# Math 5620 - Introduction to Numerical Analysis - Class Notes

Fernando Guevara Vasquez

Version 1990. Date: January 17, 2012.

# Contents

1. Disclaimer	4
Chapter 1. Iterative methods for solving linear systems	5
1. Preliminaries	5
2. Neumann series based methods	7
3. Conjugate gradient and other Krylov subspace methods	15

3

### 1. Disclaimer

These notes are work in progress and may be incomplete and contain errors. These notes are based on the following sources:

- Burden and Faires, *Numerical Analysis*, ninth ed.
- Kincaid and Cheney, *Numerical Analysis: Mathematics of scientific computing*
- Stoer and Burlisch, *Introduction to Numerical Analysis*, Springer 1992
- Golub and Van Loan, Matrix Computations, John Hopkins 1996

#### CHAPTER 1

## Iterative methods for solving linear systems

#### 1. Preliminaries

**1.1. Convergence in**  $\mathbb{R}^n$ . Let  $\|\cdot\|$  be a **norm** defined on  $\mathbb{R}^n$ , i.e. it satisfies the properties:

i.  $\forall \mathbf{x} \in \mathbb{R}^n \| \mathbf{x} \| \ge 0$  (non-negativity) ii.  $\|\mathbf{x}\| = 0 \Leftrightarrow \mathbf{x} = \mathbf{0}$  (definiteness)

iii.  $\forall \lambda \in \mathbb{R}, \mathbf{x} \in \mathbb{R}^n ||\lambda \mathbf{x}|| = |\lambda| ||\mathbf{x}||$  (multiplication by a scalar)

iv.  $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n \|\mathbf{x} + \mathbf{y}\| \le \|\mathbf{x}\| + \|\mathbf{y}\|$  (triangle inequality)

Some important examples of norms in  $\mathbb{R}^n$  are:

- (1)  $\|\mathbf{x}\|_2 = \left(\sum_{i=1}^n |x_i|^2\right)^{1/2}$  (Euclidean or  $\ell_2$  norm) (2)  $\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|$  ( $\ell_1$  norm)
- (3)  $\|\mathbf{x}\|_{\infty} = \max_{i=1,\dots,n} |x_i| \ (\ell_{\infty} \text{ or max norm})$
- (4)  $\|\mathbf{x}\|_p = \left(\sum_{i=1}^n |x_i|^p\right)^{1/p}$  (for  $p \ge 1, \ell_p$  norm)

A sequence  $\{\mathbf{v}^{(k)}\}_{k=1}^{\infty}$  in  $\mathbb{R}^n$  converges to  $\mathbf{v} \in \mathbb{R}^n$  if and only if

(1) 
$$\lim_{k \to \infty} \|\mathbf{v}^{(k)} - \mathbf{v}\|.$$

Since  $\mathbb{R}^n$  is a finite dimensional vector space the notion of convergence is independent of the norm. This follows from the fact that in a finite dimensional vector space all norms are **equivalent**. Two norms  $\|\cdot\|$ and  $\||\cdot|\|$  are equivalent if there are constants  $\alpha, \beta > 0$  such that

(2) 
$$\forall \mathbf{x} \, \alpha \| \| \mathbf{x} \| \le \beta \| \| \mathbf{x} \| \le \beta \| \| \mathbf{x} \|.$$

Another property of  $\mathbb{R}^n$  is that it is **complete**, meaning that all Cauchy sequences converge in  $\mathbb{R}^n$ . A sequence  $\{\mathbf{v}^{(k)}\}_{k=1}^{\infty}$  is said to be a **Cauchy** sequence when

(3) 
$$\forall \epsilon > 0 \exists N \forall i, j \ge N \| \mathbf{v}^{(i)} - \mathbf{v}^{(j)} \| < \epsilon.$$

In English: A Cauchy sequence is a sequence for which any two iterates can be made as close as we want provided that we are far enough in the sequence.

**1.2. Induced matrix norms.** Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  be a matrix (this notation means that the matrix has *n* rows and *n* columns) and let  $\| \cdots \|$  be a norm on  $\mathbb{R}^n$ . The operator or induced matrix norm of  $\mathbf{A}$  is defined by:

(4) 
$$\|\mathbf{A}\| = \sup_{\mathbf{x} \in \mathbb{R}^n, \mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|} = \sup_{\|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\|.$$

The induced matrix norm measures what is the maximum dilation of the image of a vector through the matrix **A**. It is a good exercise to show that it satisfies the axiom of a norm.

