

# Regularization of Ill-conditioned Linear Systems

Sheima Mohammed Eldirdiri Abueldahab (sheima@aims.ac.za)

شيماء مُحَمَّد الدرديري أبو الذهب

African Institute for Mathematical Sciences (AIMS)

Supervised by Francis Benyah  
University of the Western Cape

June 8, 2007

# Abstract

Ill-conditioned linear systems arise in many applications, for example, in the solution of integral equations, and in the solution of non-linear programming problems. In many application of linear algebra, the need arises to find a good approximation  $\hat{\mathbf{x}}$  to a vector  $\mathbf{x} \in \mathbb{R}^n$  satisfying an approximating equation  $\mathbf{Ax} \approx \mathbf{b}$  with ill-conditioned  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , given  $\mathbf{b} \in \mathbb{R}^m$ . Straightforward the computed solution  $\hat{\mathbf{x}}$  is usually meaningless approximation to  $\mathbf{x}$  due to the error in the right-hand side  $\mathbf{b}$  and the severe ill-conditioning of the matrix  $\mathbf{A}$ . In order to avoid this difficulty, one typically replaces the linear systems  $\mathbf{Ax} = \mathbf{b}$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{b} \in \mathbb{R}^m$ , by a nearby system that is less sensitive to the error in  $\mathbf{b}$  and considers the computed solution of the latter system an approximation of  $\mathbf{x}$ . This replacement is known as regularization. This essay examines various regularization methods for computing stable solution to ill-conditioned linear systems.

# Contents

<b>Abstract</b>	<b>i</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Basic concepts of Linear Systems . . . . .	1
1.2 Vector and Matrix Norms . . . . .	2
<b>2 Ill-Conditioned Linear Systems</b>	<b>4</b>
2.1 Definition of Ill-Conditioned Linear Systems . . . . .	4
2.2 The Condition Number of A Matrix . . . . .	6
2.3 Testing The Accuracy of the Solution of a Linear Systems . . . . .	7
2.4 Classical Ill-Conditioned Systems . . . . .	9
2.4.1 Polynomial Data Fitting: The Vandermonde System . . . . .	9
2.4.2 Approximation of a Function by a Polynomial: The Hilbert System . . . . .	10
<b>3 Least-squares Solutions to Linear Systems</b>	<b>13</b>
3.1 Singular Value Decomposition . . . . .	15
3.2 Solving Linear Systems Using Singular Value Decomposition . . . . .	18
3.2.1 Singular Value Decomposition and Stability of Linear Systems . . . . .	18
3.3 Numerical Rank . . . . .	19
<b>4 Regularization Methods</b>	<b>20</b>
4.1 Rank-deficient and Ill-posed Problems . . . . .	20
4.2 Regularization Methods . . . . .	20
4.2.1 The Tikhonov Method . . . . .	21
4.2.2 Truncated Singular Value Decomposition (TSVD) . . . . .	23
<b>5 Conclusion</b>	<b>26</b>

<b>A Truncated Singular Value Decomposition (TSVD)</b>	<b>27</b>
<b>Bibliography</b>	<b>30</b>

# 1. Introduction

We will begin our study of ill-conditioned linear systems with some basic concepts of linear algebra.

## 1.1 Basic concepts of Linear Systems

A linear system is an equation of the form

$$\mathbf{Ax} = \mathbf{b}, \quad (1.1)$$

where  $\mathbf{A}$  is a known  $m \times n$  matrix called the *coefficient matrix*, given by

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & a_{13} & \cdots & a_{1n} \\ a_{21} & a_{22} & a_{23} & \cdots & a_{2n} \\ a_{31} & a_{32} & a_{33} & \cdots & a_{3n} \\ \vdots & & & & \\ a_{m1} & a_{m2} & a_{m3} & \cdots & a_{mn} \end{bmatrix}, \quad (1.2)$$

$\mathbf{b}$  is a known  $m$ -vector of constants, given by

$$\mathbf{b} = [ b_1 \ b_2 \ b_3 \ \cdots \ b_m ]^T, \quad (1.3)$$

and  $\mathbf{x}$  is an unknown  $n$ -vector, to be determined, given by

$$\mathbf{x} = [ x_1 \ x_2 \ x_3 \ \cdots \ x_n ]^T. \quad (1.4)$$

Linear systems of equations arise from many situations including the following:

1. Solving systems of non-linear equations by Newton's method.
2. Curve fitting or polynomial interpolation.
3. Numerical solutions of differential equations by the finite difference method.
4. Polynomial approximations to continuous functions on a finite interval  $[0, 1]$  results in solving a linear system  $\mathbf{Ax} = \mathbf{b}$ , where  $\mathbf{A}$  is a Hilbert matrix.

For any linear systems the **augmented matrix** given by

$$[\mathbf{A} \mid \mathbf{b}] = \begin{bmatrix} a_{11} & a_{12} & a_{13} & \cdots & a_{1n} & b_1 \\ a_{21} & a_{22} & a_{23} & \cdots & a_{2n} & b_2 \\ a_{31} & a_{32} & a_{33} & \cdots & a_{3n} & b_3 \\ \vdots & & & & & \vdots \\ a_{m1} & a_{m2} & a_{m3} & \cdots & a_{mn} & b_m \end{bmatrix} \quad (1.5)$$

**Theorem 1.1.**<sup>1</sup> Consider the system  $\mathbf{Ax} = \mathbf{b}$ , with coefficient matrix  $\mathbf{A}$  and augmented matrix  $[\mathbf{A} \mid \mathbf{b}]$ , of size  $m \times (n + 1)$ . Then the linear system,

1.  $\mathbf{Ax} = \mathbf{b}$  is inconsistent (i.e., no solution exists) if and only if  $\text{rank}(\mathbf{A}) < \text{rank}[\mathbf{A} \mid \mathbf{b}]$ .
2.  $\mathbf{Ax} = \mathbf{b}$  has a unique solution if and only if  $\text{rank}(\mathbf{A}) = \text{rank}[\mathbf{A} \mid \mathbf{b}] = n$ .
3.  $\mathbf{Ax} = \mathbf{b}$  has infinitely many solutions if and only if  $\text{rank}(\mathbf{A}) = \text{rank}[\mathbf{A} \mid \mathbf{b}] < n$ .

## 1.2 Vector and Matrix Norms

**Definition 1.2. (Vector Norm)** A vector norm is a function  $\mathbf{V} \rightarrow \mathbb{R}$  satisfies the following three conditions:

1.  $\|\mathbf{x}\| > 0$ , for any nonzero vector  $\mathbf{x}$ .
2.  $\|\alpha\mathbf{x}\| = |\alpha| \|\mathbf{x}\|$ , for any real scalar  $\alpha$ .
3.  $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ , for any vectors  $\mathbf{x}$  and  $\mathbf{y}$ .

The three most important vector norms are :

$$\|\mathbf{x}\|_1 = \sum_{i=1}^n |\mathbf{x}_i| = |\mathbf{x}_1| + \cdots + |\mathbf{x}_n|, \quad (1.6)$$

$$\|\mathbf{x}\|_2 = \left[ \sum_{i=1}^n |\mathbf{x}_i|^2 \right]^{\frac{1}{2}} = (|\mathbf{x}_1|^2 + \cdots + |\mathbf{x}_n|^2)^{\frac{1}{2}} = (\mathbf{x}^T \mathbf{x})^{\frac{1}{2}}, \quad (1.7)$$

and

$$\|\mathbf{x}\|_\infty = \max_{1 \leq i \leq n} \{|\mathbf{x}_i|\}. \quad (1.8)$$

Where  $\mathbf{x}$  is  $n$ -dimensional vector

$$\mathbf{x} = [x_1 \quad x_2 \quad \cdots \quad x_n]^T. \quad (1.9)$$

**Example 1.3.** Find  $\|\mathbf{x}\|_\infty$  of the following vector  $\mathbf{x}$ .

$$\mathbf{x} = \begin{bmatrix} 10 \\ -3 \\ 5 \end{bmatrix}. \quad (1.10)$$

---

<sup>1</sup>Theorem 1.1 can be found in Ref. [7]

**Solution 1.3.**

$$\begin{aligned}
\| \mathbf{x} \|_{\infty} &= \max_{1 \leq i \leq 3} | \mathbf{x}_i | \\
&= \max [ | 10 |, | -3 |, | 5 | ] \\
&= \max [10, 3, 5] \\
&= 10.
\end{aligned}$$

**Definition 1.4. (Matrix Norm)** For any real matrix  $\mathbf{A}$ , a matrix norm  $\| \mathbf{A} \|$  is a nonnegative number associated with  $\mathbf{A}$  having the properties

1.  $\| \mathbf{A} \| \geq 0$ ,  $\| \mathbf{A} \| = 0$  iff  $\mathbf{A} = \mathbf{0}$ .
2.  $\| \alpha \mathbf{A} \| = |\alpha| \| \mathbf{A} \|$ , for any scalar  $\alpha$ .
3.  $\| \mathbf{A} + \mathbf{B} \| \leq \| \mathbf{A} \| + \| \mathbf{B} \|$ .
4.  $\| \mathbf{AB} \| \leq \| \mathbf{A} \| \| \mathbf{B} \|$ .

