathbf{x} = P(P.\mathbf{x})$.

Definition I.6.2. Orthogonal Projector on Vector.

For a fixed nonzero $\mathbf{v} \in \mathbb{C}^m$, we define that orthogonal projector onto \mathbf{v} as:

$$P_{\mathbf{v}} = \frac{1}{\|\mathbf{v}\|^2} \mathbf{v} \mathbf{v}^* = \frac{\mathbf{v} \mathbf{v}^*}{\mathbf{v}^* \mathbf{v}}.$$

Geometrically, we may represent the projector as:

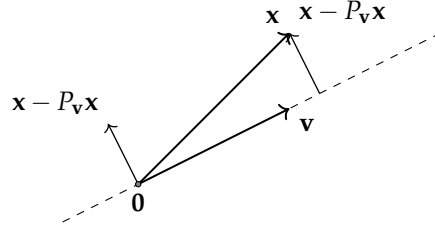


Figure 1.7. Geometric feature of the Orthogonal Projector.

Theorem I.6.3. Geometric Properties with Orthogonal Projection.

For all $\mathbf{x} \in \mathbb{C}^m$, we have:

- (i) $P_{\mathbf{v}}\mathbf{x} \in \text{span}\{\mathbf{v}\}$, i.e., $\mathbf{y} \parallel \mathbf{v}$, and
- (ii) $\mathbf{x} - P_{\mathbf{v}}\mathbf{x} \perp \mathbf{v}$.

Proof. (i) For the projector, we have:

$$P_{\mathbf{v}}\mathbf{x} = \frac{1}{\|\mathbf{v}\|^2} \mathbf{v} \mathbf{v}^* \mathbf{x} = \frac{\mathbf{v}^* \mathbf{x}}{\|\mathbf{v}\|^2} \mathbf{v} = k \mathbf{v} = \text{span}\{\mathbf{v}\}.$$

(ii) For the inner product:

$$\mathbf{v}^* (\mathbf{x} - P_{\mathbf{v}}\mathbf{x}) = \mathbf{v}^* \mathbf{x} - \mathbf{v}^* \left(\frac{1}{\|\mathbf{v}\|^2} \mathbf{v} \mathbf{v}^* \mathbf{x} \right) = \mathbf{v}^* \mathbf{x} - \frac{1}{\|\mathbf{v}\|^2} \mathbf{v}^* \mathbf{v} \mathbf{v}^* \mathbf{x} = 0.$$

□

Proposition I.6.4. Orthogonal Projector is Idempotent.

Let $P_{\mathbf{v}}$ be the orthogonal projector onto \mathbf{v} , $P_{\mathbf{v}}^2 = P_{\mathbf{v}}$.

Proof. We may deduce that:

$$\begin{aligned} P_{\mathbf{v}}^2 &= \left(\frac{1}{\|\mathbf{v}\|^2} \mathbf{v} \mathbf{v}^* \right) \left(\frac{1}{\|\mathbf{v}\|^2} \mathbf{v} \mathbf{v}^* \right) = \frac{1}{\|\mathbf{v}\|^4} \mathbf{v} \mathbf{v}^* \mathbf{v} \mathbf{v}^* \\ &= \frac{\mathbf{v}^* \mathbf{v}}{\|\mathbf{v}\|^4} \mathbf{v} \mathbf{v}^* = \frac{1}{\|\mathbf{v}\|^2} \mathbf{v} \mathbf{v}^* = P_{\mathbf{v}}. \end{aligned}$$

□

Recall that if $\beta = \{\mathbf{v}_1, \dots, \mathbf{v}_r\}$ is an orthogonal basis of a subspace $V \subset \mathbb{C}^m$, then for all $\mathbf{x} \in V$, we can write:

$$\mathbf{x} = c_1 \mathbf{v}_1 + c_2 \mathbf{v}_2 + \dots + c_r \mathbf{v}_r,$$

where each $c_i = \frac{\mathbf{v}_i^* \mathbf{x}}{\|\mathbf{v}_i\|^2}$ for all $i = 1, 2, \dots, r$.

Only when we have an orthogonal basis, we have the following sum:

$$\mathbf{x} = \frac{\mathbf{v}_1^* \mathbf{x}}{\|\mathbf{v}_1\|^2} \mathbf{v}_1 + \frac{\mathbf{v}_2^* \mathbf{x}}{\|\mathbf{v}_2\|^2} \mathbf{v}_2 + \dots + \frac{\mathbf{v}_r^* \mathbf{x}}{\|\mathbf{v}_r\|^2} \mathbf{v}_r = P_{\mathbf{v}_1} \mathbf{x} + P_{\mathbf{v}_2} \mathbf{x} + \dots + P_{\mathbf{v}_r} \mathbf{x}.$$

However, on non orthogonal basis, this might not be true.

Let $\{\mathbf{v}_1, \mathbf{v}_2\}$ be an orthogonal basis and let $\{\mathbf{w}_1, \mathbf{w}_2\}$ be a non-orthogonal basis in \mathbb{R}^2 , we can present the following projections of $\mathbf{x} \in \mathbb{R}^2$.

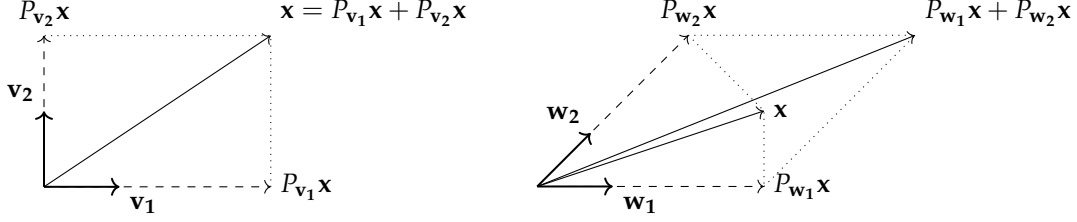


Figure I.8. Projection with orthogonal basis (left) and non-orthogonal basis (right).

Proposition I.6.5. Complementary Projector is Project.

If $P \in \mathbb{C}^{m \times m}$ is a projector, then $\text{Id} - P$ is also a projector.

Proof. Note that:

$$(\text{Id} - P)^2 = (\text{Id} - P)(\text{Id} - P) = \text{Id}^2 - P - P + P^2 = \text{Id}^2 - P - P + P = \text{Id}^2 - P.$$

□

Remark I.6.6. Kernel and Image of Complementary Projector.

Consider a orthogonal projector on a line $\mathbf{v} \neq \mathbf{0}$, and for $\mathbf{v} \in \mathbb{C}^m$, with $P_{\mathbf{v}} = \mathbf{v}\mathbf{v}^* / \|\mathbf{v}\|^2$, so we have:

$$(\text{Id} - P)\mathbf{x} = \mathbf{x} - P_{\mathbf{v}}\mathbf{x}.$$

Recall the kernel and image of the projector:

$$\text{im}(P_{\mathbf{v}}) = \text{span}\{\mathbf{v}\},$$

$$\ker(P_{\mathbf{v}}) = \{\mathbf{w} : \mathbf{v}^*\mathbf{w} = 0\} = (\text{span}\{\mathbf{v}\})^\perp.$$

However, for the complementary projector, we have:

$$\text{im}(\text{Id} - P_{\mathbf{v}}) = \{\mathbf{w} : \mathbf{v}^*\mathbf{w} = 0\} = \ker(P_{\mathbf{v}}),$$

$$\ker(\text{Id} - P_{\mathbf{v}}) = \text{span}\{\mathbf{v}\} = \text{im}(P_{\mathbf{v}}).$$

For any *Idempotent* matrices $P^2 = P$, we first find $\ker P$ such that:

$$\mathbf{w} = \mathbf{x} - P.\mathbf{x},$$

We find that:

$$P.\mathbf{w} = P.\mathbf{x} - P^2.\mathbf{x} = P.\mathbf{x} - P.\mathbf{x} = \mathbf{0}.$$

Hence, $\mathbf{w} \in \ker P$.

Also, we have $\mathbf{w} = (\text{Id} - P).\mathbf{x} \in \text{im}(\text{Id} - P)$.

Theorem I.6.7. Idempotent \implies Image and Kernel Relation.

If $P = P^2$, then:

$$\text{im}(\text{Id} - P) = \ker P.$$

Moreover:

$$\text{im } P = \ker(\text{Id} - P).$$

Proof. (i) ($\text{im}(\text{Id} - P) \subseteq \ker P$;) We let $\mathbf{w} \in \text{im}(\text{Id} - P)$ be generic, then there exists \mathbf{x} such that:

$$\mathbf{w} = (\text{Id} - P).\mathbf{x} = \mathbf{x} - P.\mathbf{x},$$

hence:

$$P.\mathbf{w} = P.\mathbf{x} - P^2.\mathbf{x} = P.\mathbf{x} - P.\mathbf{x} = \mathbf{0}.$$

Therefore, $w \in \ker P$.

($\text{im}(\text{Id} - P) \supseteq \ker P$;) Let $\mathbf{w} \in \ker P$ be generic, i.e., $P.\mathbf{w} = \mathbf{0}$, then we have:

$$(\text{Id} - P).\mathbf{w} = \mathbf{w} - P.\mathbf{w},$$

which results in $\mathbf{w} = (\text{Id} - P).\mathbf{w}$, so $w \in \text{im}(\text{Id} - P)$, as desired.

(ii) Let $Q = \text{Id} - P$, by the previous part, we have:

$$\text{im}(\text{Id} - Q) = \ker Q,$$

hence, by $\text{Id} - Q = \text{Id} - \text{Id} + P = P$, hence:

$$\text{im } P = \ker(\text{Id} - Q).$$

□

Proposition I.6.8. Idempotent \implies Intersection of Image and Kernel is Trivial.

Let $P = P^2$ be idempotent:

$$\text{im } P \cap \ker P = \{\mathbf{0}\}.$$

Proof. We let $\mathbf{w} \in \text{im } P \cap \ker P$ be trivial, then:

$$\begin{cases} \mathbf{w} \in \text{im } P = \ker(\text{Id} - P), \\ \mathbf{w} \in \ker P. \end{cases}$$

This implies that:

$$\begin{cases} (\text{Id} - P).\mathbf{w} = \mathbf{0}, \\ P.\mathbf{w} = \mathbf{0}, \end{cases}$$

which implies that $\mathbf{w} - P.\mathbf{w} = \mathbf{0}$, so $\mathbf{w} = \mathbf{0}$.

□

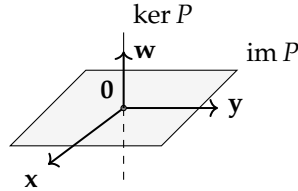


Figure I.9. Kernel and image of P as idempotent.

Note that suppose $\dim(\text{im } P) = r$, by rank-nullity, we have $\dim(\ker P) = m - r$, hence all dimensions are

captured by the image and the kernel, that can be represented by direct sum that:

$$\text{im } P \oplus \ker P = \mathbb{C}^m.$$

Thus, any $\mathbf{x} \in \mathbb{C}^m$ can be uniquely decomposed into:

$$\mathbf{x} = \mathbf{y} + \mathbf{w},$$

where $\mathbf{y} \in \text{im } P$ and $\mathbf{w} \in \ker P$.

To compute the \mathbf{y} and \mathbf{w} above, we have:

$$\mathbf{y} \in \text{im } P = \ker(\text{Id} - P),$$

hence resulting in:

$$(\text{Id} - P).\mathbf{y} = \mathbf{0},$$

hence $\mathbf{y} = P.\mathbf{y}$, and by acting P on \mathbf{x} , we have:

$$P.\mathbf{x} = P.\mathbf{y} + P.\mathbf{w} = \mathbf{y},$$

hence resulting in:

$$\begin{cases} \mathbf{y} = P.\mathbf{x}, \\ \mathbf{w} = \mathbf{x} - \mathbf{y}. \end{cases}$$

Remark I.6.9. Eigenvalues and Diagonalizability.

Assume $P = P^2$ is idempotent:

- (i) For all $\mathbf{x} \in \ker P$, we have $P.\mathbf{x} = \mathbf{0}$, hence $P.\mathbf{x} = 0.\mathbf{x}$.
- (ii) For all $\mathbf{x} \in \text{im } P = \ker(\text{Id} - P)$, we have $(\text{Id} - P).\mathbf{x} = \mathbf{0}$, so $\mathbf{x} = P.\mathbf{x}$ hence $P.\mathbf{x} = 1.\mathbf{x}$.
- (iii) Hence, for $\lambda = 0$, the eigenvalue has geometric multiplicity of r , and the eigenvector has geometric multiplicity of $m - r$.

In conclusion, idempotent implies diagonalizability. J

Proposition I.6.10. Gram-Schmidt Process.

Let $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_p\}$ be a basis of a subspace $V \subset \mathbb{C}^m$, the Gram-Schmidt Process allows us to construct an orthogonal basis $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_p\}$ such that $\text{span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_p\} = \text{span}\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_p\}$.

Proof. The construction of the Gram-Schmidt is as follows:

- (i) Let $\mathbf{u}_1 = \mathbf{v}_1$.
- (ii) We construct \mathbf{u}_2 that is orthogonal to \mathbf{v}_1 , let:

$$\mathbf{u}_2 = \mathbf{v}_2 - \text{component of } \mathbf{v}_2 \text{ that is parallel to } \mathbf{u}_1 = \mathbf{v}_2 - \frac{\mathbf{u}_1^* \mathbf{v}_2}{\|\mathbf{u}_1\|^2} \mathbf{u}_1.$$

Note that readers can verify that $\text{span}\{\mathbf{u}_1, \mathbf{u}_2\} = \text{span}\{\mathbf{v}_1, \mathbf{v}_2\}$.

