

Tridiagonal matrices If no pivoting is needed, then solving systems with a tridiagonal matrix is very efficient ($O(n)$ flops)

Here is an example of such a system:

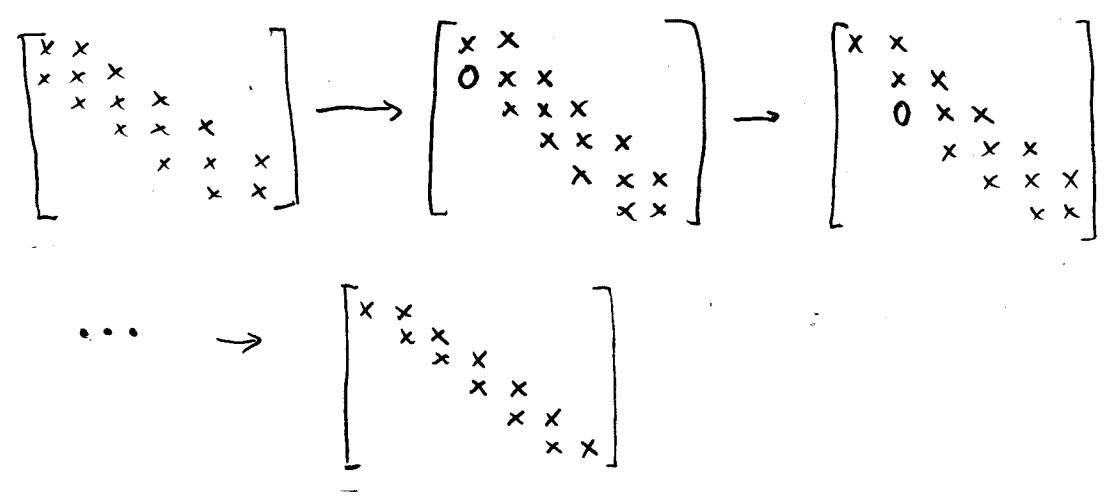
$$\begin{bmatrix}
 d_1 & c_1 & & & \\
 a_1 & d_2 & c_2 & & \\
 & \ddots & \ddots & \ddots & \\
 & & a_{n-1} & d_n & c_{n-1} \\
 & & & & a_{n-1} & d_n
 \end{bmatrix}
 \begin{bmatrix}
 x_1 \\
 x_2 \\
 \vdots \\
 x_n
 \end{bmatrix}
 =
 \begin{bmatrix}
 b_1 \\
 b_2 \\
 \vdots \\
 b_n
 \end{bmatrix}$$

In matlab we can write A using vectors a, d, c:

$$A = \text{diag}(a, -1) + \text{diag}(d, 0) + \text{diag}(c, 1);$$

In real life we only need to store 3 vectors.

Here is how Gaussian elimination looks like:



Introducing zeros below first equation:

$$d_2 = d_2 - \frac{a_1 c_1}{d_1}$$

$$b_2 = b_2 - \frac{a_1 b_1}{d_1}$$

(note: here we carry operations for L at the same time on RHS. we are not computing factor L !!)

for $i = 2, \dots, n$

$$\begin{cases} d_i = d_i - \frac{a_{i-1} c_{i-1}}{d_{i-1}} \\ b_i = b_i - \frac{a_{i-1} b_{i-1}}{d_{i-1}} \end{cases}$$

(update formulas have only one term)

Then to find x we do back substitution, which also simplifies:

$$x_n = b_n / d_n$$

for $i = n-1 : -1 : 1$

$$x_i = (b_i - c_i x_{i+1}) / d_i$$

Similar algorithms exist for band matrices (ie. only a few sub, super diagonals are non zero)

DIRECT SPARSE SOLVERS (brief intro)

A sparse matrix is a matrix with most of its entries being zero.

- saves memory as only non-zero elements are stored.
- saves also computation time as algorithms need only to access non-zero elements.

Here are two common ways of storing a sparse matrix

Coordinate format

iA = row index of nz entries

jA = column " " " "

vA = values of nz entries.

$$a_{iA(k), jA(k)} = vA(k)$$

Compressed sparse column format

PA = points to all nz elements in a column.

iA = row indices of nz entries in column

vA = values " " " " "

Example: $A = \begin{bmatrix} 4 & 0 & 3 & 0 \\ 0 & 1 & 0 & 1 \\ 0 & 2 & 3 & 0 \\ 5 & 3 & 0 & 6 \end{bmatrix}$

Coordinate

iA	jA	vA
1	1	4
4	1	5
2	2	1
3	2	2
4	2	3
1	3	3
3	3	3
2	4	1
4	4	6

CSC

PA	iA	vA
1	1	4
3	4	5
6	2	1
8	3	2
10	4	3
	1	3
	3	3
	2	1
	4	6

Some popular direct sparse solvers:

Super LU (Lawrence Berkeley National Lab)

MUMPS (French group: CERFACS, CNRS, INPT, INRIA)

UMFPACK (University of Florida, what is used by Matlab's backslash for solving sparse systems)