Some important examples of induced matrix norms:

- (1)  $\|\mathbf{A}\|_1 = \max_{j=1,\dots,n} \|\mathbf{c}_j\|_1$ , where  $\mathbf{c}_j$  is the *j*-th column of **A**.
- (2)  $\|\mathbf{A}\|_{\infty} = \max_{i=1,\dots,n} \|\mathbf{r}_i\|_1$ , where  $\mathbf{r}_i$  is the *i*-th row of **A**.
- (3)  $\|\mathbf{A}\|_2 =$  square root of largest eigenvalue of  $\mathbf{A}^T \mathbf{A}$ . When  $\mathbf{A}$  is symmetric we have  $\|\mathbf{A}\|_2 = |\lambda_{max}|$ , with  $\lambda_{max}$  being the largest eigenvalue of  $\mathbf{A}$  in magnitude.

A matrix norm that is not an induced matrix norm is the Frobenius norm:

(5) 
$$\|\mathbf{A}\|_F = \left(\sum_{i,j=1}^n |a_{ij}|^2\right)^{1/2}.$$

Some properties of induced matrix norms:

(1) 
$$\|\mathbf{A}\mathbf{x}\| \le \|\mathbf{A}\| \|\mathbf{x}\|.$$

(2) 
$$\|\mathbf{A}\mathbf{B}\| \le \|\mathbf{A}\| \|\mathbf{B}\|$$

(3) 
$$\|\mathbf{A}^k\| \le \|\mathbf{A}\|^k$$
.

**1.3. Eigenvalues.** The **eigenvalues** of a  $n \times n$  matrix **A** are the roots of the characteristic polynomial

(6) 
$$p(\lambda) = \det(\lambda \mathbf{I} - \mathbf{A})$$

This is a polynomial of degree *n*, so it has at most *n* complex roots. An **eigenvector v** associated with an eigenvalue  $\lambda$  is a nonzero vector such that  $A\mathbf{v} = \lambda \mathbf{v}$ . Sometimes it is convenient to refer to an eigenvalue  $\lambda$  and corresponding eigenvector  $\mathbf{v} \neq \mathbf{0}$  as an **eigenpair** of **A**.

The **spectral radius**  $\rho(\mathbf{A})$  is defined as the magnitude of the largest eigenvalue in magnitude of a matrix  $\mathbf{A}$  i.e.

(7) 
$$\rho(\mathbf{A}) = \max\{|\lambda| \mid \det(\lambda \mathbf{I} - \mathbf{A}) = 0\}.$$

The spectral radius is the radius of the smallest circle in  $\mathbb{C}$  containing all eigenvalues of **A**.

Two matrices **A** and **B** are said to be similar if there is an invertible matrix **X** such that

(8)	AX = XB.	
Rev: 1990, January 17, 2012	6	na.tex

Similar matrices have the same characteristic polynomial and thus the same eigenvalues. This follows from the properties of the determinant since

(9)

 $det(\lambda \mathbf{I} - \mathbf{A}) = det(\mathbf{X}^{-1}) det(\lambda \mathbf{I} - \mathbf{A}) det(\mathbf{X}) = det(\lambda \mathbf{I} - \mathbf{X}^{-1}\mathbf{A}\mathbf{X}) = det(\lambda \mathbf{I} - \mathbf{B}).$ 

#### 2. Neumann series based methods

**2.1. Neumann series.** We start by a fundamental theorem that is the theoretical basis for several methods for solving the linear system Ax = b. This theorem can be seen as a generalization to matrices of the geometric series identity  $(1 - x)^{-1} = 1 + x + x^2 + ...$  for |x| < 1.

THEOREM 1. Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  be such that  $\|\mathbf{A}\| < 1$  for some induced matrix norm  $\|\cdot\|$ , then:

*i.* 
$$\mathbf{I} - \mathbf{A}$$
 *is invertible*  
*ii.*  $(\mathbf{I} - \mathbf{A})^{-1} = \mathbf{I} + \mathbf{A} + \mathbf{A}^2 + \dots = \sum_{k=0}^{\infty} \mathbf{A}^k$ .