- (iii) For \mathbf{u}_3 we think of it as:

$$\mathbf{u}_3 = \mathbf{v}_3 - \frac{\mathbf{u}_1^* \mathbf{v}_3}{\|\mathbf{u}_1\|^2} \mathbf{u}_1 - \frac{\mathbf{u}_2^* \mathbf{v}_3}{\|\mathbf{u}_2\|^2} \mathbf{u}_2.$$

This can be illustrated in \mathbb{R}^3 here, as follows:

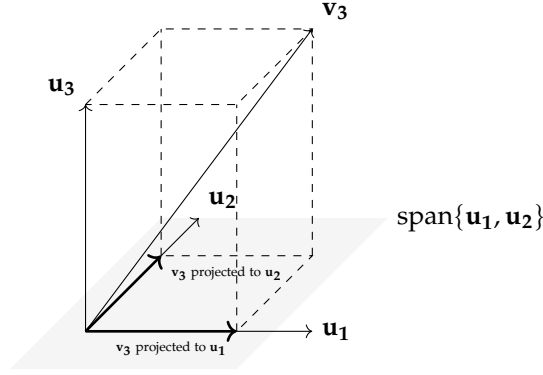


Figure I.10. Subtracting the orthogonal projections gives the orthogonal component.

(k) At step k , we have:

$$\mathbf{u}_k = \mathbf{v}_k - \sum_{j=1}^{k-1} \frac{\mathbf{u}_j^* \mathbf{v}_k}{\|\mathbf{u}_j\|^2} \mathbf{u}_j.$$

Hence, we have completed the construction for finitely dimensional basis. \square

Remark I.6.11. 3rd Step of Gram Schmidt.

Consider the built:

$$Q_2 = [\mathbf{q}_1 \quad \mathbf{q}_2]$$

Here, the projector on $\text{span}\{\mathbf{q}_1, \mathbf{q}_2\}$ has:

$$T_2 = Q_2 Q_2^*$$

which, by block matrix multiplication, leads to:

$$T_2 \mathbf{A}_3 = Q_2 Q_2^* \mathbf{A}_3 = (\mathbf{q}_1^* \mathbf{A}_3) \mathbf{q}_1 + (\mathbf{q}_2^* \mathbf{A}_3) \mathbf{q}_2,$$

Hence, we have:

$$\mathbf{u}_3 = \mathbf{A}_3 - T_2 \mathbf{A}_3 = (\text{Id} - T_2) \mathbf{A}_3,$$

so we can define:

$$P_3 = \text{Id} - T_2, \text{ so } \mathbf{u}_3 = P_3 \mathbf{A}_3, \text{ hence } \mathbf{q}_3 = \frac{\|\mathbf{u}_3\|}{\mathbf{u}_3}.$$

Theorem I.6.12. For Idempotent, Hermitian $\iff \text{im } P \perp \ker P$.

Suppose $P = P^2$ is idempotent:

$$P^* = P \iff \ker P \perp \text{im } P.$$

Proof. (\implies .) Note that for any linear map P , $\ker P \perp \text{im}(P^*)$, hence $P = P^*$ implies $\ker P \perp \text{im } P$.

(\impliedby .) We choose:

$$\beta = \{\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_r\},$$

$$\gamma = \{\mathbf{q}_{r+1}, \mathbf{q}_{r+1}, \dots, \mathbf{q}_m\},$$

as the orthonormal basis of $\text{im } P$ and $\ker P$, respectively. From assumption that $\ker P \perp \text{im } P$, then $\beta \perp \gamma$, so we have:

$$P \cdot \mathbf{q}_1 = \mathbf{q}_1, P \cdot \mathbf{q}_2 = \mathbf{q}_2, \dots, P \cdot \mathbf{q}_r = \mathbf{q}_r, P \cdot \mathbf{q}_{r+1} = \mathbf{0}, \dots, P \cdot \mathbf{q}_m = \mathbf{0}.$$

hence, our matrix can be represented as:

$$\begin{bmatrix} P.\mathbf{q}_1 & P.\mathbf{q}_2 & \cdots & P.\mathbf{q}_r & P.\mathbf{q}_{r+1} & \cdots & P.\mathbf{q}_m \end{bmatrix} = \begin{bmatrix} \mathbf{q}_1 & \cdots & \mathbf{q}_r & \mathbf{0} & \cdots & \mathbf{0} \end{bmatrix}$$

Hence, we have:

$$Q^*PQ = \begin{bmatrix} \mathbf{q}_1^* \\ \vdots \\ \mathbf{q}_r^* \\ \mathbf{q}_{r+1}^* \\ \vdots \\ \mathbf{q}_m^* \end{bmatrix} \begin{bmatrix} \mathbf{q}_1 & \cdots & \mathbf{q}_r & \mathbf{0} & \cdots & \mathbf{0} \end{bmatrix} = \begin{bmatrix} 1 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 1 & 0 & \cdots & 0 \\ 0 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & 0 & \cdots & 0 \end{bmatrix},$$

in which the top $r \times r$ rows are non-trivial. Thus:

$$Q^*PQ = D,$$

hence $P = QDQ^*$, so $P^* = (QDQ^*)^* = (Q^*)^*D^*Q^* = QDQ^* = P$. \square

In fact, we may simplify the above prove by noticing orthogonality. Since the eigenvalues are non-negative, we have $P = QDQ^*$ being the SVD of P where $U = Q$, $\Sigma = D$ and $V = Q$.

If we are considering the parts of the *nonzero* entries of the matrices, we have $P = \tilde{Q}\tilde{\Sigma}\tilde{Q}^* = \tilde{Q}\text{Id}\tilde{Q}^* = \tilde{Q}\tilde{Q}^*$. Note that:

$$\tilde{Q} = \begin{bmatrix} \mathbf{q}_1 & \cdots & \mathbf{q}_r \end{bmatrix}$$

is the orthonormal basis of the range, hence inducing a rank-one decomposition, that is:

$$P = \mathbf{q}_1\mathbf{q}_1^* + \mathbf{q}_2\mathbf{q}_2^* + \cdots + \mathbf{q}_r\mathbf{q}_r^*.$$

Recall that due to the projection, we have:

$$P_{\mathbf{q}_i} = \frac{\mathbf{q}_i\mathbf{q}_i^*}{\|\mathbf{q}_i\|^2} = \mathbf{q}_i\mathbf{q}_i^*.$$

Again, we would ask the question. *What if the basis we have for $\text{im } P$ is not orthonormal?*

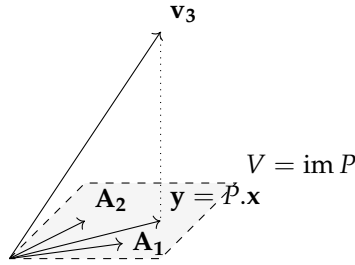


Figure I.11. Projections when the basis is not orthogonal.

In such case, we have:

$$V = \text{span}\{\mathbf{A}_1, \mathbf{A}_2, \cdots, \mathbf{A}_r\},$$

which are linearly independent but not necessarily orthonormal. Hence, we want to find an $m \times n$ matrix P with the property that for all $\mathbf{x} \in \mathbb{C}^m$:

$$\mathbf{y} = P.\mathbf{x} \in V \setminus \backslash \mathbf{w} = \mathbf{x} - \mathbf{y} \perp V.$$

Proposition I.6.13. Matrix with full rank has Hermitian compose itself Invertible.

For A being a $m \times r$ matrix with $\dim(\text{im } A) = r$, A^*A is invertible.

Definition I.6.14. Pseudo Inverse of a Matrix.

Suppose A is $m \times r$ and $\dim(\text{im } A) = r$, we want to solve that:

$$A.x = \mathbf{b} \text{ with } \mathbf{b} \in \text{im } A.$$

Hence, for $A^*Ax = A^*\mathbf{b}$ since A^*A is square and invertible, then $\mathbf{x} = (A^*A)^{-1}A^*\mathbf{b}$, and we have:

$$A^+ = (A^*A)^{-1}A^*$$

as the pseudo-inverse of A .

Theorem I.6.15. Pseudo Inverse Conditions.

For a matrix $P \in \mathbb{C}^{m \times m}$ that satisfies for all $\mathbf{x} \in \mathbb{C}^m$ that:

(i) $\mathbf{y} = P.\mathbf{x} \in V$, and

(ii) $\mathbf{w} = \mathbf{x} - \mathbf{y} \perp V$,

is $P = A(A^*A)^{-1}A^*$ for A being the matrix composed of the linearly independent vectors.

Proof. By (i), for any $\mathbf{x} \in \mathbb{C}^m$, we want $\mathbf{y} = P.\mathbf{x} \in V$, i.e.:

$$\mathbf{y} = P.\mathbf{x} = c_1\mathbf{A}_1 + c_2\mathbf{A}_2 + \cdots + c_r\mathbf{A}_r = \underbrace{\begin{bmatrix} \mathbf{A}_1 & \cdots & \mathbf{A}_r \end{bmatrix}}_A \underbrace{\begin{bmatrix} c_1 \\ \vdots \\ c_r \end{bmatrix}}_{\mathbf{c}} = A.\mathbf{c}.$$

By (ii), we have:

$$\begin{aligned} \mathbf{x} - P.\mathbf{x} \perp \mathbf{A}_1 &\implies \mathbf{A}_1^*(\mathbf{x} - P.\mathbf{x}) = 0 \\ \vdots &\implies \vdots \\ \mathbf{x} - P.\mathbf{x} \perp \mathbf{A}_r &\implies \mathbf{A}_r^*(\mathbf{x} - P.\mathbf{x}) = 0 \end{aligned}$$

Hence, we can get:

$$\begin{bmatrix} \mathbf{A}_1^* \\ \vdots \\ \mathbf{A}_r^* \end{bmatrix} . (\mathbf{x} - P.\mathbf{x}) = \mathbf{0}.$$

Hence, we have $A^*(\mathbf{x} - P.\mathbf{x}) = \mathbf{0}$, then $A^*(\mathbf{x} - A.\mathbf{c}) = \mathbf{0}$, xp $A^*.\mathbf{x} = A^*A\mathbf{c}$, which leads to $\mathbf{c} = (A^*A)^{-1}A^*.\mathbf{x}$, so:

$$P.\mathbf{x} = A.\mathbf{c} = A(A^*A)^{-1}A^*\mathbf{x} \text{ for all } \mathbf{x}.$$

□

Remark I.6.16. Special Cases with Projector.

If $\{\mathbf{A}_1, \dots, \mathbf{A}_r\}$ are orthonormal, then $A^*A = r \text{Id}$ and $P = AA^*$.

If $V = \text{span}\{\mathbf{v}\}$, so $A = \begin{bmatrix} \mathbf{v} \end{bmatrix}$, then:

$$P = \mathbf{v}(\mathbf{v}^* \mathbf{v})^{-1} \mathbf{v}^* = \frac{\mathbf{v} \mathbf{v}^*}{\|\mathbf{v}\|^2},$$

which the orthogonal projection on a line.

I.7 QR Decomposition

Proposition I.7.1. Gram Schmidt for Normalization.

Let $\beta = \{\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_n\}$ be a basis, the Gram-Schmidt process, gives $\gamma = \{\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_n\}$, so we have:

$$\begin{aligned} \mathbf{u}_1 &= \mathbf{A}_1 \longrightarrow \mathbf{q}_1 = \frac{\mathbf{u}_1}{\|\mathbf{u}_1\|} \\ \mathbf{u}_2 &= \mathbf{A}_2 - \frac{\mathbf{u}_1^* \mathbf{A}_2}{\|\mathbf{u}_1\|} \mathbf{u}_1 \longrightarrow \mathbf{q}_2 = \frac{\mathbf{u}_2}{\|\mathbf{u}_2\|} \\ \mathbf{u}_3 &= \mathbf{A}_3 - \frac{\mathbf{u}_1^* \mathbf{A}_3}{\|\mathbf{u}_1\|} \mathbf{u}_1 - \frac{\mathbf{u}_2^* \mathbf{A}_3}{\|\mathbf{u}_2\|} \mathbf{u}_2 \longrightarrow \mathbf{q}_3 = \frac{\mathbf{u}_3}{\|\mathbf{u}_3\|}. \end{aligned}$$

For this part, our goal is to have given a matrix $A \in \mathbb{C}^{m \times n}$, we want to factor it as $A = QR$, where $Q \in \mathbb{C}^{m \times m}$ is unitary and $R \in \mathbb{C}^{m \times n}$ is upper triangular, i.e.:

$$\underbrace{\begin{bmatrix} \mathbf{A}_1 & \cdots & \mathbf{A}_n \end{bmatrix}}_{\substack{A \\ \cap \\ \mathbb{C}^{m \times n}}} = \underbrace{\begin{bmatrix} \mathbf{q}_1 & \cdots & \mathbf{q}_n & \mathbf{q}_{n+1} & \cdots & \mathbf{1}_n \end{bmatrix}}_{\substack{Q \\ \cap \\ \mathbb{C}^{m \times m}}} \circ \underbrace{\begin{bmatrix} * & * & \cdots & * \\ 0 & * & \cdots & * \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & * \\ 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 \end{bmatrix}}_{\substack{R \\ \cap \\ \mathbb{C}^{m \times n}}}.$$

Remark I.7.2. Reprise to Gram Schmidt.

We can first use Gram Schmidt to make a basis orthogonal, and normality is trivial as normalizing in the induced normed vector space is straightforward.