These are sophisticated libraries that can use a graph representation of a matrix to run the LU (or Cholesky) factorization symbolically and automatically switch between sparse and dense algorithms -

We are not going to see these algorithms in detail. Consider only the following examples to see how sparse structure reduces amount of computations. (these come from <http://www.cise.ufl.edu/~ndavis> the creator of UMFPACK homepage)

Matrix vector product: $y = Ax$

Dense version:

```
y = 0
for j = 1..n
  for i = 1..n
    | yi = yi + aijxj
```

$O(n^2)$

Sparse version:

```
y = 0
for j = 1..n
  for each i for which aij ≠ 0:
    | yi = yi + aijxj
```

$O(n \times (nz \text{ in a column}))$

Triangular solve: $Lx = b$ no log: $L_{ii} = 1$ (unit LT.)

(21)

$$(\Delta) | = 1$$

Dense version:

$$x = b$$

for $j = 1 \dots n$

$$\left[\begin{array}{l} \text{for } i = j+1 \dots n \\ x_i = x_i - l_{ij} x_j \end{array} \right.$$

$$O(n^2)$$

Sparse version #1

$$x = b$$

for $j = 1 \dots n$

$$\left[\begin{array}{l} \text{if } x_j \neq 0 \\ \text{for each } i > j \text{ with } l_{ij} \neq 0 \\ x_i = x_i - l_{ij} x_j \end{array} \right.$$

$$O(n + |b|)$$

Can be made even more efficient if we knew in advance the sparsity pattern of x .

let $\mathcal{K} = \{i \mid x_i \neq 0\} \equiv$ sparsity pattern of sol x .

$\mathcal{D} = \{i \mid b_i \neq 0\} \equiv$ sparsity pattern of RHS b .

If we knew \mathcal{K} then algo. becomes:

Sparse version #2

$$x = b$$

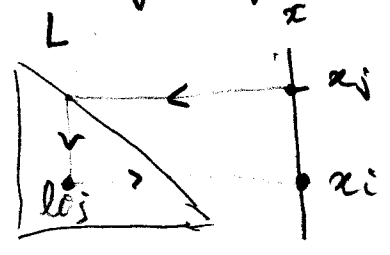
for $j \in \mathcal{K}$

$$\left[\begin{array}{l} \text{for each } i > j \text{ w/ } l_{ij} \neq 0 \\ x_i = x_i - l_{ij} x_j \end{array} \right.$$

$$O(|b|) \text{ flops}$$

(which can be better than version 1).

The idea to obtain X is to run triangular solve symbolically in graph representing L .



- ① $b_j \neq 0 \Rightarrow x_j \neq 0$
- ② $x_j \neq 0$ and $l_{ij} \neq 0 \Rightarrow x_i \neq 0$

- Replace L by an oriented graph with an edge $j \rightarrow i$ if $l_{ij} \neq 0$.
- $X \equiv$ closure (i.e. fixed point) of neighbor (\mathcal{R})

The enemy of sparse direct solvers is fill-in.

for example:

A sparse matrix may not necessarily have a sparse Choleky factor L .

- See sparse matrices Matlab demos.
- See also: "Direct Methods for Sparse Linear Systems"
T. Davis
→ has sample code which is simple enough to read.

Norms and the analysis of errors (§7.1 in textbook)

How are we sure that after doing $O(m^3)$ operations to solve a linear system using LU factor + backward substitution that errors we make don't pile up?

no we need notion of error in \mathbb{R}^n .

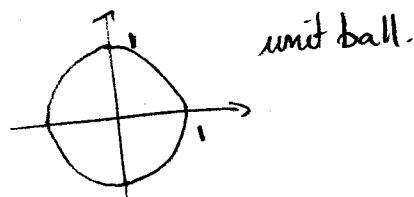
Def (norm) $\|x\|$ is a norm if

- i) $\|x\| \geq 0 \quad \forall x \in \mathbb{R}^n$
- ii) $\|x\| = 0 \Rightarrow x = 0$
- iii) $\|\lambda x\| = |\lambda| \|x\|, \quad \forall x \in \mathbb{R}^n, \lambda \in \mathbb{R}$
- iii) $\|x + y\| \leq \|x\| + \|y\|$ (triangle inequality)

Examples:

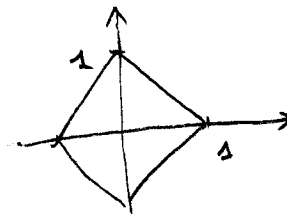
Euclidean norm

$$\|x\|_2 = \left(\sum_{i=1}^n x_i^2 \right)^{\frac{1}{2}}$$



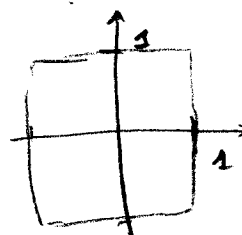
l_1 norm

$$\|x\|_1 = \sum_{i=1}^n |x_i|$$



l_∞ norm

$$\|x\|_\infty = \max_{i=1..n} |x_i|$$



Induced matrix norm

Let $\|\cdot\|$ be a norm in \mathbb{R}^n then we defined the induced matrix norm