The three most important matrix norms are :

$$\| \mathbf{A} \|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^m | \mathbf{a}_{ij} |. \quad (1.11)$$

$$\| \mathbf{A} \|_2 = [ \sigma_{\max} (\mathbf{A}^T \mathbf{A}) ]^{\frac{1}{2}}, \quad (1.12)$$

$\sigma_{\max}$  is the largest eigenvalue of the positive-semidefinite matrix  $\mathbf{A}^T \mathbf{A}$ .

$$\| \mathbf{A} \|_{\infty} = \max_{1 \leq i \leq m} \sum_{j=1}^n | \mathbf{a}_{ij} |. \quad (1.13)$$

**Example 1.5.** <sup>2</sup>Find  $\| \mathbf{A} \|_{\infty}$  of the following matrix  $\mathbf{A}$ .

$$\mathbf{A} = \begin{bmatrix} 10 & -7 & 0 \\ -3 & 2.099 & 6 \\ 5 & -1 & 5 \end{bmatrix}. \quad (1.14)$$

**Solution 1.5.** The  $\| \mathbf{A} \|_{\infty}$  given by

$$\begin{aligned}
\| \mathbf{A} \|_{\infty} &= \max_{1 \leq i \leq 3} \sum_{j=1}^3 | \mathbf{a}_{ij} | \\
&= \max [ ( | 10 | + | -7 | + | 0 | ), ( | -3 | + | 2.099 | + | 6 | ), ( | 5 | + | -1 | + | 5 | ) ] \\
&= \max [ (10 + 7 + 0), (3 + 2.099 + 6), (5 + 1 + 5) ] \\
&= \max [17, 11.099, 11] \\
&= 17.
\end{aligned}$$

---

<sup>2</sup>Example 1.5 can be found in Ref. [5]

## 2. Ill-Conditioned Linear Systems

This chapter devote to the study of ill-conditioned linear system  $\mathbf{Ax} = \mathbf{b}$ . We will show how small change in the augmented matrix  $[\mathbf{A} \mid \mathbf{b}]$  affect on the computed solution  $\hat{\mathbf{x}}$  of linear system to be ill-conditioned or well-conditioned. Some classical applications to ill-conditioned linear systems will also discuss (e.g., the Hilbert and the Vandermonde matrices). Since many discussion about ill-conditioned linear systems required knowledge about singular value decomposition (SVD), we will discuss this method and we will show how it can be use to solve the linear systems.

### 2.1 Definition of Ill-Conditioned Linear Systems

The linear system  $\mathbf{Ax} = \mathbf{b}$  is said to be **ill-conditioned** if relatively small changes in the entries of augmented matrix  $[\mathbf{A} \mid \mathbf{b}]$  can cause relatively large changes in the solution. If the system  $\mathbf{Ax} = \mathbf{b}$  is ill-conditioned, then the computed solution to  $\mathbf{Ax} = \mathbf{b}$  will generally not be very accurate. On the other hand, the system is said to be **well-conditioned** if relatively small changes in the entries of  $[\mathbf{A} \mid \mathbf{b}]$  results in relatively small changes in the solutions to  $\mathbf{Ax} = \mathbf{b}$ . In that case one should be able to compute solutions more accurately.

**Example 2.1.**<sup>3</sup> Consider the system

$$\begin{bmatrix} 1 & 3 \\ 2 & 5.999 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 4 \\ 7.999 \end{bmatrix}. \quad (2.1)$$

Denoting the system of equations as

$$\mathbf{Ax} = \mathbf{b},$$

and making a small change in the right hand side vector of the equation 2.1, gives

$$\begin{bmatrix} 1 & 3 \\ 2 & 5.999 \end{bmatrix} \begin{bmatrix} x'_1 \\ x'_2 \end{bmatrix} = \begin{bmatrix} 4 \\ 8 \end{bmatrix} \quad (2.2)$$

. The exact solution of equation 2.2, is

$$\mathbf{x}' = \begin{bmatrix} x'_1 \\ x'_2 \end{bmatrix} = \begin{bmatrix} 4 \\ 0 \end{bmatrix}. \quad (2.3)$$

The system 2.1 has exact solution

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}. \quad (2.4)$$

---

<sup>3</sup>Example 2.1 can be found in Ref. [5]



The relative error in  $\mathbf{b}$  is

$$\frac{\|\mathbf{b} - \mathbf{b}'\|_{\infty}}{\|\mathbf{b}\|_{\infty}} = 1.250 \times 10^{-4}.$$

However, the corresponding relative error in the solution vector is

$$\frac{\|\mathbf{x} - \mathbf{x}'\|_{\infty}}{\|\mathbf{x}\|_{\infty}} = 3.0.$$

We see the small relative error in the right hand side vector of  $1.250 \times 10^{-4}$  results in a large relative error in the solution vector as 3.0.

This system is **ill-conditioned** since a small change in the right hand side resulted in a large change in the solution vector.

**Example 2.2.**<sup>4</sup> Consider the system

$$\begin{bmatrix} 1 & 2 \\ 2 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 4.001 \\ 7.001 \end{bmatrix}. \quad (2.5)$$

Denoting the system of equations as

$$\mathbf{Ax} = \mathbf{b},$$

and making a small change in the right hand side vector of the equation 2.5, gives

$$\begin{bmatrix} 1 & 2 \\ 2 & 3 \end{bmatrix} \begin{bmatrix} x'_1 \\ x'_2 \end{bmatrix} = \begin{bmatrix} 4 \\ 7 \end{bmatrix}. \quad (2.6)$$

The exact solution of equation 2.6, is

$$\mathbf{x} = \begin{bmatrix} x'_1 \\ x'_2 \end{bmatrix} = \begin{bmatrix} 2 \\ 1 \end{bmatrix}. \quad (2.7)$$

The system 2.5 has exact solution

$$\mathbf{x}' = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1.999 \\ 1.001 \end{bmatrix}. \quad (2.8)$$

The relative error in  $\mathbf{b}$  is

$$\frac{\|\mathbf{b} - \mathbf{b}'\|_{\infty}}{\|\mathbf{b}\|_{\infty}} = \frac{0.001}{7.001} = 1.4284 \times 10^{-4}. \quad (2.9)$$

However, the corresponding relative error in the solution vector is

$$\frac{\|\mathbf{x} - \mathbf{x}'\|_{\infty}}{\|\mathbf{x}\|_{\infty}} = \frac{0.001}{1.9990} = 5.0025 \times 10^{-4}.$$

---

<sup>4</sup>Example 2.2 can be found in Ref. [5]

We see the small relative error in the right hand side vector of  $1.4284 \times 10^{-4}$  results in a small relative error in the solution vector as  $5.0025 \times 10^{-4}$ .

This system is **well-conditioned** as a small changes in the right hand side resulted in small change in the solution vector.

## 2.2 The Condition Number of A Matrix

The condition number is a measure of how close a matrix  $\mathbf{A}$  is to being singular.

**Definition 2.3. (Condition Number)** Let  $\mathbf{A}$  be an  $n \times n$  nonsingular matrix. The condition number of  $\mathbf{A}$ , defined by  $\text{cond}(\mathbf{A})$ , denoted by  $\text{cond}(\mathbf{A})$ , is defined as

$$\text{cond}(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|. \quad (2.10)$$

Where  $\|\cdot\|$  is a given matrix norm.

**Example 2.4.**<sup>5</sup> Let

$$\mathbf{A} = \begin{bmatrix} 3 & 3 \\ 4 & 5 \end{bmatrix},$$

$$\mathbf{A}^{-1} = \frac{1}{3} \begin{bmatrix} 5 & -3 \\ -4 & 3 \end{bmatrix}.$$

And

$$\|\mathbf{A}\|_{\infty} = 9, \quad \|\mathbf{A}^{-1}\|_{\infty} = \frac{8}{3},$$

Then

$$\text{cond}(\mathbf{A}) = 9 \times \frac{8}{3} = 24.$$

If  $\text{cond}(\mathbf{A})$  is small (close to 1) then the matrix is said to be **well-conditioned**. On the other hand, if  $\text{cond}(\mathbf{A})$  is large, that is, if it is significantly larger than one, then the matrix is said to be **ill-conditioned**. The condition number of a matrix  $\mathbf{A}$  associated with the linear system  $\mathbf{Ax} = \mathbf{b}$  gives a bound on how inaccurate the solution  $\mathbf{x}$  will be after approximate solution. In particular, one should think of the condition number as being the rate at which the solution,  $\mathbf{x}$ , will change with respect to a change in  $\mathbf{b}$ . Thus, if the condition number is large, even a small error in  $\mathbf{b}$  may cause a large error in  $\mathbf{x}$ . On the other hand, if the condition number is small then the error in  $\mathbf{x}$  will not be much bigger than the error in  $\mathbf{b}$ .