PROOF. Assume for contradiction that  $\mathbf{I} - \mathbf{A}$  is singular. This means there is a  $\mathbf{x} \neq \mathbf{0}$  such that  $(\mathbf{I} - \mathbf{A})\mathbf{x} = \mathbf{0}$ . Taking  $\mathbf{x}$  such that  $\|\mathbf{x}\| = 1$ , we have

(10) 
$$1 = \|\mathbf{x}\| = \|\mathbf{A}\mathbf{x}\| \le \|\mathbf{A}\| \|\mathbf{x}\| = \|\mathbf{A}\|,$$

which contradicts the hypothesis  $||\mathbf{A}|| < 1$ .

We now need to show convergence to  $(\mathbf{I} - \mathbf{A})^{-1}$  of the partial series

(11) 
$$\sum_{k=0}^{m} \mathbf{A}^{k}.$$

Observe that:

(12) 
$$(\mathbf{I} - \mathbf{A}) \sum_{k=0}^{m} \mathbf{A}^{k} = \sum_{k=0}^{m} \mathbf{A}^{k} - \mathbf{A}^{k+1} = \mathbf{A}^{0} - \mathbf{A}^{m+1}.$$

Therefore

(13) 
$$\|(\mathbf{I} - \mathbf{A}) \sum_{k=0}^{m} \mathbf{A}^{k} - \mathbf{I}\| = \|\mathbf{A}^{m+1}\| \le \|\mathbf{A}\|^{m+1} \to 0 \text{ as } m \to \infty.$$

Here is an application of this theorem to estimate the norm

(14) 
$$\|(\mathbf{I} - \mathbf{A})^{-1}\| \le \sum_{k=0}^{\infty} \|\mathbf{A}^k\| \le \sum_{k=0}^{\infty} \|\mathbf{A}\|^k = \frac{1}{1 - \|\mathbf{A}\|}.$$

Here is a generalization of the Neumann series theorem.

THEOREM 2. If **A** and **B** are  $n \times n$  matrices such that  $||\mathbf{I} - \mathbf{AB}|| < 1$  for some induced matrix norm then

*i.* **A** and **B** are invertible.

*ii. The inverses are:* 

(15)  
$$\mathbf{A}^{-1} = \mathbf{B} \sum_{k=0}^{\infty} (\mathbf{I} - \mathbf{A}\mathbf{B})^{k},$$
$$\mathbf{B}^{-1} = \left[\sum_{k=0}^{\infty} (\mathbf{I} - \mathbf{A}\mathbf{B})^{k}\right] \mathbf{A}.$$

PROOF. By using the Neumann series theorem, AB = I - (I - AB) is invertible and

(16) 
$$(\mathbf{AB})^{-1} = \sum_{k=0}^{\infty} (\mathbf{I} - \mathbf{AB})^k.$$

Therefore

(17)

$$\mathbf{A}^{-1} = \mathbf{B}(\mathbf{A}\mathbf{B})^{-1} = \mathbf{B}\sum_{k=0}^{\infty} (\mathbf{I} - \mathbf{A}\mathbf{B})^k,$$
$$\mathbf{B}^{-1} = (\mathbf{A}\mathbf{B})^{-1}\mathbf{A} = \left[\sum_{k=0}^{\infty} (\mathbf{I} - \mathbf{A}\mathbf{B})^k\right]\mathbf{A}.$$

**2.2. Iterative refinement.** Let **A** be an invertible matrix. The iterative refinement method is a method for generating successively better approximations to the solution of the linear system Ax = b. Assume we have an invertible matrix **B** such that  $x = Bb \approx A^{-1}b$  and applying **B** is much cheaper than applying solving a system with the matrix **A** (we shall see how good the approximation needs to be later). This approximate inverse **B** may come for example from an incomplete LU factorization or from running a few steps of an iterative method to solve Ax = b. Can we use successively refine the approximations given by this method? The idea is to look at the iteration

(18) 
$$\mathbf{x}^{(0)} = \mathbf{B}\mathbf{b}$$
  
 $\mathbf{x}^{(k)} = \mathbf{x}^{(k-1)} + \mathbf{B}(\mathbf{b} - \mathbf{A}\mathbf{x}^{(k-1)}).$ 

If this iteration converges, then the limit must satisfy

$$\mathbf{x} = \mathbf{x} + \mathbf{B}(\mathbf{b} - \mathbf{A}\mathbf{x})$$

i.e. if the method converges it converges to a solution of Ax = b.