Since with orthonormal basis, we have $\mathbf{v} \in \text{span}(\beta)$ as for all $\mathbf{A}_i \in \text{span} \beta$ that:

$$\mathbf{A}_i = \sum_{j=1}^n (\mathbf{q}_j^* \mathbf{A}_j) \mathbf{q}_j = \sum_{j=1}^i (\mathbf{q}_j^* \mathbf{A}_j) \mathbf{q}_j + 0 + \cdots + 0 = (\mathbf{q}_1^* \mathbf{A}_1) \mathbf{q}_1 + (\mathbf{q}_2^* \mathbf{A}_2) \mathbf{q}_2 + \cdots + (\mathbf{q}_i^* \mathbf{A}_i) \mathbf{q}_i.$$

Using the block multiplication, we have:

$$A = \begin{bmatrix} \mathbf{A}_1 & \cdots & \mathbf{A}_n \end{bmatrix} = \begin{bmatrix} \mathbf{q}_1 & \cdots & \mathbf{q}_n \end{bmatrix} \begin{bmatrix} \mathbf{q}_1^* \mathbf{A}_1 & \mathbf{q}_2^* \mathbf{A}_2 & \mathbf{q}_3^* \mathbf{A}_3 & \cdots & \mathbf{q}_n^* \mathbf{A}_n \\ 0 & \mathbf{q}_2^* \mathbf{A}_2 & \mathbf{q}_3^* \mathbf{A}_3 & \cdots & \mathbf{q}_n^* \mathbf{A}_n \\ 0 & 0 & \mathbf{q}_3^* \mathbf{A}_3 & \cdots & \mathbf{q}_n^* \mathbf{A}_n \\ 0 & 0 & 0 & \cdots & \mathbf{q}_n^* \mathbf{A}_n \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & \mathbf{q}_n^* \mathbf{A}_n \end{bmatrix}$$

Here, we may write this as the pseudo code that:

```

input: A_1, ..., A_n
output: q_1, ..., q_n

for j = 1, ..., n:
    u_j := A_j
    for i = 1, ..., j - 1: % There is nothing for j = 1
        r_ij := q_i^* A_j
        u_j := u_j - r_ij q_i
        % The above line has risk of Catastrophic Cancellation when subtracting similar number.
        % There may be large rounding of error.
    end
    q_j := u_j / ||u_j|| % Marked Step
    r_jj := a_j^* A_j

```

Remark I.7.3. Case when rank is less.

Let A has rank $r < n$, we should be pick arbitrary unit vectors orthogonal to the other basis at the marked step to keep going with Gram Schmidt.

We get $\beta = \{\mathbf{q}_1, \dots, \mathbf{q}_n\}$ as a orthonormal set:

- when A is full rank, $\text{im } A = \text{span}(\beta)$, and
- in any case $\text{im } A \subset \text{span}(\beta)$.

Hence we have:

$$A = \tilde{Q}\tilde{R} = \begin{bmatrix} \mathbf{q}_1 & \cdots & \mathbf{q}_n \end{bmatrix} \tilde{R}.$$

- We want to find $\{\mathbf{q}_{n+1}, \dots, \mathbf{1}_n\}$ as orthonormal set orthogonal to β , with
- \tilde{R} having $m - n$ zero rows.

In particular, the (classical) Gram Schmidt will introduce large rounding error compounding up. Here, we can introduce the modified Gram Schmidt for computing purposes.

Proposition I.7.4. Modified Gram Schmidt.

Given $\{\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_n\}$, columns of $A = \begin{bmatrix} \mathbf{A}_1 & \mathbf{A}_2 & \cdots & \mathbf{A}_n \end{bmatrix}$

- (i) Initialize $\mathbf{u}_1 = \mathbf{A}_1$ and $\mathbf{q}_1 = \mathbf{u}_1 / \|\mathbf{u}_1\|$.

(ii) Have $\mathbf{u}_2^{(1)} = \mathbf{A}_2$, with:

$$\mathbf{u}_2 = \mathbf{u}_2^{(2)} = P_{\mathbf{q}_1} \mathbf{u}_2^{(1)} = (\text{Id} - \mathbf{q}_1 \mathbf{q}_1^*) \mathbf{u}_2^{(1)} = \mathbf{u}_2^{(1)} - \mathbf{q}_1 (\mathbf{q}_1^* \mathbf{u}_2^{(1)}) = \mathbf{u}_2^{(1)} - (\mathbf{q}_1^* \mathbf{u}_2^{(1)}) \mathbf{q}_1.$$

(j) For the j -th step, we have:

$$\begin{aligned} \mathbf{u}_j^{(1)} &= \mathbf{A}_j, \\ \mathbf{u}_j^{(2)} &= \mathbf{u}_j^{(1)} - (\mathbf{q}_1^* \mathbf{u}_j^{(1)}) \mathbf{q}_1, \\ &\vdots \\ \mathbf{u}_j &= \mathbf{u}_j^{(j)} = \mathbf{u}_{j-1}^{(j-1)} - (\mathbf{q}_{j-1}^* \mathbf{u}_j^{(j-1)}) \mathbf{q}_{j-1}. \end{aligned}$$

```

for i = 1, ..., n:    % first loop to initialize all u_i's
    u_i := A_i

for i = 1, ..., n:    % calculating each step
    r_ii := ||u_i||
    q_i := u_i / ||u_i||
    for j = i + 1, ..., n:
        r_ij := q_i^* * u_j
        u_j := u_j - r_ij * q_i

```

Definition I.7.5. Computational Complexity.

The complexity is measured by FLOPs: each operation is considered to operate with a unit time.

Example I.7.6. Complexity of Inner Product.

Recall $r_{ij} := q_i^* u_j$ is an inner product. In general, we consider $\mathbf{x}^* \mathbf{y}$ for any $\mathbf{x}, \mathbf{y} \in \mathbb{C}^m$, hence there will be m multiplications and m subtractions, hence it has $2m$.

Hence, if we consider the second iterated loop with j with the above pseudo-code, we have the total number of operations being:

$$\sum_{i=1}^n \sum_{j=i+1}^n 4m = 4m \sum_{i=1}^n \sum_{j=i+1}^n 1 = 4m \cdot \frac{n^2 - n}{2} \sim 2mn^2.$$

If $A \in \mathbb{C}^{m \times m}$ is Hermitian, i.e., $A^* = A$, then we have the orthogonal diagonalization that $Q^* A Q = D$, where Q is unitary. However, we want to have the decomposition for the more general A , that is:

Definition I.7.7. Householder Triangularization.

If $A \in \mathbb{C}^{m \times n}$ in which $m \geq n$, and sometimes we require $\dim(\text{im } A) = n$, then we have orthogonal

triangularization that:

$$Q^*A = \begin{bmatrix} * & * & \cdots & * \\ 0 & * & \cdots & * \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & * \\ 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 \end{bmatrix} = R,$$

where $Q \in \mathbb{C}^{m \times m}$ is unitary. (This is equivalent with $A = QQ^*A = Q$, which is just *QR decomposition*).

The idea is that let:

$$Q^* = Q_n Q_{n-1} \cdots Q_2 Q_1,$$

with each Q_k being unitary, we have:

$$Q = (Q^*)^* = Q_1^* Q_2^* \cdots Q_n^*.$$

Example I.7.8. A 5×3 Matrix with Householder Triangularization.

Let A be a 5 by 3 that:

$$A = \begin{bmatrix} * & * & * \\ * & * & * \\ * & * & * \\ * & * & * \\ * & * & * \end{bmatrix},$$

so by 3 steps of each Q_1 , Q_2 , and Q_3 to obtain that:

$$A = \begin{bmatrix} * & * & * \\ * & * & * \\ * & * & * \\ * & * & * \\ * & * & * \end{bmatrix} \xrightarrow{Q_1} \begin{bmatrix} * & * & * \\ 0 & * & * \\ 0 & * & * \\ 0 & * & * \\ 0 & * & * \end{bmatrix} \xrightarrow{Q_2} \begin{bmatrix} * & * & * \\ 0 & * & * \\ 0 & 0 & * \\ 0 & 0 & * \\ 0 & 0 & * \end{bmatrix} \xrightarrow{Q_3} \begin{bmatrix} * & * & * \\ 0 & * & * \\ 0 & 0 & * \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix},$$

where Q_1 changes entirely the first column, Q_2 changes all except the first row for the second column, and Q_3 changes all except the first two rows for the third column.

Recall that unitary ($Q^* = Q^{-1}$) $\iff \|Q \cdot \mathbf{x}\| = \|\mathbf{x}\|$ for all \mathbf{x} , which is isometry. Hence, the structure of Q_n is:

$$Q_k = \begin{bmatrix} \text{Id} & 0 \\ 0 & F \end{bmatrix}$$

so we have:

$$\begin{bmatrix} \text{Id} & 0 \\ 0 & F \end{bmatrix} \cdot \begin{bmatrix} T & B \\ 0 & X \end{bmatrix} = \begin{bmatrix} T & B \\ 0 & FX \end{bmatrix}.$$

Here, we want Q_k to be unitary, that is having orthonormal columns. So we need columns of F to be orthonormal, so equivalently, F must be an isometry on $\mathbb{C}^{m-(k-1)}$.

Here, we take $\mathbf{x} = \begin{bmatrix} * \\ \vdots \\ * \end{bmatrix} \in \mathbb{C}^{m-(k-1)}$, so we have $F \cdot \mathbf{x} = \begin{bmatrix} \|\mathbf{x}\| \\ \vdots \\ 0 \end{bmatrix}$, so we have:

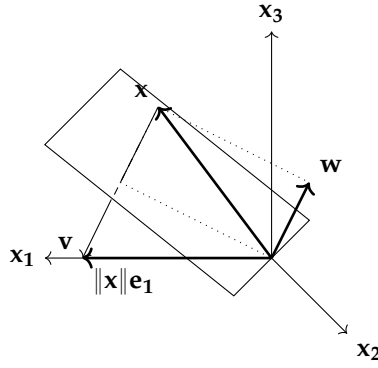


Figure I.12. Illustration of Householder reflection.

Hence, we have that $F.x - x = \|x\|e_1 - x$, hence, we have:

$$w = \text{orthogonal projection of } x \text{ onto } v = P_v x = \frac{vv^*}{\|v\|^2} x,$$

so we have that $v = -2w$. Therefore, we have:

$$F.x = x - 2w = \text{Id}.x - 2 \frac{vv^*}{\|v\|^2} x,$$

hence we have:

$$F = \text{Id} - 2 \frac{vv^*}{\|v\|^2}.$$

Theorem I.7.9. Properties of Householder Reflection.

For the F being a Householder Reflection, we have the following properties:

- (i) F is unitary,
- (ii) F is Hermitian,
- (iii) $F^{-1} = F$ or $F^2 = \text{Id}$.

Remember that we have $Q_k = \begin{bmatrix} \text{Id} & 0 \\ 0 & F \end{bmatrix}$, we have:

- (i) Q_k is unitary,
- (ii) Q_k Hermitian, and
- (iii) Q_k is involuntary.

Remark I.7.10. Computational Complexity for Householder Triangularization.

We could have defined:

$$F.x = -\|x\|e_i,$$

we choose $F.x$ that is farthest away from x , and the trick is choose $F.x = -\text{sgn}(x_1)\|x\|e_1$.

So we have computational cost is $\sim 2mn^2 - \frac{2}{3}n^3$.

Example I.7.11. Householder QR Algorithm.

For the algorithm, we do:

```

for k = 1 ... n:
    x = A_k:m,k
    v_k = sgn(x_i) * ||x|| * e_1 + x
    v_k = v_k / ||v_k||
    A_k:m,k:n = A_k:m,k:n - 2v_k(v_k^* * A_k:m,k:n)

```

In particular, we have:

$$A^{(k)} = Q_k A^{(k-1)} = Q_k Q_{k-1} \cdots Q_1 A,$$

and:

$$A^{(n)} = R = Q_n Q_{n-1} \cdots Q_1 A.$$

We want to solve that $A.x = \mathbf{b}$, which is equivalently $QR.x = \mathbf{b}$, that is $R.x = Q^*.\mathbf{b}$.

To calculate $Q^*.\mathbf{b}$, we use:

```

for k = 1 ... n:
    b_k:m = b_k:m - 2v_kv_k^* * b_k:m

```

And to compute $Q.\mathbf{u}$, we use:

```

for k = n ... 1:
    u_k:m = u_k:m' - 2v_kv_k^* * u_k:m.

```

In particular $A_a : b, c : d$ meaning to obtain the a th to b th row and c th column to d th column of A .

Remark I.7.12. Operation Count for Householder QR Decomposition.

Operation count is dominated by innermost for loop, that is the last line in Householder algorithm.

Note that $A_k:m, j$ has length $l = m - (k - 1)$, lets do operator count in terms of l , that is:

- $A_k:m, j - 2v_k(v_k^* A_k:m, j)$.
- Dot Product: l multiplications and $l - 1$ additions.
- scalar multiplication: l multiplications.
- subtraction: l subtractions.

Hence, there are a total of $4l - 1$ flops. This means that we do approximately 4 flops for every entry operated on.

Then to count the total number of entries. For the k th step, we have:

$$\begin{aligned}
 \sum_{k=1}^n (m - (k-1))(n - (k-1)) &= \sum_{k=1}^n (mn + (-m-n)(k-1) + (k-1)^2) \\
 &= mn^2 - m - n \sum_{k=1}^n (k-1) + \sum_{k=1}^n (k-1)^2 \\
 &= mn^2 + (-m-n) \frac{1}{2}n(n+1) + \frac{1}{2}n(n+1)(2n+1) \\
 &\approx mn^2 - \frac{1}{2}mn^2 - \frac{1}{2}n^3 + \frac{1}{3}n^3 = \frac{1}{2}mn^2 - \frac{1}{6}n^3.
 \end{aligned}$$

By multiplying 4, we have $2mn^2 - \frac{2}{3}n^3$. (Recall that this was $2mn^2$ for Gram-Schmidt process.)