<sup>5</sup>Example 2.4 can be found in Ref. [2]

## 2.3 Testing The Accuracy of the Solution of a Linear Systems

Once a solution  $\hat{\mathbf{x}}$  of the system  $\mathbf{Ax} = \mathbf{b}$  has been computed, it is natural to test how accurate the computed solution  $\hat{\mathbf{x}}$  is. If the exact solution  $\mathbf{x}$  is known, then one could of course compute the relative error  $\|\mathbf{x} - \hat{\mathbf{x}}\| / \|\mathbf{x}\|$ , in  $\mathbf{x}$ . However, in most practical situations, the exact solution is not known. In such cases, the most obvious thing to do is to compute the residual  $\mathbf{r} = \mathbf{b} - \mathbf{A}\hat{\mathbf{x}}$  and see how small the relative residual  $\|\mathbf{r}\| / \|\mathbf{b}\|$  is. Unfortunately, a small relative residual does not guarantee the accuracy of the solution. The following example illustrates this fact.

**Example 2.5.** <sup>6</sup> Consider the linear system  $\mathbf{Ax} = \mathbf{b}$  given by

$$\begin{bmatrix} 1 & 2 \\ 1.0001 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 3 \\ 3.0001 \end{bmatrix}.$$

Let

$$\hat{\mathbf{x}} = (3, 0)^T. \quad (2.11)$$

Then

$$\mathbf{r} = \mathbf{b} - \mathbf{A}\hat{\mathbf{x}} = \begin{bmatrix} 3 \\ 3.0001 \end{bmatrix} - \begin{bmatrix} 1 & 2 \\ 1.0001 & 2 \end{bmatrix} \begin{bmatrix} 3 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0.0002 \end{bmatrix},$$

and relative residual

$$\frac{\|\mathbf{r}\|}{\|\mathbf{b}\|} = 0.000066664, \quad (2.12)$$

which is small, even though the solution  $\hat{\mathbf{x}} = (3, 0)^T$  is nowhere near the exact solution

$$\mathbf{x} = (1, 1)^T. \quad (2.13)$$

The above phenomena can be explained by the following Theorem.

**Theorem 2.6.** <sup>7</sup> (**The Residual theorem**) *let  $\hat{\mathbf{x}}$  be the computed solution to the linear system  $\mathbf{Ax} = \mathbf{b}$ . Then*

$$\frac{\|\mathbf{x} - \hat{\mathbf{x}}\|}{\|\mathbf{x}\|} \leq \|\mathbf{A}\| \|\mathbf{A}^{-1}\| \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|}, \quad (2.14)$$

or

$$\frac{\|\mathbf{x} - \hat{\mathbf{x}}\|}{\|\mathbf{x}\|} \leq \text{cond}(\mathbf{A}) \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|}. \quad (2.15)$$

<sup>6</sup>Example 2.5 can be found in Ref. [2]

<sup>7</sup>Theorem 2.6 can be found in Ref. [2]

**Proof 2.6.**

$$\begin{aligned}\mathbf{r} &= \mathbf{b} - \mathbf{A}\hat{\mathbf{x}} \\ &= \mathbf{A}\mathbf{x} - \mathbf{A}\hat{\mathbf{x}} \\ &= \mathbf{A}(\mathbf{x} - \hat{\mathbf{x}}),\end{aligned}\tag{2.16}$$

we have

$$\mathbf{x} - \hat{\mathbf{x}} = \mathbf{A}^{-1}\mathbf{r} \quad (\text{since } \mathbf{A} \text{ is nonsingular})$$

taking norms give,

$$\|\mathbf{x} - \hat{\mathbf{x}}\| \leq \|\mathbf{A}^{-1}\| \|\mathbf{r}\|.\tag{2.17}$$

Also from

$$\mathbf{b} = \mathbf{A}\mathbf{x},\tag{2.18}$$

we have

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

that is,

$$\frac{1}{\|\mathbf{x}\|} \leq \frac{\|\mathbf{A}\|}{\|\mathbf{b}\|}.\tag{2.19}$$

This combines with 2.17 to give

$$\frac{\|\mathbf{x} - \hat{\mathbf{x}}\|}{\|\mathbf{x}\|} \leq \|\mathbf{A}\| \|\mathbf{A}^{-1}\| \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|},\tag{2.20}$$

or

$$\frac{\|\mathbf{x} - \hat{\mathbf{x}}\|}{\|\mathbf{x}\|} \leq \text{cond}(\mathbf{A}) \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|}.\tag{2.21}$$

The above Theorem 2.6 tells that the relative error in the computed solution  $\hat{\mathbf{x}}$  depends not only on the relative residual, but also on the quantity  $\text{cond}(\mathbf{A})$ . A computed solution can be guaranteed to be accurate only when the product  $\text{cond}(\mathbf{A}) \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|}$  is small.

In previous Example 2.5, we get

$$\text{cond}(\mathbf{A}) = 5.0001 \times 10^4 \quad (\text{large}),\tag{2.22}$$

and the relative residual

$$\frac{\|\mathbf{r}\|}{\|\mathbf{b}\|} = 0.000066664.$$

So

$$\begin{aligned}\frac{\|\mathbf{x} - \hat{\mathbf{x}}\|}{\|\mathbf{x}\|} &\leq \text{cond}(\mathbf{A}) \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|} = (5.0001 \times 10^4) \cdot (0.000066664) \\ &= 3.3333 \quad (\text{large}).\end{aligned}$$

## 2.4 Classical Ill-Conditioned Systems

Ill-conditioned linear systems occur in many applications of mathematics. We give some examples below.

### 2.4.1 Polynomial Data Fitting: The Vandermonde System

The interpolating  $n$ -th degree polynomial  $p_n(x)$  at given  $(n + 1)$  data points  $(x_i, y_i)$ , is obtained by solving the  $(n + 1) \times (n + 1)$  linear system

$$y_i = a_0 + a_1x_i + a_2x_i^2 + a_3x_i^3 + \dots + a_nx_i^n, \quad i = 0, 1, \dots, n \quad \text{for } a_0, a_1, \dots, a_n.$$

In matrix form we have

$$\begin{bmatrix} 1 & x_0 & x_0^2 & \cdots & x_0^n \\ 1 & x_1 & x_1^2 & \cdots & x_1^n \\ 1 & x_2 & x_2^2 & \cdots & x_2^n \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \cdots & x_n^n \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} = \begin{bmatrix} y_0 \\ y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}. \quad (2.23)$$

The coefficient matrix of the system (2.23) is called the **Vandermonde matrix**. If the  $x_i$ 's are distinct, then the Vandermonde matrix is nonsingular, and so the system has a unique solution. That is, there is a unique polynomial that interpolates the given points.

For a small  $n$  the Vandermonde system can be solved quite easily. However, as  $n$  increases the system becomes increasingly ill-conditioned.

**Example 2.7.**<sup>8</sup> Let the  $x_i$ 's be the  $(n + 1)$  equally-spaced points in the interval  $[0, 1]$ . For example the Vandermonde matrix for  $n = 4$  is

$$\mathbf{V}_4 = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 1 & (1/4) & (1/4)^2 & (1/4)^3 & (1/4)^4 \\ 1 & (1/2) & (1/2)^2 & (1/2)^3 & (1/2)^4 \\ 1 & (3/4) & (3/4)^2 & (3/4)^3 & (3/4)^4 \\ 1 & 1 & 1 & 1 & 1 \end{bmatrix} \quad (2.24)$$

The condition number,  $\text{cond}(\mathbf{V}_4) = 686.4349$ . The table below gives the condition numbers of the Vandermonde matrices for  $n = 4, 5, 6, 7, 8, 9, 10, 11, 12$ .

<sup>8</sup>Example 2.7 can be found in Ref. [2]

$n$	$\text{cond}(\mathbf{V}_n)$
4	6.8643e+02
5	4.9244e+03
6	3.6061e+04
7	2.6782e+05
8	2.0094e+06
9	1.5193e+07
10	1.1558e+08
11	8.8348e+08
12	6.7806e+09

Table 1: Condition Numbers of Vandermonde Matrices.