THEOREM 3. The iterative refinement method (18) generates iterates of the form

(20) 
$$\mathbf{x}^{(m)} = \mathbf{B} \sum_{k=0}^{m} (\mathbf{I} - \mathbf{A}\mathbf{B})^{k} \mathbf{b}, \ m \ge 0.$$

Rev: 1990, January 17, 2012

na.tex

 $\square$ 

Thus by the generalized Neumann series theorem, the method converges to the solution of  $\mathbf{A}\mathbf{x} = \mathbf{b}$  provided  $\|\mathbf{I} - \mathbf{A}\mathbf{B}\| < 1$  for some induced matrix norm (*i.e.* provided **B** is sufficiently close to being an inverse of **A**).

PROOF. We show the form of the iterates in the iterative refinement method by induction on *m*. First the case m = 0 is trivial since  $\mathbf{x}^{(0)} = \mathbf{Bb}$ . Assuming the m-th case holds:

(21)  

$$\mathbf{x}^{(m+1)} = \mathbf{x}^{(m)} + \mathbf{B}(\mathbf{b} - \mathbf{A}\mathbf{x}^{(m)})$$

$$= \mathbf{B} \sum_{k=0}^{m} (\mathbf{I} - \mathbf{A}\mathbf{B})^{k} \mathbf{b} + \mathbf{B}\mathbf{b} - \mathbf{A}\mathbf{B} \sum_{k=0}^{m} (\mathbf{I} - \mathbf{A}\mathbf{B})^{k} \mathbf{b}$$

$$= \mathbf{B} \left[ \mathbf{b} + (\mathbf{I} - \mathbf{A}\mathbf{B}) \sum_{k=0}^{m} (\mathbf{I} - \mathbf{A}\mathbf{B})^{k} \mathbf{b} \right]$$

$$= \mathbf{B} \sum_{k=0}^{m+1} (\mathbf{I} - \mathbf{A}\mathbf{B})^{k} \mathbf{b}.$$

2.3. Matrix splitting methods. In order to solve the linear system Ax = b, we introduce a **splitting matrix** Q and use it to define the iteration:

(22) 
$$\mathbf{x}^{(0)} = \text{given},$$
$$\mathbf{Q}\mathbf{x}^{(k)} = (\mathbf{Q} - \mathbf{A})\mathbf{x}^{(k-1)} + \mathbf{b}, \ k \ge 1$$

Since we need to solve for  $\mathbf{x}^{(k)}$  the matrix  $\mathbf{Q}$  needs to be invertible and solving systems with **Q** needs to be a cheap operation (for example **Q** could be diagonal or triangular). If the iteration converges, the limit  $\mathbf{x}$ must satisfy

$$\mathbf{Q}\mathbf{x} = (\mathbf{Q} - \mathbf{A})\mathbf{x} + \mathbf{b}.$$

In other words: if the iteration (22) converges, the limit solves the linear system Ax = b. The next theorem gives a sufficient condition for convergence.

THEOREM 4. If  $\|\mathbf{I} - \mathbf{Q}^{-1}\mathbf{A}\| < 1$  for some matrix induced norm, then the iterates (22) converge to the solution to Ax = b regardless of the initial *iterate*  $\mathbf{x}^{(0)}$ .

**PROOF.** Subtracting the equations

na.tex

we obtain a relation between the error at step k and the error at step k-1:

(25) 
$$\mathbf{x}^{(k)} - \mathbf{x} = (\mathbf{I} - \mathbf{Q}^{-1}\mathbf{A})(\mathbf{x}^{(k-1)} - \mathbf{x}).$$

Taking norms we get:

(26) 
$$\|\mathbf{x}^{(k)} - \mathbf{x}\| \le \|\mathbf{I} - \mathbf{Q}^{-1}\mathbf{A}\| \|\mathbf{x}^{(k-1)} - \mathbf{x}\| \le \dots \le \|\mathbf{I} - \mathbf{Q}^{-1}\mathbf{A}\|^k \|\mathbf{x}^{(0)} - \mathbf{x}\|.$$
  
Thus if  $\|\mathbf{I} - \mathbf{Q}^{-1}\mathbf{A}\| < 1$  we have  $\mathbf{x}^{(k)} \to \mathbf{x}$  as  $k \to \infty$ .