II Applications with Computer Programming

II.1 MATLAB Preliminaries

The MATLAB programs can be helpful in conducting computations, and its embedded arrays allow linear algebra computations.

```
x = (-128:128)'/128;      % Create a column from -128 to 128 (inclusively) and normalize
A = [x.^0 x.^1 x.^2 x.^3]; % Create Matrix with each column being from -1 to 1 with 257 steps.
[Q,R] = qr(A,0);          % QR factorization

scale = Q(257,:);         % Calculate the scale
Q = Q*diag(1./scale);     % Modify Q via diagonal matrix
plot(x,Q);                % Plot x against Q.
```

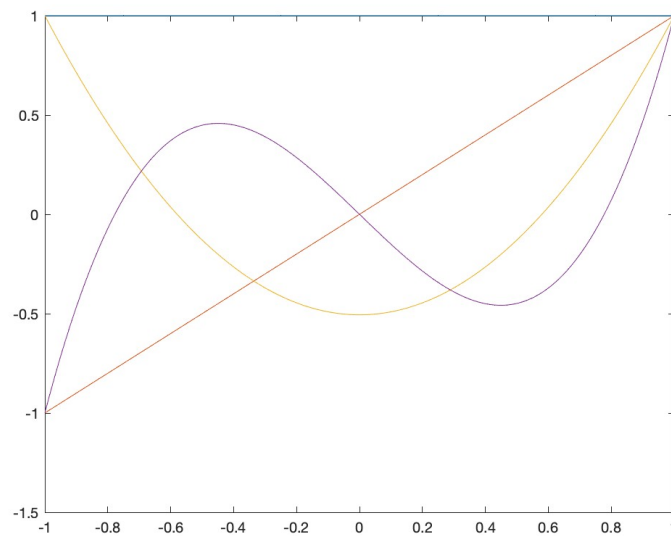


Figure II.1. MATLAB plot on the above code snippet.

Example II.1.1. Computation Error of Gram Schmidt.

First, we have the matrix as:

$$A = \begin{bmatrix} 0.70000 & 0.70711 \\ 0.70001 & 0.70711 \end{bmatrix}.$$

By keeping 5 digits, we have:

$$Q = \begin{bmatrix} 0.70710 & 1.0000 \\ 0.80811 & 0.0000 \end{bmatrix},$$

which is not very accurate.

Of course, we may implement the (traditional) Gram Schmidt algorithm through MATLAB function. This function on MATLAB takes in a matrix A and returns the factorization of Q and R .

```

function [Q,R] = clgs(A)
n = length(A);
R = zeros(n);
Q = zeros(n);
V = zeros(n);
    for j = 1:n
        V(:,j) = A(:,j);
        for i = 1:j-1
            R(i,j) = Q(:,i)'*A(:,j);
            V(:,j) = V(:,j) - R(i,j)*Q(:,i);
        end
        R(j,j) = norm(V(:,j),2);
        Q(:,j) = V(:,j)./R(j,j);
    end
end

```

II.2 Representation of Numbers

Consider the base 10 representation of 273, we have:

$$(273)_{10} = 2 \times 10^2 + 7 \times 10^1 + 3 \times 10^0.$$

Definition II.2.1. Base 2 Representation.

In base 2 representation, all numbers are represented by 0's and 1's.

For example, consider 100101 in binary, we have:

$$(100101)_2 = 1 \times 2^5 + 0 \times 2^4 + 0 \times 2^3 + 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0 = 32 + 4 + 1 = 37.$$

Example II.2.2. Converting Base 10 to Base 2.

Consider $(156)_{10}$ and we want to make it base 2, we have:

$$\begin{aligned}
 156 &= 78 \times 2 + 0, \\
 78 &= 39 \times 2 + 0, \\
 39 &= 19 \times 2 + 1, \\
 19 &= 9 \times 2 + 1, \\
 9 &= 4 \times 2 + 1, \\
 4 &= 2 \times 2 + 0, \\
 2 &= 1 \times 2 + 0, \\
 1 &= 0 \times 2 + 1.
 \end{aligned}$$

Hence, we have $(156)_{10} = (10011100)_2$.

Then, we consider the representation of floating point number. Without loss of generality, we represent numbers in $[0, 1)$.

Consider the base 10 number, we have:

$$(0.345)_{10} = 3 \times 10^{-1} + 4 \times 10^{-2} + 5 \times 10^{-3}.$$

For base 2 number, we consider:

$$(0.1101)_2 = 1 \times 2^{-1} + 1 \times 2^{-2} + 1 \times 2^{-4} = 0.5 + 0.25 + 0.625 = (0.8125)_{10}.$$

Example II.2.3. Converting Base 10 Floating to Base 2.

Here, we consider converting base 10 floating points to base 2 floating points, namely $(0.1)_{10}$:

0.1	0.
$0.1 \times 2 = 0.2 < 1$	0.0
$0.2 \times 2 = 0.4 < 1$	0.00
$0.4 \times 2 = 0.8 < 1$	0.000
$0.8 \times 2 = 1.6 \geq 1$	0.0001
$0.6 \times 2 = 1.2 \geq 1$	0.00011
$0.2 \times 2 = 0.4 < 1$	0.000110
\vdots	\vdots

Notice that there is a repeating pattern, so we have:

$$(0.1)_{10} = (0.0001100110011 \dots)_2 = (0.\overline{00011})_2.$$

This is a repeating decimals. If we were to reconvert, we have:

$$\begin{aligned} \left(\frac{1}{2}\right)^4 + \left(\frac{1}{2}\right)^5 + \left(\frac{1}{2}\right)^8 + \left(\frac{1}{2}\right)^9 + \dots &= \left[\left(\frac{1}{2}\right)^4 + \left(\frac{1}{2}\right)^5\right] \left[1 + \left(\frac{1}{2}\right)^4 + \left(\frac{1}{2}\right)^8 + \dots\right] \\ &= \left(\frac{1}{2}\right)^4 \left[1 + \frac{1}{2}\right] \frac{1}{1 - (1/2)^4} = \frac{3}{2} \times \frac{1}{16} \times \frac{1}{15/16} = \frac{1}{10}. \end{aligned}$$

The above example can account for some issues with the floating point inaccuracies in the representation of numbers.

Definition II.2.4. IEEE Floating Point Representation.

For single precision, there are 32 bits, which is 4 bytes.

1 sign bit	#e 8 exponent bits	#f 23 mantissa bits
------------	--------------------	---------------------

For double precision, there are 64 bits, which is 8 bytes.

1 sign bit	#e 11 exponent bits	#f 52 mantissa bits
------------	---------------------	---------------------

Here we have the bits being stored in the continue memory locations. For converting to a base 10 number, we have:

$$N = (-1)^s (1 + f) 2^{e-127}.$$

Here, 127 is the bias for single precision, note that $1 + f \in [1.2)$.

Let $f = (0.m_1 \dots m_{23})_2$, we have it as:

$$m_1 \left(\frac{1}{2}\right)^1 + \dots + m_{23} \left(\frac{1}{2}\right)^{23}.$$

- For $e = 127$, we just have $\pm(1.m_1 \dots m_{23})_2$,

- For $e = 128$, we just have $\pm(1m_1m_2 \cdots m_{23})_2$, which has one less precision, that is 2^{-22} , and it is in $[2, 4)$,
- For $e = 129$, we just have $\pm(1m_1m_2m_3 \cdots m_{23})_2$, which has two less precision, and it is in $[4, 8)$.

For the smaller numbers, we have:

- For $e = 126$, we have $(0.1m_1 \cdots m_{23})_2$, so the precision is one more, that is 2^{-24} , the size of interval is $[1/2, 1)$.
- For $e = 127$, we have $(0.1m_1 \cdots m_{23})_2$, so the precision is two more, , the size of interval is $[1/4, 1/2)$.

Example II.2.5. Converting Base 10 to FP.

Consider converting from base 10 to FP, we have:

$$(15)_{10} = (1111)_2, (0.1)_{10} = (0.0\overline{1111})_2,$$

thus we have:

$$(15.1)_{10} = 1111.000\overline{11} = 1.111000\overline{11} \times 2^3$$

Hence the exponent is $127 + 3 = 130$, and the mantissa being $111000\overline{11}$, up to the correct number of digits:

1 sign bit	#e 8 exponent bits	#f 23 mantissa bits
0	10000010	11100011001100110011001

In this case, we truncated all the digits afterwards, causing imprecisions.

In particular, the catastrophic cancellation since we are not considering the digits afterwards.

Even with double precisions points, the we still have the base as:

$$(-1)^s(1 + f) \underbrace{2^{e-1023}}_{\text{bias}},$$

Definition II.2.6. Machine Representable Number.

A Machine representable number (MRN) is a number that can be represented exactly $f_p(x) = x$.

Let $f = (0.m_1m_2 \cdots m_{52})_2$, and we have the following case:

- For $e = 1023$, we have the values being in $[1, 2)$, with the coefficient of $(1/2)^{52} \simeq 2.220 \times 10^{-16}$.
- For $e = 1024$, we have the values being in $[2, 4)$, with the coefficient of $(1/2)^{51} \simeq 4.441 \times 10^{-16}$.
- For $e = 1025$, we have the values being in $[4, 8)$, with the coefficient of $(1/2)^{50} \simeq 8.882 \times 10^{-16}$.
- For $e = 1077$, we have the values being in $[2^{52}, 2^{53})$, with the coefficient of $(1/2)^0 = 1$.

For any generic e , we have the following:

- Interval: $[2^{e-1023}, 2^{e-1022})$,
- Width of interval: $2^{e-1022} - 2^{e-1023} = 2^{e-1023}$,
- Step size: $2^{-52+e-1023} = 2^{e-1075}$, and

- Number of MRNs: 2^{-52} .

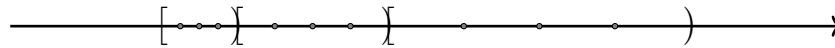
Example II.2.7. Largest and Smallest (Positive) MRN.

The largest MRN is $\underbrace{111 \dots 1}_{52 \text{ ones}} \underbrace{000 \dots 0}_{972 \text{ zeros}} \simeq 10^{304}$. Here, the step size is $10^{1024-52}$.

The smallest MRN is $\underbrace{0.000 \dots 0}_{1022 \text{ zeros}} \underbrace{1}_{52 \text{ ones}} \underbrace{000 \dots 0}_{972 \text{ zeros}} = 2^{-1023} \simeq 2.470 \times 10^{-324}$. (This is evaluated to zero on computers).

The next smallest (real smallest) is $2^{-1023} + 2^{-(1023+52)}$.

For each fixed e , it represents an interval, which got larger when e grows larger, with the numbers inside being uniformly distributed along each interval:


Remark II.2.8. Zero, Infinity, and NaN in Machine Representation.

For floating numbers, 0, $\pm\infty$, and NaN (not a number) cannot be represented conventionally, but they are distributed with a special slot.

- For 0, we have:

1 sign bit	#e 11 exponent bits	#f 52 mantissa bits
0 or 1	000...0	000...0

- For ∞ , we have:

1 sign bit	#e 11 exponent bits	#f 52 mantissa bits
0 or 1	111...1	000...0

For the $-\infty$, it makes the sign bit as negative.

- For NaN, we have:

1 sign bit	#e 11 exponent bits	#f 52 mantissa bits
0 or 1	111...1	111...1

Definition II.2.9. Machine Epsilon.

The machine epsilon, denoted $\epsilon_{\text{machine}}$, is the distance between 1 and the next larger MRN.

Example II.2.10. Machine Epsilon for 1.

For the double precision floating point, we represent the number 1 as:

1 sign bit	#e 11 exponent bits	#f 52 mantissa bits
0	01111111111	000...0

where the exponent is 1023 in base 10.

Therefore, the next MRN is:

1 sign bit	#e 11 exponent bits	#f 52 mantissa bits
0	01111111111	000...01

which is exactly $1 + 2^{-52}$, or the step size for $e = 1023$ is always 2^{-52} . Therefore:

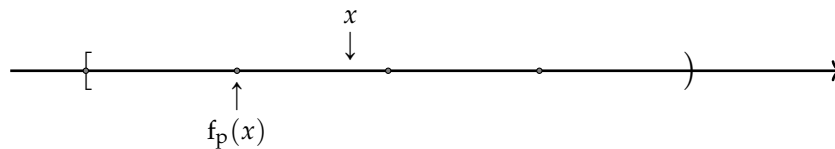
$$\epsilon_{\text{machine}} = 2^{-52}.$$

Remark II.2.11. Absolute and Relative Error.

Let x be a positive real number, we have:

- Let $f_p(x)$ be the floating point representation of x .
- We assume that we are in double precision.

From x , you first determine e , and then we find the number on the interval (for simplicity, we just represent 4 points on the number line):



Suppose we have just truncation, we have:

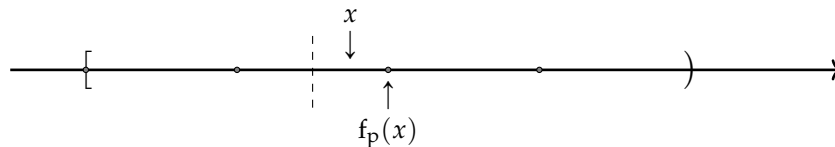
- The absolute error is $|x - f_p(x)| \leq 2^{-52+e-1024}$, which is the step size.
- The relative error is $\frac{|x - f_p(x)|}{|f_p(x)|} \leq \frac{2^{-52+e-1024}}{|(-1)^s(1+f)2^{e-1024}|} = \frac{2^{-52}}{1+f} \leq \epsilon_{\text{machine}}.$

The key conclusion is that the relative error is bounded above by the machine epsilon.