## 2.4.2 Approximation of a Function by a Polynomial: The Hilbert System

**Definition 2.8.**<sup>9</sup> **(Hilbert Matrix)** The  $n \times n$  matrix  $\mathbf{H}_n$ , with entries  $h_{ij} = \frac{1}{i+j-1}$ ,  $1 \leq i \leq n$ ,  $1 \leq j \leq n$  is called the Hilbert matrix of order  $n$ .

The Hilbert matrix arises in least-squares polynomial approximation of continuous functions on the interval  $[0, 1]$ , using the standard basis  $\{1, x, x^2, \dots, x^n\}$ , for  $\mathbb{P}^n$ . Suppose a continuous function  $f(x)$  defined on the interval  $[0, 1]$  is to be approximated by a polynomial of degree  $n-1$ :

$$p_{n-1}(x) = \sum_{i=1}^n a_i x^{i-1}, \quad (2.25)$$

such that the error

$$E = \|p_{n-1} - f\|_2^2 = \int_0^1 \left[ \sum_{i=1}^n a_i x^{i-1} - f(x) \right]^2 dx \quad (2.26)$$

is minimized. The coefficients  $a_i$  of the polynomial are easily determined by setting

$$\frac{\partial E}{\partial a_i} = 0, \quad i = 1, 2, \dots, n$$

$$\frac{\partial E}{\partial a_i} = 2 \int_0^1 \left[ \sum_{j=1}^n a_j x^{j-1} - f(x) \right] x^{i-1} dx = 0, \quad i = 1, \dots, n$$

or

$$\sum_{j=1}^n a_j \int_0^1 x^{i+j-2} dx = \int_0^1 f(x) x^{i-1} dx, \quad i = 1, \dots, n.$$

<sup>9</sup>Definition 2.8 can be found in Ref. [2]

To obtain the latter form, we have interchanged the summation and integration. Letting

$$h_{ij} = \int_0^1 x^{i+j-2} dx \quad (2.27)$$

and

$$b_i = \int_0^1 x^{i-1} dx, \quad i = 1, \dots, n. \quad (2.28)$$

We have

$$\sum_{j=1}^n h_{ij} a_j = b_i, \quad i = 1, \dots, n. \quad (2.29)$$

This is exactly the linear system  $\mathbf{H}\mathbf{a} = \mathbf{b}$ , given by

$$\begin{bmatrix} h_{11} & h_{12} & \cdots & h_{1,n} \\ h_{21} & h_{22} & \cdots & h_{2,n} \\ \vdots & \vdots & & \vdots \\ h_{n,1} & h_{n,2} & \cdots & h_{n,n} \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix},$$

$$\text{With } \mathbf{H} = [h_{ij}], \quad \mathbf{a} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} \quad \text{and} \quad \mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix}.$$

For example, the Hilbert matrix of order 5, denoted by  $\mathbf{H}_5$  is given by

$$\mathbf{H}_5 = \begin{bmatrix} 1 & 1/2 & 1/3 & 1/4 & 1/5 \\ 1/2 & 1/3 & 1/4 & 1/5 & 1/6 \\ 1/3 & 1/4 & 1/5 & 1/6 & 1/7 \\ 1/4 & 1/5 & 1/6 & 1/7 & 1/8 \\ 1/5 & 1/6 & 1/7 & 1/8 & 1/9 \end{bmatrix}. \quad (2.30)$$

The inverse of  $\mathbf{H}_5$  is given by

$$\mathbf{H}_5^{-1} = \begin{bmatrix} 25 & -300 & 1050 & -1400 & 630 \\ -300 & 4800 & -18900 & 26880 & -12600 \\ 1050 & -18900 & 79380 & -117600 & 56700 \\ -1400 & 26880 & -117600 & 179200 & -88200 \\ 630 & -12600 & 56700 & -88200 & 44100 \end{bmatrix}. \quad (2.31)$$

The entries are large compared with the entries of  $\mathbf{H}_5$ . Therefore some perturbation in the right-hand vector  $\mathbf{b}$ , gets magnified when multiplied by  $\mathbf{H}_5^{-1}$ . The condition number of  $\mathbf{H}_5$  is  $\text{cond}(\mathbf{H}_5) = 4.7661e + 05$  which increases rapidly with  $n$ .

---

$n$	$\text{cond}(\mathbf{H}_n)$
5	4.766072e+05
6	1.495106e+07
7	4.753674e+08
8	1.525758e+10
9	4.931544e+11
10	1.602529e+13
11	5.223946e+14
12	1.794510e+16

Table 2: Condition Numbers of Hilbert Matrices

This table shows the condition numbers of the Hilbert matrices for  $n = 5, 6, 7, 8, 9, 10, 11, 12$ . We see the value of the condition number increase when  $n$  also increase.



### 3. Least-squares Solutions to Linear Systems

In several practical situations, we need to solve a linear system  $\mathbf{Ax} = \mathbf{b}$  where the matrix  $\mathbf{A}$  is nonsquare and/or singular. In such cases, solutions may not exist at all; in cases where there are solutions, there may be infinitely many solutions. For example, when  $\mathbf{A}$  is  $m \times n$  and  $m > n$ , we have an overdetermined system (that is, the number of equations is greater than the number of unknowns), and an overdetermined system typically has no solution. In contrast, an underdetermined system ( $m < n$ ) typically has an infinite number of solutions.

In these cases, the best we can hope for is to find a vector  $\mathbf{x}$  that will make  $\mathbf{Ax}$  as close as possible to the vector  $\mathbf{b}$ . In other words, we seek a vector  $\mathbf{x}$  such that  $\|\mathbf{r}\| = \|\mathbf{Ax} - \mathbf{b}\|$  is minimized. When the 2-norm is used, this solution is referred to as a **least-squares solution** to the linear system  $\mathbf{Ax} = \mathbf{b}$ . The problem of finding least-squares solutions to the linear system  $\mathbf{Ax} = \mathbf{b}$  is known as the **linear least-squares problem** (LSP). The linear least-squares problem is formally defined as follows.

**Definition 3.1.**<sup>10</sup> Given an  $m \times n$  matrix  $\mathbf{A}$  and  $m$ -vector  $\mathbf{b}$ , the least-squares problem is to find an  $n$ -vector  $\mathbf{x}$  such that the norm of the residual vector,  $\|\mathbf{Ax} - \mathbf{b}\|_2$ , is as small as possible.

If the least-squares problem has more than one solution, the one having the 2-norm is called the **minimum norm solution**.

The least-squares solution  $\hat{\mathbf{x}}$ , for the the linear system

$$\mathbf{Ax} = \mathbf{b}, \tag{3.1}$$

is given by minimizing the 2-norm square of the residual  $\mathbf{Ax} - \mathbf{b}$ , that is, the quantity

$$\|\mathbf{Ax} - \mathbf{b}\|_2^2 = (\mathbf{Ax} - \mathbf{b})^T (\mathbf{Ax} - \mathbf{b}). \tag{3.2}$$

So

$$\|\mathbf{Ax} - \mathbf{b}\|_2^2 = (\mathbf{Ax})^T (\mathbf{Ax}) - \mathbf{b}^T \mathbf{Ax} - (\mathbf{Ax})^T \mathbf{b} + \mathbf{b}^T \mathbf{b}. \tag{3.3}$$

The two middle terms are equal, and the minimum is found at the zero of the derivative with respect to  $\mathbf{x}$ ,

$$2\mathbf{A}^T \mathbf{Ax} - 2\mathbf{A}^T \mathbf{b} = 0. \tag{3.4}$$

Then

$$\mathbf{A}^T \mathbf{Ax} = \mathbf{A}^T \mathbf{b}. \tag{3.5}$$

Thus the least-squares solution is

$$\hat{\mathbf{x}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b}. \tag{3.6}$$

---

<sup>10</sup>Definition 3.1 can be found in Ref. [3]

**Example 3.2.** <sup>11</sup> Find the least-squares solution of  $\mathbf{Ax} = \mathbf{b}$ , where

$$\mathbf{A} = \begin{bmatrix} 1 & 274 & 2450 \\ 1 & 180 & 3254 \\ 1 & 375 & 3802 \\ 1 & 205 & 2838 \\ 1 & 86 & 2347 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 162 \\ 120 \\ 223 \\ 131 \\ 67 \end{bmatrix}. \quad (3.7)$$

**Solution 3.2.** this least-squares solution is given by

$$\text{Step 1: Form } \mathbf{A}^T \mathbf{A} = \begin{bmatrix} 5 & 1120 & 14,691 \\ 1120 & 297,522 & 3466,420 \\ 14,691 & 3466,402 & 44608,873 \end{bmatrix}. \quad (3.8)$$

$$\text{Step 2: Form } \mathbf{A}^T \mathbf{b} = \begin{bmatrix} 703 \\ 182,230 \\ 2164,253 \end{bmatrix}. \quad (3.9)$$

$$\text{Step 3: Solve the equation: } \mathbf{A}^T \mathbf{Ax} = \mathbf{A}^T \mathbf{b}. \quad (3.10)$$