As a stopping criterion we can look at the difference between two consecutive iterates  $\|\mathbf{x}^{(k-1)} - \mathbf{x}^{(k)}\|$ . To see this, let  $\delta = \|\mathbf{I} - \mathbf{Q}^{-1}\mathbf{A}\| < 1$ . Then by the proof of the previous theorem we must have

(27) 
$$\|\mathbf{x}^{(k)} - \mathbf{x}\| \le \delta \|\mathbf{x}^{(k-1)} - \mathbf{x}\| \le \delta (\|\mathbf{x}^{(k-1)} - \mathbf{x}^{(k)}\| + \|\mathbf{x}^{(k)} - \mathbf{x}\|).$$

Hence by isolating  $\|\mathbf{x}^{(k)} - \mathbf{x}\|$  we can bound the error by the difference between two consecutive iterates:

(28) 
$$\|\mathbf{x}^{(k)} - \mathbf{x}\| \leq \frac{\delta}{1-\delta} \|\mathbf{x}^{(k-1)} - \mathbf{x}^{(k)}\|.$$

Of course there can be issues if  $\delta$  is very close to 1.

We now look at examples of matrix splitting methods. Let us first introduce a standard notation for partitioning the matrix **A** into its diagonal elements **D**, strictly lower triangular part  $-\mathbf{E}$  and strictly upper triangular part  $-\mathbf{F}$  so that

$$\mathbf{A} = \mathbf{D} - \mathbf{E} - \mathbf{F}.$$

2.3.1. *Richardson method*. Here the splitting matrix is  $\mathbf{Q} = \mathbf{I}$ , so the iteration is

(30) 
$$\mathbf{x}^{(k)} = (\mathbf{I} - \mathbf{A})\mathbf{x}^{(k-1)} + \mathbf{b} = \mathbf{x}^{(k-1)} + \mathbf{r}^{(k-1)},$$

where the residual vector is  $\mathbf{r} = \mathbf{b} - \mathbf{A}\mathbf{x}$ . Using the theorem on convergence of splitting methods, we can expect convergence when  $\|\mathbf{I} - \mathbf{A}\| < 1$  in some matrix induced norm, or in other words if the matrix  $\mathbf{A}$  is sufficiently close to the identity.

2.3.2. *Jacobi method*. Here the splitting matrix is  $\mathbf{Q} = \mathbf{D}$ , so the iteration is

(31) 
$$\mathbf{D}\mathbf{x}^{(k)} = (\mathbf{E} + \mathbf{F})\mathbf{x}^{(k-1)} + \mathbf{B}.$$

We can expect convergence when  $\|\mathbf{I} - \mathbf{D}^{-1}\mathbf{A}\| < 1$  for some matrix induced norm. If we choose the  $\|\cdot\|_{\infty}$  norm, we can get an easy to check sufficient condition for convergence of the Jacobi method. Indeed:

(32) 
$$\mathbf{D}^{-1}\mathbf{A} = \begin{bmatrix} 1 & a_{12}/a_{11} & a_{13}/a_{11} & \dots & a_{1n}/a_{11} \\ a_{21}/a_{22} & 1 & a_{23}/a_{22} & \dots & a_{2n}/a_{22} \\ \vdots & & & \vdots \\ a_{n1}/a_{nn} & a_{n2}/a_{nn} & \dots & a_{nn-1}/a_{nn} & 1 \end{bmatrix}.$$

*Rev: 1990, January 17, 2012* 

Hence

(33) 
$$\|\mathbf{I} - \mathbf{Q}^{-1}\mathbf{A}\|_{\infty} = \max_{i=1,\dots,n} \sum_{j=1, j\neq i}^{n} \frac{|a_{ij}|}{|a_{ii}|}.$$

A matrix satisfying the condition  $\|\mathbf{I} - \mathbf{Q}^{-1}\mathbf{A}\|_{\infty} < 1$  is said to be **diagonally dominant** since it is equivalent to saying that in every row the diagonal element is larger than the sum of all the other ones (in magnitude)

(34) 
$$|a_{ii}| > \sum_{j=1, j \neq i}^{n} |a_{ij}|, \text{ for } i = 1, \dots, n.$$

This result can be summarized as a theorem:

THEOREM 5. If **A** is diagonally dominant, then the Jacobi method converges regardless of the initial iterate to the solution of Ax = b.

We emphasize that this is only a sufficient condition for convergence. The Jacobi method may converge for matrices that are not diagonally dominant.

The pseudocode for the Jacobi algorithm is

for k = 1,2,...
x = x + (b-A\*x)./diag(diag(A))
and

end

Each iteration involves a multiplication by **A** and division by the diagonal elements of **A**.