Alternatively, if we consider the numerator of the relative error as x , we have:

$$\frac{|x - f_p(x)|}{|x|} = \frac{|x - f_p(x)|}{|f_p(x)|} \cdot \frac{|f_p(x)|}{|x|} \leq \epsilon_{\text{machine}}.$$

Then, we consider the rounding off, so we have:



Now, the absolute error is halved, and relative error is bounded by $\epsilon_{\text{machine}}/2$.

Here, we shall be concerned on how we can trust the computer.

For double precisions floating point, we have:

$$\epsilon_{\text{machine}} = 2^{-52} \simeq 2.22045 \times 10^{-16}.$$

So we can keep about 15 to 16 digits accurate

Example II.2.12. MATLAB Illustration on Numbers.

Below is a code segment of the illustration on MATLAB

```
>> 1.23456789012345678901234567890
```

```
ans = 1.23456789012346
```

```
>> ans + 100000000
```

```
ans = 1.0000000012345679e+008
```

```
>> ans - 100000000
```

```
ans = 1.234567890553165
```

Note that for the second operation, the operation truncates the digits after the first 9, and when the same number is subtracted, it results in some junk digits.

More specifically, look at the following code snippet:

```
>> 1.2345678901234567890 + 10000000 - 10000000
```

```
ans = 1.234567890553165
```

```
>> 1.2345678901234567890 + (10000000 - 10000000)
```

```
ans = 1.234567890123456
```

Here, we can observe that associativity does not hold over computer level computations.

This is caused by the similarity of significant digits, so only the first 15-th digits are exact, and from 16-th digits afterwards, they may be affected by round-off.

This is an example of the *Catastrophic Cancellation*, which occurs when computing $x - y$, where $x > y$ but $x \simeq y$.

There, $x - y$ can result in fewer significant digits than x and/or y .

Example II.2.13. Trick to Eliminate Subtraction.

When we compute:

$$y = \frac{\sqrt{x^2 + 4} - x}{2} \text{ for } x \text{ very large, and } x > 0.$$

For $x = 10^{10.5}$, MATLAB returns $y = 0$. However, we can have:

$$y = \frac{\sqrt{x^2 + 4} - x}{2} \cdot \frac{\sqrt{x^2 + 4} + x}{\sqrt{x^2 + 4} + x} = \frac{x^2 + 4 - x^2}{2(\sqrt{x^2 + 4} + x)} = \frac{2}{\sqrt{x^2 + 4} + x},$$

and without subtractions inside, we have the output approximately as 3.162278×10^{-11} .

Remark II.2.14. Real Number and Floating Point.

Let $x \in \mathbb{R}$, and $f_p(x)$ be the floating point representation of x , which is a MRN, and it is defined as:

$$\epsilon_{\text{machine}} = (-1)^s (1 + f) 2^{e-1023}.$$

The machine epsilon is determined as the absolute value of the difference between 1 and the next larger MRP, that is 2^{-52} for double precision floating point number.

The relative error in floating point approximating is:

$$\frac{|x - f_p(x)|}{|x|}.$$

Whereas for truncation, we redefine:

$$\epsilon_{\text{machine}} = \frac{1}{2(1 - 2^{-52})} 2^{-52}.$$

Thus, in both cases, we have:

$$\frac{|x - f_p(x)|}{|x|} \leq \epsilon_{\text{machine}}.$$

Below, we introduce a “axiom” of floating point representations, but it turns out to be direct from the previous remark.

Proposition II.2.15. Property of Floating Point Representation.

For all $x \in \mathbb{R}$, there exists ε (positive or negative) with $|\varepsilon| \leq \epsilon_{\text{machine}}$ such that:

$$f_p(x) = x(1 + \varepsilon).$$

I.e., the relative distance between x and $f_p(x)$ is always smaller than $\epsilon_{\text{machine}}$.

Proof. Without loss of generality, we let $x > 0$, since $x < 0$ is a similar case. We have:

$$\begin{aligned} -\epsilon_{\text{machine}}x &\leq f_p(x) - x \leq \epsilon_{\text{machine}}x, \\ x - \epsilon_{\text{machine}}x &\leq f_p(x) \leq x + \epsilon_{\text{machine}}x, \\ x(1 - \epsilon_{\text{machine}}) &\leq f_p(x) \leq x(1 + \epsilon_{\text{machine}}), \\ 1 - \epsilon_{\text{machine}} &\leq \frac{f_p(x)}{x} \leq 1 + \epsilon_{\text{machine}}. \end{aligned}$$

Therefore, there exists ε with $|\varepsilon| \leq \epsilon_{\text{machine}}$ such that:

$$\frac{f_p(x)}{x} = 1 + \varepsilon,$$

so we have $f_p(x) = x(1 + \varepsilon)$, as desired. □

Then, we think about the floating point arithmetic.

Remark II.2.16. Notation for Arithmetics.

Here, we have \mathbb{R} denote the real numbers, and \mathcal{F} denote the machine representable numbers in floating point with double precision. We think of $f_p : \mathbb{R} \rightarrow \mathcal{F}$.

The four main arithmetic operations on \mathbb{R} are $+$, $-$, \times , and \div .

The four main arithmetic operations on \mathcal{F} are \oplus , \ominus , \otimes , and \odiv .

Here, let $*$ be the generic operator, since we have $\mathcal{F} \subset \mathbb{R}$, we defined the maps in \mathcal{F} naturally pre-composing the inclusion $\iota : \mathcal{F} \hookrightarrow \mathbb{R}$ and post-compose with the f_p . Therefore, for any $x', y' \in \mathcal{F} \subset \mathbb{R}$, we defined the operation as:

$$x' \circledast y' = f_p(\iota(x') * \iota(y')) = f_p(x' * y').$$

Here, by the property of floating point representation, we have:

$$x' \circledast y' = f_p(x' * y') = (x' * y') \cdot (1 + \varepsilon),$$

where $|\varepsilon| \leq \epsilon_{\text{machine}}$.

Definition II.2.17. Big O Notation.

Given two functions of real-valued inputs, $a(t)$ and $b(t) > 0$, we have $a(t) = \mathcal{O}(b(t))$ as $t \rightarrow 0$ when there exists $C > 0$ such that $|a(t)| \leq C \cdot b(t)$ in a neighborhood of $t = 0$, i.e., there exists some $\delta > 0$ such that the statement holds for all $|t| < \delta$. J

Here, we can have an example with some function.

Example II.2.18. Sine Function in Big O.

$\sin t = \mathcal{O}(|t|)$ as $t \rightarrow 0$, since we have $|\sin t| \leq |t|$ in $B_\delta(0)$ for some $\delta > 0$. J

Remark II.2.19. Floating Point with Machine Epsilon.

If $x' = f_p(x)$, then:

$$\frac{|x - x'|}{|x|} < 1 \cdot \epsilon_{\text{machine}},$$

thus, we have:

$$\frac{|x - x'|}{|x|} = \mathcal{O}(\epsilon_{\text{machine}}).$$

We say that the relative error is of the order of $\epsilon_{\text{machine}}$. J

III Computational Methods

III.1 Least Square Approximations

Consider $A \in \mathbb{C}^{m \times n}$ with $A = \begin{bmatrix} \mathbf{A}_1 & \mathbf{A}_2 & \cdots & \mathbf{A}_n \end{bmatrix}$.

Remark III.1.1. Equivalence Conditions to Trivial Kernel.

For the above A , we have $\ker A = \{\mathbf{0}\} \iff \mathbf{A}_1, \dots, \mathbf{A}_n$ are linearly independent $\iff \dim(\text{im } A) = n$.
In such case, we have $m \geq n$, that is, it must be square or having more rows.

Remark III.1.2. Existence and Uniqueness of Solutions to Linear Equation.

From linear algebra, with $A\mathbf{x} = \mathbf{b}$ for matrix A and vectors \mathbf{x}, \mathbf{b} , we have:

- (i) If $\mathbf{b} \in \text{im } A$, then there exists (at least) a solution to the linear equation.
- (ii) If $\ker A = \{\mathbf{0}\}$, the uniqueness is guaranteed.

If $\mathbf{b} \notin \text{im } A$, then for all \mathbf{x} , we have $A\mathbf{x} \neq \mathbf{b}$.

Our goal is to choose $\mathbf{x} \in \mathbb{C}^n$ so that $A\mathbf{x} \simeq \mathbf{b}$, i.e., $A\mathbf{x}$ is as close to \mathbf{b} as possible.

This leads to the least square problem: Given a matrix $A \in \mathbb{C}^{m \times n}$ and a vector $\mathbf{b} \in \mathbb{C}^m$, we find \mathbf{x} that minimizes the residue:

$$\epsilon = \|\mathbf{b} - A\mathbf{x}\|_2$$

We are having the following assumptions:

- (i) $m \geq n$, and
- (ii) $\dim(\text{im } A) = n \iff \ker A = \{\mathbf{0}\}$.

Recall that we have $\ker A = \ker(A^*A)$, so $\ker A = \{\mathbf{0}\}$ implies that A^*A is invertible.

For the first case, we consider $\mathbf{b} \in \text{im } A$, so we have $A\mathbf{x} = \mathbf{b}$ being consistent but overdetermined, i.e., there are more equations than variables. The following is a illustration when $n = 2$:

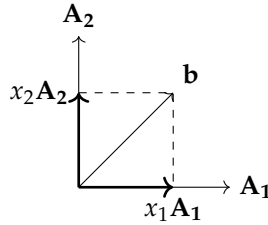


Figure III.1. Overdetermined equation in $n = 2$ plane for image.

Here, we then have $A\mathbf{x} = \mathbf{b}$ implying that $A^*A\mathbf{x} = A^*\mathbf{b}$, in which A^*A is invertible, hence $\mathbf{x} = (A^*A)^{-1}A^*\mathbf{b}$, in which we have:

$$A^+ = (A^*A)^{-1}A^* \in \mathbb{C}^{n \times m},$$

which is the *pseudo-inverse* of A .

Proposition III.1.3. Properties of Pseudo-inverse.

For a *pseudo-inverse* of A , denoted A^+ , the following properties hold:

- (i) $A^+A = (A^*A)^{-1}A^*A = \text{Id}_n$,
- (ii) $AA^+ = A(A^*A)^{-1}A^* = P$.

Then, we shall consider the other case, *i.e.*, if $\mathbf{b} \notin \text{im } A$, then the minimizer of $\inf_{\mathbf{x} \in \mathbb{C}^n} \|\mathbf{b} - A\mathbf{x}\|_2$ is still $\mathbf{x} = A^+\mathbf{b}$, where \mathbf{x} is the vector of coordinates of \mathbf{y} with respect to $\mathbf{A}_1, \dots, \mathbf{A}_n$.

- For a generic $\mathbf{x} \in \mathbb{C}^n$, we have $\mathbf{z} = A\mathbf{x} \in \text{im } A$, and
- For any $\mathbf{z} \in \text{im } A = \text{span}\{\mathbf{A}_1, \dots, \mathbf{A}_n\}$, there exists $\mathbf{x} \in \mathbb{C}^n$ such that $\mathbf{z} = A\mathbf{x}$.

Remember that we want to minimize $\inf_{\mathbf{x} \in \mathbb{C}^n} \|\mathbf{b} - A\mathbf{x}\|_2$, which is equivalent to minimizing $\inf_{\mathbf{z} \in \text{im } A} \|\mathbf{b} - \mathbf{z}\|_2$, which can be considered as follows:

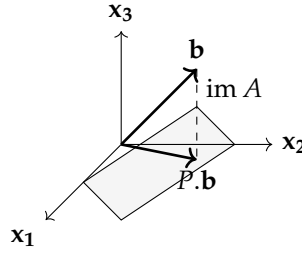


Figure III.2. The projection onto the image span.

Proposition III.1.4. Minimizer of the Approximation Problem.

Let $A \in \mathbb{C}^{m \times n}$ with $m \geq n$, we have $\dim(\text{im } A) = n$. We fix $\mathbf{b} \in \mathbb{C}^m$ arbitrarily. The minimizer of the problem is:

$$\inf_{\mathbf{z} \in \text{im } A} \|\mathbf{b} - \mathbf{z}\|_2$$

is $\mathbf{z} = \mathbf{y}$, where $\mathbf{y} = P\mathbf{b} = AA^+\mathbf{b}$.

Equivalently, the minimizer of $\inf_{\mathbf{x} \in \mathbb{C}^n} \|\mathbf{b} - A\mathbf{x}\|_2$ is $\mathbf{x} = A^+\mathbf{b}$.

Such situation can be applied onto linear regression in statistics.

Example III.1.5. Linear Regression Example.

Let data points in \mathbb{R}^2 be:

$$\{(x_1, y_1), \dots, (x_m, y_m)\},$$

so the linear model can be considered as:

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 = \begin{bmatrix} 1 & x & x^2 \end{bmatrix} \cdot \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix}.$$

When $\beta_0, \beta_1, \beta_2$ are fixed, for each x_i , the i -th prediction is:

$$\hat{y}_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2,$$

and the i -th residue is:

$$\epsilon_i = y_i - \hat{y}_i = y_i - (\beta_0 + \beta_1 x_i + \beta_2 x_i^2),$$

so the sum of the square errors is:

$$\text{SSE}(\beta_0, \beta_1, \beta_2) = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n \left[y_i - (\beta_0 + \beta_1 x_i + \beta_2 x_i^2) \right]^2.$$

Here, we use the trick by:

$$A = \begin{bmatrix} 1 & x_1 & x_1^2 \\ \vdots & \vdots & \vdots \\ 1 & x_m & x_m^2 \end{bmatrix},$$

$$\boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}.$$

Here, we have:

$$\mathbf{b} - A \cdot \boldsymbol{\beta} = \begin{bmatrix} y_1 - (\beta_0 + \beta_1 x_1 + \beta_2 x_1^2) \\ \vdots \\ y_m - (\beta_0 + \beta_1 x_m + \beta_2 x_m^2) \end{bmatrix}.$$

Therefore, we have:

$$\|\mathbf{b} - A \cdot \boldsymbol{\beta}\|^2 = \sum_{i=1}^n \left[y_i - (\beta_0 + \beta_1 x_i + \beta_2 x_i^2) \right]^2 = \text{SSE}(\beta_0, \beta_1, \beta_2),$$

so we are minimizing $\hat{\boldsymbol{\beta}} = A^+ \cdot \mathbf{b}$.