Then

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 7.0325 \\ 0.5044 \\ 0.0070 \end{bmatrix}. \quad (3.11)$$

This system is artificially ill-conditioned because the column of  $\mathbf{A}$  are out of scale.

$$\text{cond}(\mathbf{A}) = 1.7576 \times 10^4 \quad (\text{large}). \quad (3.12)$$

$$\text{cond}(\mathbf{A}^T \mathbf{A}) = 3.0891 \times 10^8 \quad (\text{large}). \quad (3.13)$$

To see how the computed least-squares solution agree with the data of the vector  $\mathbf{b}$ , we compute

$$\begin{bmatrix} 1 & 274 & 2450 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = 162.4043. \quad (3.14)$$

(True value = 162)

$$\begin{bmatrix} 1 & 180 & 3254 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = 120.6153. \quad (3.15)$$

(True value = 120)

---

<sup>11</sup>Example 3.2 can be found in Ref. [3]

$$\begin{bmatrix} 1 & 375 & 3802 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = 222.8193. \quad (3.16)$$

(True value = 223)

## 3.1 Singular Value Decomposition

The **singular value decomposition** (SVD) of a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  is of great theoretical and practical importance. It is a matrix factorization with applications in many algorithms. Its history goes back more than a century, but its use in numerical computations is much more recent [1].

**Definition 3.3. (Singular Value Decomposition)** For any real  $m \times n$  matrix  $\mathbf{A}$ , there exists orthogonal matrices an  $m \times m$  matrix  $\mathbf{U}$  and  $n \times n$  matrix  $\mathbf{V}$ , and  $\mathbf{\Sigma}$  is  $m \times n$  “diagonal” matrix, with entries  $\sigma_1, \sigma_2, \dots, \sigma_n$  called the singular values of  $\mathbf{A}$  such that

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T. \quad (3.17)$$

**Comment:**  $\mathbf{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_n)$  has non-negative diagonal elements appearing in non-increasing order such that  $\sigma_1 \geq \dots \geq \sigma_n \geq 0$ .

From above equation 3.17 we can get the form

$$\mathbf{A}\mathbf{V} = \mathbf{U}\mathbf{\Sigma}, \quad (3.18)$$

such that

$$\mathbf{A}[\mathbf{v}_1 \mathbf{v}_2 \dots \mathbf{v}_n] = [\mathbf{u}_1 \mathbf{u}_2 \dots \mathbf{u}_n] \begin{bmatrix} \sigma_1 & 0 & \dots & 0 \\ 0 & \sigma_2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & \sigma_n \\ 0 & \dots & \dots & 0 \\ \vdots & & & \vdots \\ 0 & \dots & \dots & 0 \end{bmatrix}. \quad (3.19)$$

So

$$\mathbf{A}\mathbf{v}_i = \sigma_i \mathbf{u}_i \quad i = 1, 2, \dots, n. \quad (3.20)$$

Where the vectors  $\mathbf{u}_i$  and  $\mathbf{v}_i$  are called the right and left singular vectors of  $\mathbf{A}$ , respectively.

$$\mathbf{A} = \sum_{i=1}^n \sigma_i \mathbf{u}_i \mathbf{v}_i^T \quad (3.21)$$

and

$$\mathbf{A}^T = \sum_{i=1}^n \sigma_i^T \mathbf{v}_i \mathbf{u}_i^T, \quad (3.22)$$

where  $n = \text{numerical rank}(\mathbf{A})$  and  $\mathbf{u}_i$  and  $\mathbf{v}_i$  are respectively the  $i$ -th columns of  $\mathbf{U}$  and  $\mathbf{V}$ .

**Comment:** The numerical rank of  $\mathbf{A}$  in this case is equal to the number of non-zero singular values.

**Theorem 3.4.** <sup>12</sup>Let  $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$  be the singular value decomposition of an  $m \times n$  matrix  $\mathbf{A}$  ( $m \geq n$ ). Let  $r$  be the rank of the matrix  $\mathbf{A}$ . Then

1.  $\mathbf{V}^T(\mathbf{A}^T\mathbf{A})\mathbf{V} = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_r^2, 0, \dots, 0)_{n \times n}$ .
2.  $\mathbf{U}^T(\mathbf{A}\mathbf{A}^T)\mathbf{U} = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_r^2, 0, \dots, 0)_{m \times m}$ .

From this Theorem 3.4 it immediately follows that

1. The left singular vectors  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$  are the unit eigenvectors of the matrix  $\mathbf{A}^T\mathbf{A}$ .
2. The right singular vectors  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n$  are the unit eigenvectors of the matrix  $\mathbf{A}\mathbf{A}^T$ .
3.  $\sigma_1^2, \dots, \sigma_r^2$  are the nonzero eigenvalues of both  $\mathbf{A}^T\mathbf{A}$  and  $\mathbf{A}\mathbf{A}^T$ .

**Example 3.5.** Let

$$\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 0 & 0 \\ 1 & 1 \end{bmatrix}. \quad (3.23)$$

The singular value decomposition gives

$$\Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}. \quad (3.24)$$

Then

$$\mathbf{A}\mathbf{A}^T = \begin{bmatrix} 2 & 0 & 2 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}, \quad (3.25)$$

and it's corresponding vector  $\mathbf{U}$  is

$$\mathbf{U} = \begin{bmatrix} 0.00000 & 0.70711 & 0.70711 \\ -1.00000 & -0.00000 & 0.00000 \\ 0.00000 & -0.70711 & 0.70711 \end{bmatrix}. \quad (3.26)$$

---

<sup>12</sup>Theorem 3.4 can be found in Ref. [3]

Also

$$\mathbf{A}^T \mathbf{A} = \begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix}, \quad (3.27)$$

and it's corresponding vector  $\mathbf{V}$  is

$$\mathbf{V} = \begin{bmatrix} -0.70711 & 0.70711 \\ -0.70711 & -0.70711 \end{bmatrix}. \quad (3.28)$$

So

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \begin{bmatrix} 1 & 1 \\ 0 & 0 \\ 1 & 1 \end{bmatrix}. \quad (3.29)$$

**Example 3.6.** Let

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & -1 \\ 4 & 3 & 1 \\ 6 & 7 & -1 \end{bmatrix}. \quad (3.30)$$

The singular value decomposition gives

$$\mathbf{\Sigma} = \begin{bmatrix} 10.67963 & 0.00000 & 0.00000 \\ 0.00000 & 1.98632 & 0.00000 \\ 0.00000 & 0.00000 & 0.00000 \end{bmatrix}. \quad (3.31)$$

Then

$$\mathbf{A}\mathbf{A}^T = \begin{bmatrix} 6 & 9 & 21 \\ 9 & 26 & 44 \\ 21 & 44 & 86 \end{bmatrix}, \quad (3.32)$$

and it's corresponding vector  $\mathbf{U}$  is

$$\mathbf{U} = \begin{bmatrix} -0.20627 & 0.53925 & -0.81650 \\ -0.45417 & -0.79187 & -0.40825 \\ -0.86671 & 0.28662 & 0.40825 \end{bmatrix}. \quad (3.33)$$

Also

$$\mathbf{A}^T \mathbf{A} = \begin{bmatrix} 53 & 56 & -3 \\ 56 & 62 & -6 \\ -3 & -6 & 3 \end{bmatrix}, \quad (3.34)$$

and it's corresponding vector  $\mathbf{V}$  is

$$\mathbf{V} = \begin{bmatrix} -0.676353 & -0.457399 & 0.577350 \\ -0.734296 & 0.357039 & -0.577350 \\ 0.057943 & -0.814438 & -0.577350 \end{bmatrix}. \quad (3.35)$$

So

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \begin{bmatrix} 1 & 2 & -1 \\ 4 & 3 & 1 \\ 6 & 7 & -1 \end{bmatrix}. \quad (3.36)$$

## 3.2 Solving Linear Systems Using Singular Value Decomposition

The idea of using **SVD** in the solution of the least-squares problem is to determine if a linear system  $\mathbf{Ax} = \mathbf{b}$  has a solution and, if so, how to compute it.

The least squares solution to the problem  $\min \|\mathbf{Ax} - \mathbf{b}\|_2$  is given by

$$\hat{\mathbf{x}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b}, \quad (3.37)$$

substituting  $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$ , we have

$$\hat{\mathbf{x}} = \mathbf{V}\Sigma^+ \mathbf{U}^T \mathbf{b} \quad (3.38)$$

$$\begin{aligned} \hat{\mathbf{x}} &= \sum_{i=1}^n \sigma_i^+ \mathbf{u}_i^T \mathbf{b} \mathbf{v}_i \\ &= \sum_{i=1}^n \frac{\mathbf{u}_i^T \mathbf{b}}{\sigma_i} \mathbf{v}_i, \end{aligned} \quad (3.39)$$

where  $n =$  the numerical  $\mathbf{rank}(\mathbf{A})$  and  $\mathbf{u}_i$  and  $\mathbf{v}_i$  are respectively the  $i$ -th columns of  $\mathbf{U}$  and  $\mathbf{V}$ , [3].