2.3.3. *Gauss-Seidel method*. Here  $\mathbf{Q} = \mathbf{D} - \mathbf{E}$ , i.e. the lower triangular part of **A**. The iterates are:

(35) 
$$(\mathbf{D} - \mathbf{E})\mathbf{x}^{(k)} = \mathbf{F}\mathbf{x}^{(k-1)} + \mathbf{b}$$

Each iteration involves multiplication by the strictly upper triangular part of **A** and solving a lower triangular system (forward substitution). Here is an easy to check sufficient condition for convergence.

THEOREM 6. If **A** is diagonally dominant, then the Gauss-Seidel method converges regardless of the initial iterate to the solution of Ax = b.

The proof of this theorem is deferred to later, when we will find a necessary and sufficient condition for convergence of matrix splitting methods. Gauss-Seidel usually outperforms the Jacobi method.

2.3.4. Successive Over Relaxation (SOR) method. Here  $\mathbf{Q} = \omega^{-1}(\mathbf{D} - \omega \mathbf{E})$  and  $\omega$  is a parameter that needs to be chosen ahead of time. For symmetric positive definite matrices choosing  $\omega \in (0, 2)$  gives convergence. The iterates are:

(36) 
$$(\mathbf{D} - \omega \mathbf{E})\mathbf{x}^{(k)} = \omega(\mathbf{F}\mathbf{x}^{(k-1)} + \mathbf{b}) + (1 - \omega)\mathbf{D}\mathbf{x}^{(k-1)}.$$
  
*na.tex* 11 *Rev:* 1990, *January* 17, 2012

The cost of an iteration is similar to that of Gauss-Seidel and with a good choice of the relaxation parameter  $\omega$ , SOR can outperform Gauss-Seidel.

**2.4. Convergence of iterative methods.** The goal of this section is to give a sufficient and necessary condition for the convergence of the iteration

$$\mathbf{x}^{(k)} = \mathbf{G}\mathbf{x}^{(k-1)} + \mathbf{c}.$$

The matrix splitting methods with iteration  $\mathbf{Q}\mathbf{x}^{(k)} = (\mathbf{Q} - \mathbf{A})\mathbf{x}^{(k-1)} + \mathbf{b}$  from previous section can be written in the form (37) with  $\mathbf{G} = (\mathbf{I} - \mathbf{Q}^{-1}\mathbf{A})$  and  $\mathbf{c} = \mathbf{Q}^{-1}\mathbf{b}$ .

If the iteration (37) converges its limit satisfies

$$\mathbf{x} = \mathbf{G}\mathbf{x} + \mathbf{c},$$

that is  $\mathbf{x} = (\mathbf{I} - \mathbf{G})^{-1}\mathbf{c}$ , assuming the matrix  $\mathbf{I} - \mathbf{G}$  is invertible. We will show the following theorem.

THEOREM 7. The iteration  $\mathbf{x}^{(k)} = \mathbf{G}\mathbf{x}^{(k-1)} + \mathbf{c}$  converges to  $(\mathbf{I} - \mathbf{G})^{-1}\mathbf{c}$  if and only if  $\rho(\mathbf{G}) < 1$ .

To prove theorem 7 we need the following result.

THEOREM 8. The spectral radius satisfies:

$$\rho(\mathbf{A}) = \inf_{\|\cdot\|} \|\mathbf{A}\|,$$

where the inf is taken over all induced matrix norms.

This theorem means that the smallest possible induced matrix norm is the 2-norm, if the matrix **A** is symmetric. The proof of this theorem is deferred to the end of this section. Let us first prove theorem 7.

PROOF OF THEOREM 7. We first show that  $\rho(\mathbf{G}) < 1$  is sufficient for convergence. Indeed if  $\rho(\mathbf{G}) < 1$ , then there is an induced matrix norm  $\|\cdot\|$  for which  $\|\mathbf{G}\| < 1$ . The iterates (37) are:

$$\mathbf{x}^{(1)} = \mathbf{G}\mathbf{x}^{(0)} + \mathbf{c}$$
  

$$\mathbf{x}^{(2)} = \mathbf{G}^{2}\mathbf{x}^{(0)} + \mathbf{G}\mathbf{c} + \mathbf{c}$$
  

$$\mathbf{x}^{(3)} = \mathbf{G}^{3}\mathbf{x}^{(0)} + \mathbf{G}^{2}\mathbf{c} + \mathbf{G}\mathbf{c} + \mathbf{c}$$