Theorem III.1.6. Projections Minimize the Least Square Problem.

Let $A \in \mathbb{C}^{m \times n}$, with $m \geq n$ have $\dim(\text{im } A) = n$ (full rank). Let $\mathbf{b} \in \mathbb{R}^m$ be arbitrary, the minimizer of the problem is:

$$\min_{\mathbf{z} \in \text{im } A} \|\mathbf{b} - \mathbf{z}\|_2,$$

is $\mathbf{z} = \mathbf{y}$, where:

$$\mathbf{y} = P\mathbf{b} = AA^+ \mathbf{b} = A(A^*A)^{-1}A^* \mathbf{b}.$$

Equivalently, the minimizer of the problem $\min_{\mathbf{x} \in \mathbb{C}^n} \|\mathbf{b} - A \cdot \mathbf{x}\|_2$ is $\mathbf{x} = A^+ \cdot \mathbf{b}$.

The full rank condition is helpful when trying to solve $A \cdot \mathbf{x} = \mathbf{b}$, when $\mathbf{b} \notin \text{im } A$, so for all \mathbf{x} , $A \cdot \mathbf{x} \neq \mathbf{b}$.

Remark III.1.7. Computational Costs for Matrix Operations.

The computational cost to compute A^+ is:

- $C(A^*A) \sim 2mn^2$,
- $C((A^*A)^{-1}) \sim \frac{8}{3}n^2$, and
- $C((A^*A)^{-1}A^*) \sim 2mn^2$.

Hence, the total computation cost for the projection for least square problem is $4mn^2 + \frac{8}{3}n^3$.

Note that there are other ways to compute $\mathbf{x} = A^+ \cdot \mathbf{b}$ by:

- (i) Cholesky factorization,
- (ii) QR factorization, and
- (iii) SVD.

In particular, we have the QR factorization as:

$$\begin{array}{ccccc} A & = & Q & \circ & R, \\ \mathfrak{m} & & \mathfrak{m} & & \mathfrak{m} \\ \mathbb{C}^{m \times n} & & \mathbb{C}^{m \times m} & & \mathbb{C}^{m \times n} \end{array}$$

where R has the top n rows being upper triangular.

Moreover, we may use the reduced QR decomposition that $A = \tilde{Q}\tilde{R}$ where we reduce to the n columns for Q and R . However, we have:

$$\tilde{Q}^* \tilde{Q} = \begin{bmatrix} \mathbf{q}_1^* \\ \vdots \\ \mathbf{q}_n^* \end{bmatrix} \begin{bmatrix} \mathbf{q}_1^* & \cdots & \mathbf{q}_n^* \end{bmatrix} = \text{Id}.$$

Also, we have $\text{span}\{\mathbf{q}_1, \dots, \mathbf{q}_n\} = \text{span}\{\mathbf{A}_1, \dots, \mathbf{A}_n\}$, so $P = \tilde{Q}(\tilde{Q}^* \tilde{Q})^{-1} \tilde{Q}^* = \tilde{Q} \tilde{Q}^*$.

We let \mathbf{y} be the orthogonal projection of \mathbf{b} onto $\text{im } A$, that is:

$$\mathbf{y} = P\mathbf{b} = \tilde{Q} \tilde{Q}^* \mathbf{b},$$

and let \mathbf{x} be the coordinates of \mathbf{y} with respect to $\{\mathbf{A}_1, \dots, \mathbf{A}_n\}$, so:

$$\mathbf{y} = \sum_{i=1}^n x_i \mathbf{A}_i = A \cdot \mathbf{x}.$$

Therefore, we have $A \cdot \mathbf{x} = \tilde{Q} \tilde{Q}^* \mathbf{b}$, that is $\tilde{Q} \tilde{R} \mathbf{x} = \tilde{Q} \tilde{Q}^* \mathbf{b}$, by post-compose \tilde{Q}^* , we can get $\tilde{R} \cdot \mathbf{x} = \tilde{Q}^* \mathbf{b}$, which is easy to compute.

Remark III.1.8. Algorithm and Cost of the Minimization Process.

The inputs are A and \mathbf{b} , and the output is \mathbf{x} , where $\mathbf{x} = A^+ \cdot \mathbf{b}$, we do the following:

- (i) compute the reduced QR decomposition of A (with the householder triangularization, its complexity is $2mn^2 - \frac{2}{3}n^3$),
- (ii) compute $\mathbf{c} = \tilde{Q}^* \mathbf{b}$ (there are n dot products in \mathbb{C}^m , so a total of about $2mn$), and then
- (iii) solve $\tilde{R} \cdot \mathbf{x} = \tilde{Q}^* \cdot \mathbf{b}$ for \mathbf{x} (for this part, note that the number of FLOP each block substitution is $1 + 3 + 5 + \dots$, that is n^2).

Therefore, the computational cost is $\sim 2mn^2 - \frac{2}{3}n^3$.

III.2 Conditioning and Condition Number

Definition III.2.1. Ill Conditioned.

A problem is “ill conditioned” when a small variation of data causes large variation of solution.

Example III.2.2. Ill Conditioned Problem.

Consider $A\mathbf{x} = \mathbf{b}$, it is either:

- (i) $A \in \mathbb{C}^{m \times m}$ is invertible, then $\mathbf{x} = A^{-1} \cdot \mathbf{b}$, or
- (ii) $A \in \mathbb{C}^{m \times n}$, $\dim(\text{im } A) = n$, and $\mathbf{b} \in \text{im } A$, so $\mathbf{x} = A^+ \cdot \mathbf{b}$.

Suppose we compute a solution $\tilde{\mathbf{x}}$ (not quite correct), and the error is:

$$\mathbf{e} = \delta\mathbf{x} = \mathbf{x} - \tilde{\mathbf{x}},$$

which is the difference between the real and computed solution to $A\mathbf{x} = \mathbf{b}$. In reality, this cannot be computed since we do not have access to the real solution (that is why we are computing it.)

Hence, the goal is to find an upper bound for relative size of error, *i.e.*:

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|}.$$

What we can compute is the residue:

$$\mathbf{r} = \delta\mathbf{b} = \mathbf{b} - A\tilde{\mathbf{x}} = A\mathbf{x} - A\tilde{\mathbf{x}} = A(\mathbf{x} - \tilde{\mathbf{x}}) = A\delta\mathbf{x}.$$

Therefore, we have $\delta\mathbf{x} = A^{-1}\delta\mathbf{b}$, hence its norm is:

$$\|\delta\mathbf{x}\| = \|A^{-1}\delta\mathbf{b}\| \leq \|A^{-1}\| \cdot \|\delta\mathbf{b}\|.$$

Then, since we have $\mathbf{b} = A\mathbf{x}$, we have:

$$\|\mathbf{b}\| = \|A\mathbf{x}\| \leq \|A\| \cdot \|\mathbf{x}\|,$$

Therefore, we may obtain the upper bound of the residual as:

$$\underbrace{\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|}}_{\text{relative size of error}} \leq \underbrace{\|A\| \cdot \|A^{-1}\|}_{\text{conditional number of the matrix } A} \cdot \underbrace{\frac{\|\delta\mathbf{b}\|}{\|\mathbf{b}\|}}_{\text{relative size of residual}}.$$

Definition III.2.3. Condition Number of Matrix.

Let A be an invertible matrix, we have the condition number of A as:

$$\kappa(A) = \|A\| \cdot \|A^{-1}\|.$$

The large condition number means ill-conditioned, and small condition number means well conditioned. If A^{-1} is hard to compute, then one eigenvalue is $\simeq 0$, which implies large $\kappa(A)$.

When $\|\bullet\| = \|\bullet\|_2$, we have:

$$\kappa(A) = \frac{\sigma_1}{\sigma_m}.$$

Here are some motivation of the condition number of a square, invertible matrix A :

$$\kappa(A) = \underbrace{\|A\| \cdot \|A^{-1}\|}_{\text{norm induced by vector norm}}.$$

For example, when having the 2×2 case, we want to solve that $A \cdot \mathbf{x} = \mathbf{b}$, and the solution is $\mathbf{x} = A^{-1} \cdot \mathbf{b}$, where as the inverse is:

$$A = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \Rightarrow A^{-1} = \frac{1}{ad - bc} \underbrace{\begin{bmatrix} d & -b \\ -c & a \end{bmatrix}}_{\text{adj}(A)}.$$

Note that when A is close to singularity, that is, $\det A$ is close to zero, then we have catastrophic cancellation, which happens when subtracting 2 numbers very close to each other. *This is a large round-off error.*

In the $m \times m$ case, we do the inverse as:

$$\left[A \mid \text{Id}_n \right] \xrightarrow{\text{RREF}} \left[\text{Id}_n \mid A^{-1} \right]$$

The gist is that when A is close to singularity, A^{-1} cannot be computed accurately, so $\mathbf{x} = A^{-1} \mathbf{b}$ cannot be computed accurately.

Let $\tilde{\mathbf{x}}$ be the estimate of the solution \mathbf{x} , we can compute the residual as:

$$\delta \mathbf{b} = \mathbf{b} - A \cdot \tilde{\mathbf{x}} = A \cdot \mathbf{x} - A \cdot \tilde{\mathbf{x}} = A \cdot (\underbrace{\mathbf{x} - \tilde{\mathbf{x}}}_{\delta \mathbf{x} \text{ error}}).$$

Remark III.2.4. Computation of $\delta \mathbf{x}$.

Here, we have $\delta \mathbf{x} = A^{-1} \delta \mathbf{b}$, since we cannot compute $\delta \mathbf{x}$ accurately.

The general rule for numerical linear algebra is to avoid computing $\det A$ and A^{-1} , because they are:

- (i) computationally intensive, and
- (ii) computationally inaccurate when A is close to singularity.

Example III.2.5. Pseudo Inverse Case.

If $A \in \mathbb{C}^{m \times n}$ such that $m \geq n$, with full rank and $\mathbf{b} \in \mathbb{C}^m$, then the minimizer of $\|\mathbf{b} - A \cdot \mathbf{x}\|_2$ is:

$$\mathbf{x} = A^+ \cdot \mathbf{b} = (A^* A)^{-1} A^* \cdot \mathbf{b}.$$

We avoided inverting $A^* A$ by computing \mathbf{x} via QR decomposition by solving $\tilde{R} \cdot \mathbf{x} = \tilde{Q}^* \cdot \mathbf{b}$.

Later we will see techniques to compute:

- the eigenvalues of square matrices, and
- the SVD of any matrix,

that avoid determinant and inverse.

We have proven that:

$$\underbrace{\frac{\|\delta \mathbf{x}\|}{\|\mathbf{x}\|}}_{\text{relative size of error (cannot compute)}} \leq \underbrace{\|A\| \cdot \|A^{-1}\|}_{\kappa(A) \text{ conditional number (typically can computed without } A^{-1})}} \cdot \underbrace{\frac{\|\delta \mathbf{b}\|}{\|\mathbf{b}\|}}_{\text{relative size of residual (can compute)}}.$$

Remark III.2.6. Conditional Number based on Choice of Vector Norm.

$\kappa(A)$ depends on the choice of the vector norm. For example, choose $\|\bullet\|_2$ as a 2-norm, we have:

$$\kappa_2(A) = \|A\|_2 \cdot \|A^{-1}\|_2.$$

Theorem III.2.7. Conditional Number as Fraction of Singular Value.

For any invertible $A \in \mathbb{C}^{m \times m}$:

$$\kappa(A) = \frac{\sigma_1}{\sigma_m},$$

where σ_1 and σ_m are, respectively, the largest and smallest singular value of A .

Proof. For any $A \in \mathbb{C}^{m \times m}$, we have $\|A\|_2 = \sigma_1$, that is the largest eigenvalue.

When writing SVD as:

$$A = U \Sigma V^* = U \begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_m \end{bmatrix} V^*.$$

Then:

$$A^{-1} = (V^*)^{-1} \Sigma^{-1} U^{-1} = \bar{V} \begin{bmatrix} 1/\sigma_1 & 0 & \cdots & 0 \\ 0 & 1/\sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1/\sigma_m \end{bmatrix} (\bar{U})^* = U_1 \begin{bmatrix} 1/\sigma_1 & 0 & \cdots & 0 \\ 0 & 1/\sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1/\sigma_m \end{bmatrix} V_1^*.$$

Note that since $V^{-1} = V^*$, then $V = \bar{V}$, so their complex conjugates are the same, now, we may invert the order of the rows to obtain that:

$$A^{-1} = U_2 \begin{bmatrix} 1/\sigma_m & 0 & \cdots & 0 \\ 0 & 1/\sigma_{m-1} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1/\sigma_1 \end{bmatrix} V_2^*,$$

hence we have the largest singular value is $1/\sigma_m$. The conditional number follows as $\sigma_1 \cdot \frac{1}{\sigma_m} = \frac{\sigma_1}{\sigma_m}$. \square

In particular, the fraction is called the eccentricity of ellipsoid with semi-axes $\sigma_1, \dots, \sigma_m$.