### 3.2.1 Singular Value Decomposition and Stability of Linear Systems

We will now analyze the influence of the coefficient matrix of the linear system  $\mathbf{Ax} = \mathbf{b}$  on the solution using singular value decomposition (SVD) [9].

$$\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T = \mathbf{A} = \sum_{i=1}^n \sigma_i \mathbf{u}_i \mathbf{v}_i^T, \quad (3.40)$$

and

$$\hat{\mathbf{x}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b} = \sum_{i=1}^n \sigma_i^+ \mathbf{u}_i^T \mathbf{b} \mathbf{v}_i. \quad (3.41)$$

Hence if  $\sigma_L$  (The largest singular value of  $\mathbf{A}$ ) is small then small changes in  $\mathbf{A}$  or  $\mathbf{b}$  will cause a significant change in  $\mathbf{x}$ . As we showed before in section 2.2, the conditioned number of a matrix  $\mathbf{A}$  is defined as:

$$\text{cond}(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|, \quad (3.42)$$

and for the matrix 2-norm,

$$\text{cond}(\mathbf{A}) = \|\mathbf{A}\|_2 \|\mathbf{A}^{-1}\|_2 = \frac{\sigma_L(\mathbf{A})}{\sigma_S(\mathbf{A})}. \quad (3.43)$$

Where  $\sigma_L$  and  $\sigma_S$  is the largest and smallest singular values of  $\mathbf{A}$  respectively.

( $\|\mathbf{A}\|_2$  is the square root of the largest eigenvalue of  $\mathbf{A}^T \mathbf{A}$  which is the largest singular value of  $\mathbf{A}$ , see equation 1.12).

### 3.3 Numerical Rank

In practical application that need singular values, we have to know when to accept a computed singular value to be “zero”. Accept a computed singular value to be zero if it is less than or equal to  $10^{-t} \|\mathbf{A}\|_\infty$ , where the entries of  $\mathbf{A}$  are correct to  $t$  digits. Having defined a tolerance  $\delta = 10^{-t} \|\mathbf{A}\|_\infty$  for a zero singular value, we can make the following convention (Golub and Van Loan 1989, 247):

$\mathbf{A}$  has “numerical rank”  $n$  if the computed singular values  $\sigma_1, \sigma_2, \dots, \sigma_r$  satisfy

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \delta \geq \sigma_{n+1} \geq \sigma_r \quad (3.44)$$

Thus, to determine the numerical rank of a matrix  $\mathbf{A}$ , count the “large” singular values only. If this number is  $n$ , then  $\mathbf{A}$  has numerical rank  $n$ .

**Remark:** Note that finding the numerical rank of a matrix will be “tricky” if there is no suitable gap between a set of singular values [3].

## 4. Regularization Methods

In this chapter, we give a brief introduction to Rank-deficient and discrete ill-posed problems. We introduce the regularization methods to improve the accuracy of the solution of ill-conditioned linear systems. We focussed in two methods, e.g. truncated singular value decomposition (**TSVD**) and Tikhonov methods.

### 4.1 Rank-deficient and Ill-posed Problems

The numerical treatment of system of equations with an ill-conditioned coefficient matrix depends on the type of ill-conditioning of matrix  $\mathbf{A}$ .

There are two important classes of problems to consider, and many practical problems belong to one of these two classes [4].

1. **Rank-deficient problems** are characterized by the matrix  $\mathbf{A}$  having a cluster of small singular values, and there is well-determined gap between large and small singular values. This implies that one or more rows and columns of  $\mathbf{A}$  are nearly linear combinations of some or all of the remaining rows and columns.
2. **Discrete ill-posed problem.** Here, all the the singular values of  $\mathbf{A}$  gradually decay to zero, also the ratio between the largest and the smallest nonzero singular values is large and there is no notion of a numerical rank for these matrices. The primary difficulty with the discrete ill-posed problems is that they are essentially underdetermined due to the cluster of small singular values of the matrix  $\mathbf{A}$ . Hence it is necessary to incorporate further information about the desired solution in order to stabilise the problem and to single out a useful and stable solution. This is the purpose of regularization.

### 4.2 Regularization Methods

One way of getting a useful solution to discrete ill-posed or rank-deficient linear problems is by regularization. The main idea of regularization is to find a new problem or method that damps the effect of the noise in the input data and also to make a solution of ill-conditioning linear system  $\mathbf{Ax} = \mathbf{b}$  more “regular” or smooth. For the solution of ill-conditioned problems the ordinary methods give unacceptable solutions, so that regularization methods are necessary. Since many discussion of ill-conditioned matrices requires knowledge of the **SVD** of the matrix  $\mathbf{A}$ , our essay is mainly focused on **Truncated singular value decomposition** (TSVD) and **Tikhonov** regularization methods.



### 4.2.1 The Tikhonov Method

Tikhonov regularization method is the most commonly used method of regularization of ill-posed problem. The Tikhonov regularization method involves obtaining the exact or a least-squares solution of linear systems by minimizing the functional

$$\Phi(\mathbf{x}) = \|\mathbf{Ax} - \mathbf{b}\|_2^2, \quad (4.1)$$

subject to  $\|L\mathbf{x}\| \leq \beta$ , where  $L$  is a differential operator.

Standard algorithms often give solutions that decay very rapidly, with large positive and negative values. To stabilize the computation, we add to equation (4.1) a term that penalizes the large components and thereby reduces them. Thus instead of equation (4.1) we minimize the expression

$$\Phi_\alpha(\mathbf{x}) = \|\mathbf{Ax} - \mathbf{b}\|_2^2 + \alpha \|L\mathbf{x}\|_2^2, \quad (4.2)$$

where  $\alpha > 0$  is called a regularization parameter. we denote it as

$$\mathbf{x}_\alpha = \mathbf{arg\ min} \{ \|\mathbf{Ax} - \mathbf{b}\|_2^2 + \alpha \|L\mathbf{x}\|_2^2 \}. \quad (4.3)$$

We consider the case when  $L = I$ , the identity operator.

$$\Phi_\alpha(\mathbf{x}) = \|\mathbf{Ax} - \mathbf{b}\|_2^2 + \alpha \|\mathbf{x}\|_2^2, \quad (4.4)$$

Regularization can be understood as a balance between two requirements [8]:

1.  $\mathbf{x}$  should give a small residual  $\mathbf{Ax}_\alpha - \mathbf{b}$ .
2.  $\mathbf{x}$  should be small in 2-norm.

This means that the Regularization should balance between these two quantities  $\|\mathbf{Ax}_\alpha - \mathbf{b}\|_2$  and  $\|\mathbf{x}_\alpha\|_2$ . Thus regularization parameter  $\alpha > 0$  can be used to “tune” the balance.

We can rewrite equation 4.4 as

$$\begin{aligned} \Phi_\alpha(\mathbf{x}) &= (\mathbf{Ax} - \mathbf{b})^T (\mathbf{Ax} - \mathbf{b}) + \alpha \mathbf{x}^T \mathbf{x} \\ &= (\mathbf{Ax})^T (\mathbf{Ax}) - \mathbf{b}^T (\mathbf{Ax}) - (\mathbf{Ax})^T \mathbf{b} + \mathbf{b}^T \mathbf{b} + \alpha \mathbf{x}^T \mathbf{x}. \end{aligned}$$

The two terms  $\mathbf{b}^T (\mathbf{Ax})$  and  $(\mathbf{Ax})^T \mathbf{b}$  are equal, so

$$\Phi_\alpha(\mathbf{x}) = (\mathbf{Ax})^T (\mathbf{Ax}) - 2(\mathbf{Ax})^T \mathbf{b} + \alpha \mathbf{x}^T \mathbf{x} + \mathbf{b}^T \mathbf{b} \quad (4.5)$$

$$= (\mathbf{Ax})^T (\mathbf{Ax}) - 2\mathbf{x}^T \mathbf{A}^T \mathbf{b} + \alpha \mathbf{x}^T \mathbf{x} + \mathbf{b}^T \mathbf{b}. \quad (4.6)$$

The minimum is found at  $\frac{\partial \Phi_\alpha(\mathbf{x})}{\partial \mathbf{x}} \Big|_{\mathbf{x}=\mathbf{x}_\alpha} = 0$ . So

$$\frac{\partial \Phi_\alpha(\mathbf{x})}{\partial \mathbf{x}} \Big|_{\mathbf{x}=\mathbf{x}_\alpha} = 2\mathbf{A}^T \mathbf{A} \mathbf{x}_\alpha - 2\mathbf{A}^T \mathbf{b} + 2\alpha \mathbf{x}_\alpha = 0. \quad (4.7)$$