(40)

$$\mathbf{x}^{(k)} = \mathbf{G}^k \mathbf{x}^{(0)} + \sum_{j=0}^{k-1} \mathbf{G}^j \mathbf{c}.$$

The term involving the initial guess goes to zero as  $k \rightarrow \infty$  because

(41) 
$$\|\mathbf{G}^{k}\mathbf{x}^{(0)}\| \le \|\mathbf{G}\|^{k}\|\mathbf{x}^{(0)}\|$$

÷

*Rev: 1990, January 17, 2012* 

Thus by the Neumann series theorem:

(42) 
$$\sum_{j=0}^{k-1} \mathbf{G}^j \mathbf{c} \to (\mathbf{I} - \mathbf{G})^{-1} \mathbf{c} \text{ as } k \to \infty.$$

Now we need to show that  $\rho(\mathbf{G}) < 1$  is necessary for convergence. Assume  $\rho(\mathbf{G}) \ge 1$  and let  $\lambda, \mathbf{u}$  be an eigenpair of  $\mathbf{G}$  for which  $|\lambda| \ge 1$ . Taking  $\mathbf{x}^{(0)} = \mathbf{0}$  and  $\mathbf{c} = \mathbf{u}$  we get:

(43) 
$$\mathbf{x}^{(k)} = \sum_{j=0}^{k-1} \mathbf{G}^j \mathbf{u} = \sum_{j=0}^{k-1} \lambda^j \mathbf{u} = \begin{cases} k\mathbf{u} & \text{if } \lambda = 1\\ \frac{1-\lambda^k}{1-\lambda}\mathbf{u} & \text{if } \lambda \neq 1. \end{cases}$$

This is an example of an iteration of the form (37) that does not converge when  $\rho(\mathbf{G}) \ge 1$ .

Theorem 8 applied to splitting matrix methods gives:

COROLLARY 1. The iteration  $\mathbf{Q}\mathbf{x}^{(k)} = (\mathbf{Q} - \mathbf{A})\mathbf{x}^{(k-1)} + \mathbf{b}$  converges to  $\mathbf{A}\mathbf{x} = \mathbf{b}$  for any initial guess  $\mathbf{x}^{(0)}$  if and only if  $\rho(\mathbf{I} - \mathbf{Q}^{-1}\mathbf{A}) < 1$ .

In order to show theorem 8 we need the following result:

THEOREM 9. Let **A** be a  $n \times n$  matrix. There is a similarity transformation **X** such that

$$(44) AX = XB$$

where **B** is an upper triangular matrix with off-diagonal components that can be made arbitrarily small.

PROOF. By the Schur factorization any matrix  $\mathbf{A}$  is similar through an unitary transformation  $\mathbf{Q}$  to an upper triangular matrix  $\mathbf{T}$ 

(45) 
$$\mathbf{A} = \mathbf{Q}\mathbf{T}\mathbf{Q}^{T}, \text{ with } \mathbf{Q}^{T}\mathbf{Q} = \mathbf{I}.$$

Let **D** = diag( $\epsilon, \epsilon^2, \dots, \epsilon^n$ ). Then

(46) 
$$(\mathbf{D}^{-1}\mathbf{T}\mathbf{D})_{ij} = t_{ij}\epsilon^{j-i}.$$

The elements below the diagonal (j < i) are zero. Those above the diagonal (j > i) satisfy

$$(47) |t_{ij}\epsilon^{j-i}| \le \epsilon |t_{ij}|$$

With  $\mathbf{X} = \mathbf{Q}\mathbf{D}$ , the matrix  $\mathbf{A}$  is similar to  $\mathbf{B} = \mathbf{D}^{-1}\mathbf{T}\mathbf{D}$ , and  $\mathbf{B}$  is upper triangular with off-diagonal elements that can be made arbitrarily small.

PROOF OF THEOREM 8. We start by proving that  $\rho(\mathbf{A}) \leq \inf_{\|\cdot\|} \|\mathbf{A}\|$ . Pick a vector norm  $\|\cdot\|$  and let  $\lambda, \mathbf{x}$  be an eigenpair of  $\mathbf{A}$  with  $\|\mathbf{x}\| = 1$ . Then

13

$$(48) \|\mathbf{A}\| \ge \|\mathbf{A}\mathbf{x}\| = \|\lambda\mathbf{x}\| = |\lambda|\|\mathbf{x}\|.$$

Since this is true for all eigenvalues  $\lambda$  of **A** we must have  $\rho(\mathbf{A}) \leq \|\mathbf{A}\|$ . Since this is true for all induced matrix norms, it must also be true for their inf, i.e.  $\rho(\mathbf{A}) \leq \inf_{\|\cdot\|} \|\mathbf{A}\|$ .