Example III.2.8. Conditional Number for Almost Singular Matrices.

In MATLAB, the command for conditional number is `cond(A)`, which is $\kappa_2(A)$. Here, we let:

$$A = \begin{bmatrix} 1 & 1 \\ 1 & 1 + \alpha \end{bmatrix}.$$

- When $\alpha = 1$, we have $\kappa_2(A) \approx 6.85410$,
- When $\alpha = 10^{-5}$, we have $\kappa_2(A) \approx 4.00002 \times 10^5$,
- When $\alpha = 10^{-12}$, we have $\kappa_2(A) \approx 3.99949 \times 10^{12}$.

Theorem III.2.9. Conditional Number of Matrix Operations.

For any matrix norm induced by a vector norm, we have:

- (i) For any nonzero constant $c \in \mathbb{C}$, we have that:

$$\kappa(cA) = \kappa(A).$$

- (ii) $\kappa(\text{Id}) = 1$, and

- (iii) for any invertible A , $\kappa(A) \geq 1$.

Proof. Recall that $\|\text{Id}\| = 1$ and $\|A^{-1}\| \geq 1/\|A\|$ for any norm.

- (i) Note that $(cA)^{-1} = 1/c \cdot A^{-1}$, so $\kappa(cA) = \|cA\| \cdot \|1/c \cdot A^{-1}\| = |c|/|c| \cdot \|A\| \cdot \|A^{-1}\| = \kappa(A)$.

- (ii) $\kappa(\text{Id}) = \|\text{Id}\| \cdot \|\text{Id}^{-1}\| = 1 \cdot 1 = 1$.

- (iii) $\kappa(A) = \|A\| \cdot \|A^{-1}\| \geq \|A\|/\|A\| = 1$, so $\kappa(A) \geq 1$. □

III.3 Stability

For the mathematical problem, we let it be defined as $f : B \rightarrow X$, where B is the set of possible data and X is the set of solutions.

Example III.3.1. Basic Linear Algebra Problem.

We are trying to solve that $A\mathbf{x} = \mathbf{b}$, where $A \in \mathbb{C}^{m \times m}$ is invertible.

The solution is:

$$\mathbf{x} = f(\mathbf{b}) = A^{-1}\mathbf{b}.$$

Typically (especially when the condition number $\kappa(A)$ is large), we can compute only an approximated version of \mathbf{x} , denoted $\tilde{\mathbf{x}}$, through an algorithm (such as QR decomposition).

Definition III.3.2. Algorithm.

Here, we define an algorithm as:

$$\tilde{f} : B \rightarrow X,$$

in which the computed solution $\tilde{\mathbf{x}} = \tilde{f}(\mathbf{b})$ is by the actual data.

A “stable” algorithm is one which computes solutions $\tilde{\mathbf{x}}$ which are approximately equal to the exact solution for slightly perturbed data:

$$\tilde{f}(\mathbf{b}) \simeq f(\mathbf{b} + \delta\mathbf{b}).$$

Definition III.3.3. Backward Stable.

Given a problem f , an algorithm \tilde{f} is called backward stable if for all set of data $\mathbf{b} \in B$, there exists a data perturbation $\delta\mathbf{b}$, where:

$$\frac{\|\delta\mathbf{b}\|}{\|\mathbf{b}\|} = \mathcal{O}(\epsilon_{\text{machine}}) \text{ such that } \tilde{f}(\mathbf{b}) = f(\mathbf{b} + \delta\mathbf{b}).$$

This stability is called backward since instead of looking at the *forward error*, we have:

$$\delta\mathbf{x} = \tilde{f}(\mathbf{b}) - f(\mathbf{b}),$$

we look backwards to see what input could have produced the computed result $\tilde{f}(\mathbf{b})$ exactly.

This definition holds, the data perturbation required to explain the computed solution is relatively small (relative to the problem's data \mathbf{b}), *i.e.*, the algorithm is numerically robust (stable), to relatively small perturbations.

Example III.3.4. Subtraction is Backwards Stable.

Suppose the mathematical problem is:

$$f : \mathbb{C}^2 \rightarrow \mathbb{C}, \mathbf{b} = \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} \mapsto b_1 - b_2.$$

The algorithm is as follows:

input: $\mathbf{b} = [b_1, b_2]$

output: $\sim x = \sim f(b_1, b_2) = \text{fp}(b_1) - \text{fp}(b_2) = \text{fp}(b_1' - b_2')$

Then, there exists ϵ_1 with $|\epsilon_1| \leq \epsilon_{\text{machine}}$ such that $b'_1 = \text{fp}(b_1) = b_1(1 + \epsilon_1)$. There exists ϵ_2 with $|\epsilon_2| \leq \epsilon_{\text{machine}}$ such that $b'_2 = \text{fp}(b_2) = b_2(1 + \epsilon_2)$. Moreover, there exists ϵ_3 with $|\epsilon_3| \leq \epsilon_{\text{machine}}$ such that $b'_3 = \text{fp}(b_3) = b_3(1 + \epsilon_3)$.

Hence, we have:

$$\begin{aligned} (b'_1 - b'_2)(1 + \epsilon_3) &= [b_1(1 + \epsilon_1) - b_2(1 + \epsilon_2)](1 + \epsilon_3) = b_1(1 + \epsilon_1)(1 + \epsilon_3) - b_2(1 + \epsilon_2)(1 + \epsilon_3) \\ &= b_1 + \underbrace{b_1(\epsilon_1 + \epsilon_3 + \epsilon_1\epsilon_3)}_{\delta b_1} - [b_2 + \underbrace{b_2(\epsilon_2 + \epsilon_3 + \epsilon_2\epsilon_3)}_{\delta b_2}] = (b_1 + \delta b_1) - (b_2 + \delta b_2) \\ &= f(b_1 + \delta b_1 + b_2 + \delta b_2) = f(\mathbf{b} + \delta\mathbf{b}). \end{aligned}$$

Just to note since for the \mathcal{O} , we are letting it $\rightarrow 0$, not infinity, so we have:

$$|\epsilon_1 + \epsilon_3 + \epsilon_1\epsilon_3| \leq |\epsilon_1| + |\epsilon_3| + |\epsilon_1| \cdot |\epsilon_3| \leq 2\epsilon_{\text{machine}} + \epsilon_{\text{machine}}^2 = \mathcal{O}(\epsilon_{\text{machine}}).$$

Hence, we have the algorithm backwards stable.

Here, we consider such operation in norm notations as well:

$$\delta\mathbf{b} = \begin{bmatrix} \delta b_1 \\ \delta b_2 \end{bmatrix} = b_1\epsilon_4 + b_2\epsilon_5.$$

Hence the square of the norm is:

$$\|\delta\mathbf{b}\|^2 = |b_1|^2|\epsilon_4|^2 + |b_2|^2|\epsilon_5|^2 = |b_1|^2\mathcal{O}(\epsilon_{\text{machine}}^2) + |b_2|^2\mathcal{O}(\epsilon_{\text{machine}}^2) = (|b_1|^2 + |b_2|^2)\mathcal{O}(\epsilon_{\text{machine}}^2),$$

which leads to that:

$$\frac{\|\delta\mathbf{b}\|^2}{\|\mathbf{b}\|^2} = \mathcal{O}(\epsilon_{\text{machine}}^2),$$

so the single power norm in $\mathcal{O}(\epsilon_{\text{machine}})$.

This makes us recall the **catastrophic cancellation**, that is when we have $b_1 \simeq b_2$, we have:

$$\frac{\|\tilde{\mathbf{x}} - \mathbf{x}\|}{\|\mathbf{x}\|} = \frac{\|\tilde{f}(\mathbf{b}) - f(\mathbf{b})\|}{\|f(\mathbf{b})\|}$$

be large. However, the backwards stability means that there exists points in the neighborhood such that the differences can be similar. Here, we have b_1 and b_2 close enough to have some point in the neighborhood to be backwards stable.

Definition III.3.5. Stability.

Given a problem f , an algorithm \tilde{f} is stable if for all $\mathbf{b} \in B$, there exists a perturbation of data $\delta\mathbf{b}$ with:

$$\frac{\|\delta\mathbf{b}\|}{\|\mathbf{b}\|} = \mathcal{O}(\epsilon_{\text{machine}}),$$

such that:

$$\frac{\|\tilde{f}(\mathbf{b}) - f(\mathbf{b} + \delta\mathbf{b})\|}{\|f(\mathbf{b} + \delta\mathbf{b})\|} = \mathcal{O}(\epsilon_{\text{machine}}).$$

The stability is a weaker statement than backwards stability.

Proposition III.3.6. Backwards Stability \implies Stability.

A backwards stability algorithm is stable. The *converse* is not necessarily true.

This relaxation is necessary, as there are algorithms that are stable but not backward stable.

Example III.3.7. Stable but not Backward Stable Problem.

Consider the computation of the outer product between 2 vectors $\mathbf{a}, \mathbf{b} \in \mathbb{C}^m$ the mathematical problem is:

$$f : \mathbb{C}^m \times \mathbb{C}^m \rightarrow \mathbb{C}^{m \times m},$$

$$(\mathbf{a}, \mathbf{b}) \mapsto A = \mathbf{a}\mathbf{b}^*.$$

The algorithm is to:

```
input: a, b in C^m
a' := fp(a) =: (a1', a2', ..., am')
b' := fp(b) =: (b1', b2', ..., bm')
~A := [~Aij]
~Aij := ai' (x) comp_conj(bj') := fp(ai' * comp_conj(bj'))
```

Note that:

$$\tilde{f}(\mathbf{a}, \mathbf{b}) = \underbrace{\tilde{A}}_{\text{not rank 1 matrix}} \neq \underbrace{(\mathbf{a} + \delta\mathbf{a})(\mathbf{b} + \delta\mathbf{b})}_{\text{rank 1 matrix}} = f(\mathbf{a} + \delta\mathbf{a}, \mathbf{b} + \delta\mathbf{b}).$$

Hence it is not backwards stable. It can be shown that the algorithm is stable, and we leave this as an exercise to the readers.

It is worth-noting that our conclusion on stability is independent from the choice of vector norm.

Definition III.3.8. Equivalent Norms.

Let X be a normed vector space, any two vector norms $\|\bullet\|_\alpha$ and $\|\bullet\|_\beta$ in X are equivalent if there exist $C_1, C_2 > 0$ such that for all $\mathbf{x} \in X$, we have:

$$C_1 \|\mathbf{x}\|_\beta \leq \|\mathbf{x}\|_\alpha \leq C_2 \|\mathbf{x}\|_\beta.$$

Proposition III.3.9. All Finite Dimensional Norms are Equivalent.

Let X be a finite dimensional normed vector space, all norms are equivalent.

Example III.3.10. $\|\bullet\|_2$ and $\|\bullet\|_\infty$ for are Equivalent.

One can prove that:

$$\|\mathbf{x}\|_\infty \leq \|\mathbf{x}\|_1 \leq \sqrt{m} \|\mathbf{x}\|.$$

There are some consequences of the equivalence norms.

Proposition III.3.11. Squeeze Theorem.

Suppose X is a finite dimensional normed vector space, and let $\|\bullet\|_\alpha$ and $\|\bullet\|_\beta$ be any norms in X . Let a sequence $\{x_n\}_{n=1}^\infty \subset X$ be such that $\|x_n\|_\alpha \rightarrow 0$ as $n \rightarrow \infty$, then $\|x_n\|_\beta \rightarrow 0$ as $n \rightarrow \infty$.

Proof. This is naturally by the above inequality:

$$0 \leq \|x_n\|_\beta \leq \frac{1}{C_1} \|x_n\|_\alpha \rightarrow 0,$$

hence, we naturally have $\|x\|_\beta \rightarrow 0$ as $n \rightarrow \infty$. □

Another consequence is on the backward stability:

Proposition III.3.12. Backward Stability.

Suppose we have $\mathbf{b} \in X$, which is a finite dimensional normed vector space, and $\|\bullet\|_\alpha, \|\bullet\|_\beta$ are two norms in X :

$$\frac{\|\delta \mathbf{b}\|_\beta}{\|\mathbf{b}\|_\beta} = \mathcal{O}(\epsilon_{\text{machine}}).$$

Then, we have:

$$\frac{\|\delta \mathbf{b}\|_\alpha}{\|\mathbf{b}\|_\alpha} = \mathcal{O}(\epsilon_{\text{machine}}).$$

Proof. From the definition, we have:

$$\|\delta \mathbf{b}\|_\alpha \leq C_2 \|\delta \mathbf{b}\|_\beta \text{ and } \frac{1}{\|\mathbf{b}\|_\alpha} \leq \frac{1}{C_1} \cdot \frac{1}{\|\mathbf{b}\|_\beta}.$$

Therefore, we must have:

$$\frac{\|\delta \mathbf{b}\|_\alpha}{\|\mathbf{b}\|_\alpha} \leq \frac{C_2}{C_1} \cdot \frac{\|\delta \mathbf{b}\|_\beta}{\|\mathbf{b}\|_\beta} \leq \frac{C_2}{C_1} C \cdot \epsilon_{\text{machine}},$$

thus $\|\delta \mathbf{b}\|_\alpha / \|\mathbf{b}\|_\alpha = \mathcal{O}(\epsilon_{\text{machine}})$. □

Then, we can think of the conditional number for squared and invertible matrices, that is:

$$\kappa = \|A\| \cdot \|A^{-1}\|.$$

Here, think of the problem of solving $A\mathbf{x} = \mathbf{b}$, the approximation solution is $\tilde{\mathbf{x}}$. Here, we have:

$$\begin{aligned} \text{Error} \quad \delta\mathbf{x} &= \mathbf{x} - \tilde{\mathbf{x}}, \\ \text{Residual} \quad \delta\mathbf{b} &= \mathbf{b} - A\tilde{\mathbf{x}}. \end{aligned}$$

Remark III.3.13. General Condition Number.