So

$$(\mathbf{A}^T \mathbf{A} + \alpha \mathbf{I}) \mathbf{x}_\alpha = \mathbf{A}^T \mathbf{b}, \quad (4.8)$$

Then

$$\mathbf{x}_\alpha = (\mathbf{A}^T \mathbf{A} + \alpha \mathbf{I})^{-1} \mathbf{A}^T \mathbf{b}. \quad (4.9)$$

Using the singular value decomposition (SVD) of  $\mathbf{A}$  in equation 4.9 becomes

$$\mathbf{x}_\alpha = \sum_{i=1}^n [(\mathbf{v}_i \sigma_i^T \mathbf{u}_i^T) (\mathbf{u}_i \sigma_i \mathbf{v}_i^T) + \alpha \mathbf{I}]^{-1} (\mathbf{v}_i \sigma_i^T \mathbf{u}_i^T) \mathbf{b}, \quad (4.10)$$

so

$$\mathbf{x}_\alpha = \sum_{i=1}^n (\sigma_i^2 + \alpha)^{-1} (\mathbf{v}_i \sigma_i \mathbf{u}_i^T) \mathbf{b}, \quad (4.11)$$

then

$$\mathbf{x}_\alpha = \sum_{i=1}^n \frac{\sigma_i}{\alpha + \sigma_i^2} (\mathbf{u}_i^T) \mathbf{b} \mathbf{v}_i. \quad (4.12)$$

The effect of the addition in equation (4.4) is to dampen the contribution of the term involving small singular values, so that instead of cutting them off altogether, we modify the method to reduce their impact. A small  $\alpha$  has very little effect on the components associated with large  $\sigma_i$ , since for  $\alpha \ll \sigma_i^2$

$$\frac{\sigma_i}{\alpha + \sigma_i^2} \cong \frac{1}{\sigma_i}. \quad (4.13)$$

On the other hand, if  $\sigma_i^2 \ll \alpha$ , then

$$\frac{\sigma_i}{\alpha + \sigma_i^2} \cong \frac{\sigma_i}{\alpha} \ll \frac{1}{\sigma_i}, \quad (4.14)$$

so that the magnification of the components associated with the small singular values is reduced. With a good choice of  $\alpha$ , one can then hope to get a relatively smooth solution that is still a reasonably good approximation to the true solution. This is called **SVD** with damping [6].

## 4.2.2 Truncated Singular Value Decomposition (TSVD)

**TSVD** is another commonly used method of regularization of Rank-deficient problems and Discrete ill-posed problems. The idea of truncating the **SVD** has been treated as a problem of determining the “**numerical rank**” of the matrix  $\mathbf{A}$  ( $\mathbf{A}$  is a noisy representation of a mathematically rank-deficient matrix). All of the computed singular values smaller than some threshold value ( $\alpha$ ) are treated as zeros which were corrupted by rounding errors into small non-zero quantities. Thus if  $\sigma_k$  is the smallest singular value greater than the truncation threshold ( $\sigma_k \geq \alpha$ ), then one replace the matrix  $\Sigma$  by a truncated matrix [10].

$$\Sigma_{tr} = \text{diag}(\sigma_1, \dots, \sigma_k, 0, \dots, 0) \quad (4.15)$$

whose inverse is given by

$$\Sigma_{tr}^+ = \text{diag}\left(\frac{1}{\sigma_1}, \dots, \frac{1}{\sigma_k}, 0, \dots, 0\right). \quad (4.16)$$

The value  $k$  is the “**numerical rank**” of  $\mathbf{A}$ . So we replace the matrix  $\mathbf{A}$  by its rank  $k$  approximation. Then  $x_{tr} = x_\alpha$  is the estimate of the solution to the truncated problem, given by

$$\mathbf{x}_{tr} = \mathbf{x}_\alpha = \sum_{i=1}^k \frac{\mathbf{u}_i^T \mathbf{b}}{\sigma_i} \mathbf{v}_i. \quad (4.17)$$

**Example 4.1.** Consider the linear system

$$\mathbf{H}_{12} \mathbf{x} = \mathbf{b}. \quad (4.18)$$

where  $\mathbf{H}_{12}$  is the  $12 \times 12$  Hilbert matrix, and

$$\mathbf{b} = \sum_{j=1}^n \mathbf{h}_{ij} \mathbf{x}_j. \quad (4.19)$$

This linear system has an exact solution  $\mathbf{x} = [1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1]^T$ . The computed solution is given by

$$\hat{\mathbf{x}} = \sum_{i=1}^n \frac{\mathbf{u}_i^T \mathbf{b}}{\sigma_i} \mathbf{v}_i. \quad (4.20)$$

Then

$$\hat{\mathbf{x}} = [1.000000004881659, 0.999999523277221, 1.000011453630246, 0.999883148383246, \\ 1.000618184166931, 0.998167348978474, 1.003001794695910, 0.997860552739088, \\ 0.999135295872923, 1.002809344038033, 0.998028938426422, 1.000484414235235]^T.$$

The maximum error is

$$\|\mathbf{x} - \hat{\mathbf{x}}\|_\infty = 0.00300179469591044. \quad (4.21)$$

The error is contaminated with small singular values. The singular values of  $H_{12}$  are

<i>numbers</i>	Singular values of $H_{12}$
1	1.79537205956200
2	$3.80275245955037 \times 10^{-1}$
3	$4.47385487521811 \times 10^{-2}$
4	$3.72231223789117 \times 10^{-3}$
5	$2.33089089021769 \times 10^{-4}$
6	$1.11633574832344 \times 10^{-5}$
7	$4.08237611034491 \times 10^{-7}$
8	$1.12286106662802 \times 10^{-8}$
9	$2.25196453231559 \times 10^{-10}$
10	$3.11134836868481 \times 10^{-12}$
11	$2.64910549271184 \times 10^{-14}$
12	$1.09849454516192 \times 10^{-16}$

However the Truncated Singular Value Decomposition (TSVD) method with  $\alpha = 10^{-8}$ , yields better solution as shown below. The maximum error ( $\| \mathbf{x} - \mathbf{x}^{tr} \|_{\infty}$ ) for  $k = 12, 11, 10, 9, 8$  and  $7$  are shown in the table below.

$k$	maximum error
12	$0.269705081406381 \times 10^{-1}$
11	$6.72224601423688 \times 10^{-3}$
10	$1.18729236775472 \times 10^{-5}$
9	$3.84094700600635 \times 10^{-6}$
8	$2.90489981684683 \times 10^{-5}$
7	$1.62703616404802 \times 10^{-4}$

From the table above we notice that the maximum error decreases from  $k = 11$  to  $9$  and after that starts increasing again. Thus we can choose the best regularized solution to be the one corresponding to  $k = 9$ .

Plugging these singular values into Equation (4.17) it is now easy to see that inclusion of the smallest singular values will cause the sum in Equation (4.17) to blow up yielding a large error in the computed solution. However if these small singular values are cut off we remove these large errors. But as we cut off more and more singular values there comes a point where the singular values is not a faithful representation of the coefficient matrix and hence we need to stop at some minimum value of  $k$ . This minimum value of  $k$  is equal to  $9$  for this example linear system.

Figure (4.1) is a plot of computed solutions obtained (by octave) using the TSVD method. The three plots in the figure represent solutions after truncating 1, 2, and 3 smallest singular values one after the other.

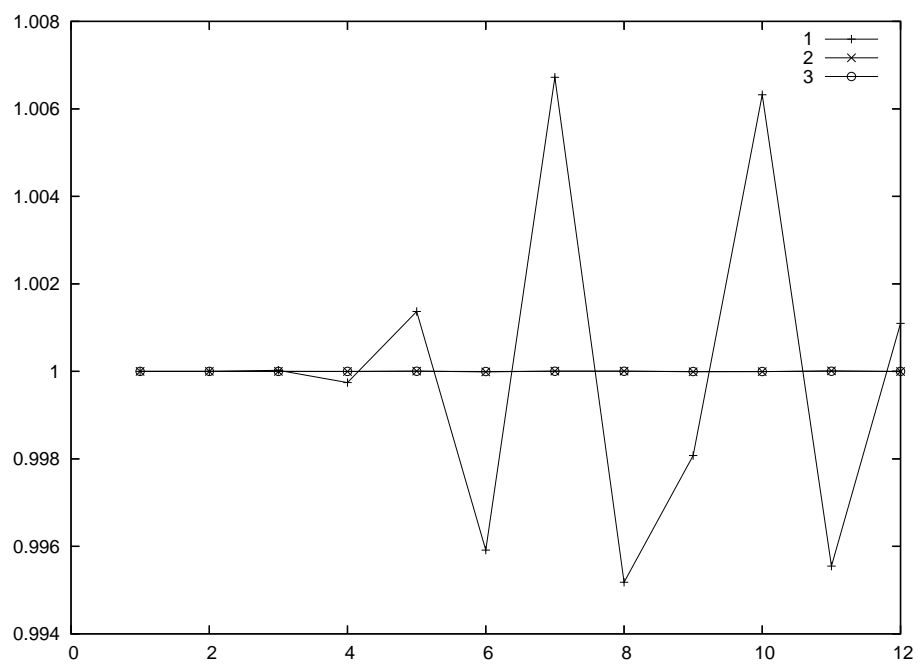


Figure 4.1: Plot of computed regularized solutions using the TSVD method. The smallest 1, 2, and 3 singular values were cut off one after the other.