We now show the reverse inequality  $\rho(\mathbf{A}) \leq \inf_{\|\cdot\|} \|\mathbf{A}\|$ . By theorem 9, for any  $\epsilon > 0$ , there is a non-singular matrix **X** such that  $\mathbf{X}^{-1}\mathbf{A}\mathbf{X} = \mathbf{D} + \mathbf{T}$ , where **D** is diagonal and **T** is strictly upper triangular with  $\|\mathbf{T}\|_{\infty} \leq \epsilon$  (why does the component by component result from 9 imply this inequality with the induced matrix norm?). Therefore:

(49) 
$$\|\mathbf{X}^{-1}\mathbf{A}\mathbf{X}\|_{\infty} = \|\mathbf{D} + \mathbf{T}\|_{\infty} \le \|\mathbf{D}\|_{\infty} + \|\mathbf{T}\|_{\infty} \le \rho(\mathbf{A}) + \epsilon.$$

It is possible to show that the norm  $\|\mathbf{A}\|_{\infty}' \equiv \|\mathbf{X}^{-1}\mathbf{A}\mathbf{X}\|_{\infty}$  is an induced matrix norm. Hence

(50) 
$$\inf_{\|\cdot\|} \|\mathbf{A}\| \le \|\mathbf{A}\|_{\infty}' \le \rho(\mathbf{A}) + \epsilon.$$

Since  $\epsilon > 0$  is arbitrary, we have  $\rho(\mathbf{A}) \leq \inf_{\|\cdot\|} \|\mathbf{A}\|$ .

2.4.1. *Convergence of Gauss-Seidel method.* As an application of the general theory above, we will determine that the Gauss-Seidel method converges when the matrix is diagonally dominant. Again, this is only a sufficient condition for convergence, the Gauss-Seidel method may converge for other matrices that are not diagonally dominant.

THEOREM 10. If **A** is diagonally dominant then the Gauss-Seidel method converges for any initial guess  $\mathbf{x}^{(0)}$ .

PROOF. We need to show that when **A** is diagonally dominant we have  $\rho(\mathbf{I} - \mathbf{Q}^{-1}\mathbf{A}) < 1$ . Let  $\lambda, \mathbf{x}$  be an eigenpair of  $\mathbf{I} - \mathbf{Q}^{-1}\mathbf{A}$  with  $\|\mathbf{x}\|_{\infty} = 1$ . Then  $(\mathbf{I} - \mathbf{Q}^{-1}\mathbf{A})\mathbf{x} = \lambda \mathbf{x}$ , or equivalently  $(\mathbf{Q} - \mathbf{A})\mathbf{x} = \lambda \mathbf{Q}\mathbf{x}$ . Written componentwise this becomes:

(51) 
$$-\sum_{j=i+1}^{n} a_{ij} x_j = \lambda \sum_{j=1}^{i} a_{ij} x_j, \ 1 \le i \le n.$$

Isolating the diagonal component:

(52) 
$$\lambda a_{ii} x_i = -\lambda \sum_{j=1}^{i-1} - \sum_{j=i+1}^n a_{ij} x_j, \ 1 \le i \le n.$$

Now pick the index *i* such that  $|x_i| = 1$  and write

(53) 
$$|\lambda||a_{ii}| \le |\lambda| \sum_{j=1}^{i-1} |a_{ij}| + \sum_{j=i+1}^{n} |a_{ij}|.$$

Isolating for  $\lambda$  and using the diagonal dominance of **A** we obtain

(54) 
$$|\lambda| \le \frac{\sum_{j=i+1}^{n} |a_{ij}|}{|a_{ii}| - \sum_{i=1}^{i-1} |a_{ij}|} < 1.$$

*Rev: 1990, January 17, 2012* 

na.tex

 $\square$ 

We conclude by noticing that this holds for all eigenvalues  $\lambda$  of  $\mathbf{I} - \mathbf{Q}^{-1}\mathbf{A}$ , and therefore must also hold for  $\rho(\mathbf{I} - \mathbf{Q}^{-1}\mathbf{A})$ .

- 2.5. Extrapolation.
  - 3. Conjugate gradient and other Krylov subspace methods