For any mathematical problem and algorithm for $f : B \rightarrow X$ as $\tilde{f} : B \rightarrow X$, we can define the condition number as:

$$\kappa = \sup_{\delta\mathbf{b} \in B} \frac{\frac{\|f(\mathbf{b}) - \tilde{f}(\mathbf{b})\|}{\|f(\mathbf{b})\|}}{\frac{\|\delta\mathbf{b}\|}{\|\mathbf{b}\|}}.$$

I.e., κ is the smallest number that is larger than all possible ratios.

Hence:

$$\frac{\|f(\mathbf{b}) - \tilde{f}(\mathbf{b})\|}{\|f(\mathbf{b})\|} \leq \kappa \frac{\|\delta\mathbf{b}\|}{\|\mathbf{b}\|} \text{ for all } \delta\mathbf{b}.$$

Therefore, we note that:

$$\frac{\|\delta\mathbf{b}\|}{\|\mathbf{b}\|} \text{ is small when } \tilde{f} \text{ is backward stable or stable,}$$

hence since κ could be large, stability causes the left-hand-side to be large.

Again, for the same example, we could have a backwards stable, but the relative error is large due to catastrophic cancellation.

Recall that for solving the problem $A\mathbf{x} = \mathbf{b}$ for a $A \in \mathbb{C}^{m \times n}$, we have the QR decomposition based method, that is:

- (i) Let $A = QR$, which is the QR decomposition of A , where Q is orthogonal and R is upper triangular. (Cost: $\sim \frac{2}{3}mn^2 - \frac{2}{3}n^3$).
- (ii) Compute $Q^{-1} \cdot \mathbf{b} = Q^* \cdot \mathbf{b}$. (Cost: $\sim 2mn$).
- (iii) Solve $R\mathbf{x} = Q^* \cdot \mathbf{b}$ with the upper triangular system. (Cost: $\sim n^2$).

Remark III.3.14. QR Decomposition Based Method is Stable.

The QR decomposition method, *i.e.*, the above three steps are backward stable.

Theorem III.3.15. Backward Stability of QR Decomposition.

Let $A = QR$ be the QR decomposition of A , and let $\tilde{Q}\tilde{R}$ be the QR decomposition computed by Householder triangularization. This algorithm is backward stable in the same sense that the computed solution $\tilde{\mathbf{x}}$ has the property:

$$\|(A + \delta A)\tilde{\mathbf{x}} - \mathbf{b}\| = \min, \quad \frac{\|\delta A\|}{\|A\|} = \mathcal{O}(\epsilon_{\text{machine}})$$

for some $\delta A \in \mathbb{C}^{m \times n}$.

III.4 Stability and Gaussian Elimination

Then, we consider the Gauss Elimination.

Remark III.4.1. LU Decomposition.

We may use Gauss elimination to form the LU decomposition, that is for $A \in \mathbb{C}^{m \times m}$, we have $A = L \circ U$,

$$\text{where } L = \begin{bmatrix} * & 0 & \cdots & 0 \\ * & * & \cdots & 0 \\ \vdots & \vdots & \ddots & 0 \\ * & * & \cdots & * \end{bmatrix} \text{ and } U = \begin{bmatrix} * & * & \cdots & * \\ 0 & * & \cdots & * \\ \vdots & \vdots & \ddots & 0 \\ 0 & 0 & \cdots & * \end{bmatrix}.$$

Here, the algorithm lies as follows to solve $A\mathbf{x} = \mathbf{b}$.

- (i) Compute $A = LU$. (Cost: $\sim \frac{2}{3}m^3$).
- (ii) Solve $L\mathbf{y} = \mathbf{b}$ for the lower triangular system for \mathbf{y} . (Cost: $\sim m^2$).
- (iii) Solve $U\mathbf{x} = \mathbf{y}$ for the upper triangular system for \mathbf{x} . (Cost: $\sim m^2$).

Note that $L\mathbf{y} = \mathbf{b}$ implies that $L(U\mathbf{x}) = \mathbf{b}$, so $A\mathbf{x} = \mathbf{b}$ so \mathbf{x} solves $A\mathbf{x} = \mathbf{b}$.

Example III.4.2. UL is Simpler than QR.

The price is it might be instable (could be corrected). Here, we consider:

$$\begin{bmatrix} a_{11} & 0 & 0 \\ a_{21} & a_{22} & 0 \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \cdot \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}.$$

Here, we can consider the solutions trivially as:

$$\begin{aligned} y_1 &= \frac{b_1}{a_{11}}, \\ y_2 &= \frac{b_2 - a_{21}y_1}{a_{22}}, \\ y_3 &= \frac{b_3 - a_{31}y_1 - a_{32}y_2}{a_{33}}. \end{aligned}$$

In particular, we consider $L_{m-1} \cdots L_2 L_1 A = U$, where L_i is the i -th set of elementary row operation, hence for $A = LU$, we have:

$$L = (L_{m-1} \cdots L_2 L_1)^{-1} = L_1^{-1} L_2^{-1} \cdots L_{m-1}^{-1}.$$

Example III.4.3. Generic 4-by-4 Matrix with Gaussian Elimination.

Consider the steps as follows:

$$\begin{array}{cccc} \begin{bmatrix} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{bmatrix} & \xrightarrow{L_1} & \begin{bmatrix} * & * & * & * \\ 0 & * & * & * \\ 0 & * & * & * \\ 0 & * & * & * \end{bmatrix} & \xrightarrow{L_2} & \begin{bmatrix} * & * & * & * \\ 0 & * & * & * \\ 0 & 0 & * & * \\ 0 & 0 & * & * \end{bmatrix} & \xrightarrow{L_3} & \begin{bmatrix} * & * & * & * \\ 0 & * & * & * \\ 0 & 0 & * & * \\ 0 & 0 & 0 & * \end{bmatrix} \\ A & & L_1 A & & L_2 L_1 A & & L_3 L_2 L_1 U \end{array}$$

More trivially for the 2-by-2 case, a matrix $A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$ is simply by $R_2 \leftarrow R_2 + kR_1$, that is $\begin{bmatrix} a & b \\ c + ka & d + kb \end{bmatrix}$, where $k = -c/a$ when $a \neq 0$ to achieve upper triangularization.

Then, we consider:

$$L_1 = \begin{pmatrix} 1 & 0 \\ k & 1 \end{pmatrix},$$

whose inverse is then:

$$L_1^{-1} = \begin{bmatrix} 1 & 0 \\ -k & 1 \end{bmatrix},$$

which is also upper triangular.

For the 3-by-3 case, let:

$$A = \begin{bmatrix} a & b & c \\ d & e & f \\ h & i & j \end{bmatrix},$$

we do the row operation to obtain that:

$$L_1 A = \begin{bmatrix} a & b & c \\ d + k_1 a & e + k_1 b & f + k_1 c \\ g + k_2 a & h + k_2 b & i + k_2 c \end{bmatrix},$$

where $k_1 = -d/a$ and $k_2 = -g/a$. Hence, we write:

$$L_1 = \begin{bmatrix} 1 & 0 & 0 \\ k_1 & 1 & 0 \\ k_2 & 0 & 1 \end{bmatrix},$$

where the inverse is:

$$L_1^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ -k_1 & 1 & 0 \\ -k_2 & 0 & 1 \end{bmatrix}.$$

Here, we let the second table to be:

$$L_1 A = \begin{bmatrix} a & b & c \\ 0 & \alpha & \beta \\ 0 & \gamma & \delta \end{bmatrix}.$$

Here, we use the similar process for row operation. Here we have:

$$L_2(L_1 A) = \begin{bmatrix} a & b & c \\ 0 & \alpha & \beta \\ 0 & \gamma + k_3 \alpha & \delta + k_3 \beta \end{bmatrix},$$

where $k_3 = -\gamma/\alpha$, and we have:

$$L_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & k_3 & 1 \end{bmatrix},$$

where the inverse is:

$$L_2^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -k_3 & 1 \end{bmatrix}.$$

Here, the inverses are:

$$L = L_1^{-1}L_2^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ -k_1 & 1 & 0 \\ -k_2 & 0 & 1 \end{bmatrix} \circ \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -k_3 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ -k_1 & 1 & 0 \\ -k_2 & -k_3 & 1 \end{bmatrix}.$$

There, we may observe that $A = LU = L_1^{-1}L_2^{-2} \cdots L_{m-1}^{-1}$ gives the Gauss Elimination process.

In the generic case, we consider the matrix:

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1m} \\ a_{21} & a_{22} & \cdots & a_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mm} \end{bmatrix}.$$

Consider the first transformation, we have:

$$\ell_{i1} = \frac{a_{i1}}{a_{11}} \text{ for all } i = 2, 3, \dots, m.$$

It is noteworthy to mention that if a_{11} is zero, we want to shuffle the rows. In fact, we would like to rearrange in a manner that has the largest coefficient at the top. But anyways, we have:

$$L_1^{-1} = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ \ell_{21} & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ \ell_{m1} & 0 & \cdots & 1 \end{bmatrix}.$$

For the second step, we assume that we have:

$$L_1A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1m} \\ 0 & a_{22} & \cdots & a_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & a_{m2} & \cdots & a_{mm} \end{bmatrix}.$$

Hence, we consider that:

$$\ell_{i2} = \frac{a_{i2}}{a_{22}}.$$

$$L_2^{-1} = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \ell_{m2} & \cdots & 1 \end{bmatrix}.$$

Remark III.4.4. Decomposition exactly as the product of the entries.

We have the matrix entry exactly as:

$$L_{m-1} \cdots L_2 L_1 = \text{Id}_m + \sum_{i=1}^m (L_i - \text{Id}_m).$$

Now, we consider the Pseudo code for the code segment as:

```

input: A in C(m*m)
output: L in C(m*m), U in C(m*m)
for j = 1, 2, ..., m-1:
    for j = 1, 2, ..., m-1:
        l[i,j] = n[i,j] / u[j,j]
        u[i,j:m] = u[i,j:m] - l[i,j] * u[j,j:m] % (*) Computationally intense
    end
end
end

```

Consider the computational complexity, we have the (*) intensive, that is:

- It multiplies a scalar by a vector of length ℓ .
- It subtracts 2 vectors of length ℓ .

Hence, there are 2ℓ flops.

At the j th step, we have $\ell = m - j + 1$.

Consider that there are less operations needed overtime, we consider the number of operations as:

$$\sum_{j=1}^{m-1} 2\ell(m-j) = 2 \sum_{j=1}^{m-1} (m-j)^2 = 2 \sum_{j=1}^{m-1} j^2 = 2 \cdot \frac{(m-1)m(2(m-1)+1)}{6} \sim \frac{2}{3}m^3.$$

Remark III.4.5. Problem with LU Decomposition.

The LU Decomposition incurs the following issues:

- We may have division by zero, such as $A = \begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}$, where $\ell_{21} = \frac{1}{0}$.
- Stability issues: For $A = \begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}$. It is invertible, and we compute $\kappa(A)$ with respect to the 2-norm, that is σ_2/σ_1 . Consider a diagonalizable matrix, we have $\sigma_1 = |\lambda_1|$ and $\sigma_2 = |\lambda_2|$. For this case, we have the eigenvalues are $\frac{1 \pm \sqrt{5}}{2}$, so we have $\sigma_1 = \frac{1+\sqrt{5}}{2}$ and $\sigma_2 = \frac{\sqrt{5}-1}{2}$, we have:

$$\frac{\sigma_1}{\sigma_2} \simeq 2.618.$$

We consider the following example, then:

Example III.4.6. Unstable for Gauss Elimination.

Let's define $\delta A = \begin{bmatrix} 10^{-20} & 0 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} \epsilon & 0 \\ 0 & 0 \end{bmatrix}$ and here we apply B as:

$$B = A + \delta A = \begin{bmatrix} 10^{-20} & 1 \\ 1 & 1 \end{bmatrix},$$

and so by the Gauss Elimination, we have:

$$L = \begin{bmatrix} 1 & 0 \\ \ell_{21} & 1 \end{bmatrix} \text{ where } \ell_{21} = \frac{1}{10^{-20}} = 10^{20}.$$

Thus, we have:

$$L = \begin{bmatrix} 1 & 0 \\ 10^{20} & 1 \end{bmatrix} \text{ and } U = \begin{bmatrix} 10^{-20} & 1 \\ 0 & 1 - 10^{20} \end{bmatrix}.$$

Here, we assume that $\text{f}_p(10^{20}) = 10^{20}$ and $\text{f}_p(10^{-20}) = 10^{-20}$, then we have machine representation as:

$$\tilde{L} = \begin{bmatrix} 1 & 0 \\ 10^{20} & 1 \end{bmatrix} \text{ and } \tilde{U} = \begin{bmatrix} 10^{-20} & 1 \\ 0 & -10^{20} \end{bmatrix}.$$

Thus, we have the matrix multiplication as:

$$\tilde{L}\tilde{U} = \begin{bmatrix} 1 & 0 \\ 10^{20} & 1 \end{bmatrix} \begin{bmatrix} 10^{-20} & 1 \\ 0 & -10^{20} \end{bmatrix} = \begin{bmatrix} 10^{-20} & 1 \\ 1 & 0 \end{bmatrix}.$$

We may observe that this is very far from the initial A , since we had a large deviation with $A = \begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}$.

Note that this is *not* a Catastrophic cancellation, but rather a *roundoff error*, hence it is unstable.