## 5. Conclusion

Finally we would like to give a brief summary to this study. In this essay we have briefly introduced the reader to various aspects of ill-conditioned linear systems and their regularization methods. In chapter one, we introduced the linear systems, some basic concepts of linear algebra., and the conditioned number of a matrix. In chapter two, we defined ill-conditioned linear systems, For the purpose of illustration, we have included two (linear) examples for ill-conditioned linear systems, (e.g. the Hilbert and the Vandermonde matrices). Also in chapter three we defined the least-squares solutions to linear systems and singular value decomposition (SVD) as a numerical “tool”, which is suited for analysis of the ill-posed problems. Also we identified that the computed solution  $\hat{\mathbf{x}}$  is usually meaningless approximation to  $\mathbf{x}$  due to the error in the right-hand side  $\mathbf{b}$  and the severe ill-conditioning of the matrix  $\mathbf{A}$ . For this reason, in chapter four we reviewed the regularization methods to improve the accuracy of the solution of ill-posed problems. We focussed in two methods, e.g. truncated singular value decomposition (TSVD) and Tikhonov mehod. One hopes that this work could be extended further to deal with other types of regularization methods.

# Appendix A. Truncated Singular Value Decomposition (TSVD)

Here is the octave code for TSVD example and we insert also the results of it as follow.

```
1;

function x = TSVD(A, b, q)
[m, n] = size(A) ;
[U, S, V] = svd(A) ;
s = diag(S) ;

x = zeros(n, 1) ;
for i = 1:q

x = x + (b' * U(:, i))/s(i) * V(:, i) ;
if x<= 10**(-8)
x=0;
endif

endfor
endfunction

N = 12
H12 = hilb(N)
[U, S, V] = svd(H12)

% The exact solution denoted by x.
x=ones(N,1)
b = H12 * x

%The computed solution denoted by xc.
xc = H12 \ b
xt = [];
for k = N:-1:1

%The TSVD solution denoted by xt.
xt = [xt TSVD(H12,b,k)];
endfor
xt
```

---

xt =

Columns 1 through 4:

1.000000035410423	1.000000010406728	1.000000000128169	0.99999999993725
0.999995761625934	0.999998969286971	0.99999989054960	0.99999999145378
1.000126836288086	1.000025022058407	1.000000216120084	1.000000032930959
0.998345365750319	0.999742665937545	0.999998289259415	0.999999654919425
1.011673175050563	1.001370185851539	1.000006491267828	1.000001539687882
0.950432454189308	0.995915975213462	0.999988127076322	0.999996861380750
1.133947503416035	1.006722246014237	1.000007150407346	1.000002059853554
0.764130473439041	0.995179805399686	1.000006868068255	1.000002031900902
1.269705081406381	0.998077220444046	0.999991908851701	0.999997470810440
0.806915643063877	1.006322162812165	0.999993629378929	0.999997927375082
1.078626086420750	0.995549233050377	1.000011456491165	1.000003840947006
0.986101573342976	1.001096511297768	0.999995873687303	0.999998580958273

Columns 5 through 8:

0.999999996461529	1.000000151890956	0.999995319603844	1.000106069091871
1.000000194093540	0.999994106412219	1.000122335866062	0.998253353494705
0.999997482649582	1.000052197954546	0.999313675155133	1.005457156224012
1.000012680488899	0.999837296383595	1.001085437561335	0.998057332941534
0.999972852028113	1.000150789068384	1.000192395006243	0.995467217244331
1.000016557348865	1.000096569609180	0.999253083336022	0.997356157419837
1.000018375576711	0.999906488482612	0.999226991564590	1.000712066947196
0.999987654689213	0.999851527117142	0.999871521112027	1.003392636567184
0.999977437565833	0.999969007644730	1.000610025132260	1.004301323177792
1.000003199234471	1.000132143083218	1.000919711186439	1.003036848655090
1.000029048998168	1.000154263950493	1.000439206328687	0.999587308252251
0.999984517595727	0.999855303160945	0.998964942748598	0.994134056971875

Columns 9 through 12:

0.998220027713585	1.021853329482553	0.811827016949171	1.968193407412721
1.016699883351630	0.901826902329911	1.267158129469560	1.211355960941289
0.978446218876212	0.999298041036143	1.258596887316880	0.907797344126588
0.984100114089535	1.053331189093677	1.182716622765847	0.735766433013591
0.999596918273331	1.070010484453778	1.099090617957318	0.622574166277303
1.011852504960719	1.064191004336744	1.020782412083546	0.541515632961854
1.017616517848767	1.045650381835735	0.950479056603555	0.480189690911242
1.017015122462759	1.020207748562225	0.888094362825359	0.431961081053743
1.011136208465709	0.991290936630181	0.832808873429604	0.392922531207807
1.001189570053022	0.960928856644816	0.783684307656640	0.360607676857949
0.988242349596960	0.930325825448542	0.739850628891383	0.333375279209793
0.973160612317457	0.900192341869936	0.700553560084826	0.310086830520352



# Acknowledgements

I would like to express my deepest thanks and gratitude to my supervisor professor Francis Benyah (Department of Math's and Appl. Math's, The University of the Western Cape ) who has directed me through this work with an admirable magnitude of skill , care and scientific patience , whose constructive comments , critical remarks and sound advice were simply indispensable for the completion of this work.

Many thanks to Mr. Henri Amuasi, Mr. Eihab basheir, Mr. Ambrose Chongo and Mr. Jan Groenewald for their help and useful discussions towards getting this work done.

\*\*\*\*\*

I dedicate this work to my family whose love, care, and encouragement helped me a lot.

إِلَى أُمِّي ... مَنْ أَرْضَعْتَنِي الْحُبَّ وَالْحَنَانَ

إِلَى أَبِي ... مَنْ عَلَّمَنِي مَعْنَى الْعِظَاءِ وَالْوَفَاءِ

إِلَى إِخْوَتِي ... رُفَقَاءَ دَرَبِي فِي الْكِفَاحِ

إِلَى كُلِّ مَنْ عَلَّمَنِي حِرْفًا وَزَوَّدَنِي بِنُورِ الْعِلْمِ

إِلَيْكُمْ جَمِيعًا ... أُهْدِي ثَمَرَةَ جُهْدِي هَذَا الْمُتَوَاضِعِ ،،

شيماء ،،

# Bibliography

- [1] ACM/ESE. <http://www.gps.caltech.edu/~tapio/acm118/handouts/svd.pdf>. Singular Value Decomposition.
- [2] Francis Benyah. *Ill-Conditioning and Regularization in the Solution of Linear Systems*. Seminar, African Institute for Mathematical Science, 2004.
- [3] Biswa Nath Data. *Numerical Linear Algebra and Applications*. Brooks/Cole Publishing Company, California, 1995.
- [4] Per Christian Hansen. *Rank-Deficient and Discrete Ill-Posed Problems*. Siam, Society for Industrial and Applied Mathematics, Philadelphia, 1998.
- [5] Autar K.Kaw. *Introduction to Matrix Algebra*. University of South Florida, USA, 2002.
- [6] Peter Linz and Richard L. C. Wang. *Exploring Numerical Methods*. Jones and Bartlett Publishers, California, 2003.
- [7] Francis J. Narcowich. <http://www.math.tamu.edu/~fnarc/psfiles/rank2005.pdf>. The Rank of a Matrix.
- [8] Tikhonov. <http://matriisi.ee.tut.fi/courses/mat-52500/tikhonov.pdf>. Tikhonov REgularization.
- [9] Website. <http://www.cs.huji.ac.il/~csip/condition.pdf>. Numerical Stability - The Condition Number of a Matrix.
- [10] Bert W.Rust. <http://math.nist.gov/~brust/pubs/truncsvd/ms-truncsvd.ps>. Truncating The Singular Value Decompositin For ill-Posed Problems.