$$\|A\| = \sup_{\|u\|=1} \|Au\| = \sup_{u \neq 0} \frac{\|Au\|}{\|u\|}$$

$\|\cdot\|$ satisfies all axioms of a norm and:

$$\|I\| = 1$$

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

Examples:

$$\begin{aligned} \|A\|_2^2 &= \sup_{x \neq 0} \frac{\|Ax\|^2}{\|x\|^2} = \sup_{x \neq 0} \frac{x^T A^T A x}{x^T x} \\ &= \lambda_1(A^T A) = \text{largest eigenvalue of } A^T A. \end{aligned}$$

$$\Rightarrow \|A\|_2 = \sqrt{\lambda_1(A^T A)}$$

(recall λ is an eigenvalue of matrix B if there exists $u \neq 0$ s.t.
 $Bu = \lambda u$.)

$$\|A\|_\infty = \sup_{\|u\|_\infty=1} \|Au\|_\infty$$

$$= \sup_{\|u\|_\infty=1} \max_{1 \leq i \leq n} |(Au)_i|$$

$$= \sup_{\|u\|_\infty=1} \max_{1 \leq i \leq n} \sum_j |a_{ij} u_j|$$

$$= \max_{1 \leq i \leq n} \sup_{\|u\|_\infty=1} \quad \quad \quad \leftarrow \text{sup attained when } u_j = \text{sgn } a_{ij}$$

$$= \max_{1 \leq i \leq n} \sum_j |a_{ij}| = \max_{1 \leq i \leq n} \ell_{i, \text{norm}} (\text{i-th row of } A)$$

Condition number

(25)

Suppose \tilde{b} is a perturbation of b .

If x, \tilde{x} satisfy:

$$Ax = b$$

$$A\tilde{x} = \tilde{b}$$

then how close are x and \tilde{x} ?

$$\begin{aligned}\|x - \tilde{x}\| &= \|A^{-1}b - A^{-1}\tilde{b}\| = \|A^{-1}(b - \tilde{b})\| \\ &\leq \|A^{-1}\| \|b - \tilde{b}\|\end{aligned}$$

We know giving a relative error is more informative ^{here it is.}

$$\begin{aligned}\|x - \tilde{x}\| &\leq \|A^{-1}\| \|b - \tilde{b}\| = \|A^{-1}\| \frac{\|Ax\|}{\|b\|} \|b - \tilde{b}\| \\ &\leq \|A^{-1}\| \|A\| \|x\| \frac{\|b - \tilde{b}\|}{\|b\|}\end{aligned}$$

$$\boxed{\frac{\|x - \tilde{x}\|}{\|x\|} \leq K(A) \frac{\|b - \tilde{b}\|}{\|b\|}}$$

Here $K(A)$ = condition number of A

$$= \|A\| \|A^{-1}\| \text{ (depends on choice of norm usually } \|\cdot\|_2 \text{ is used)}$$

$$\geq 1$$

\rightarrow quantifies how much we can trust solution x if there are errors in RHS.

\rightarrow There are ways of estimating $K(A)$ without computing A^{-1} (in a fast way)

Example :

$$A = \begin{bmatrix} 1 & 1+\epsilon \\ 1-\epsilon & 1 \end{bmatrix} \quad A^{-1} = \epsilon^{-2} \begin{bmatrix} 1 & -1-\epsilon \\ -1+\epsilon & 1 \end{bmatrix}$$

$$\|A\|_{\infty} = 2+\epsilon \quad \|A^{-1}\|_{\infty} = \epsilon^{-2} (2+\epsilon)$$

$$\Rightarrow \kappa(A) = \frac{(2+\epsilon)^2}{\epsilon^2}$$

If $\epsilon \leq 0.01$ then $\kappa(A) \geq 40000$.

\Rightarrow we lose 4 digits of precision!

As a rule of thumb:

$$\log_{10}(\kappa(A)) = \# \text{ of digits of precision lost in solution.}$$

Example : $\frac{\|b - \tilde{b}\|}{\|b\|} \approx 10^{-16}$ (Machine precision)

$$\kappa(A) \approx 10^{10}$$

$\Rightarrow \frac{\|z - \tilde{z}\|}{\|z\|} \leq 10^{-6}$ lost 10 digits of precision!

Convergence in \mathbb{R}^n

We say a sequence of vectors $v^{(k)} \in \mathbb{R}^n$ converges to $v \in \mathbb{R}^n$ as $k \rightarrow \infty$ if

$$\lim_{k \rightarrow \infty} \|v^{(k)} - v\| = 0$$

In \mathbb{R}^n the choice of norm does not matter as all norms are equivalent: (this is true for any finite dimensional space)

Two norms $\|\cdot\|$ and $\|\cdot\|'$ are equivalent if there are two constants c_1 and $c_2 > 0$ s.t.

$$c_1 \|x\| < \|x\|' \leq c_2 \|x\|$$

Also \mathbb{R}^n is complete meaning any sequence satisfying the Cauchy criterion converges.

$$\forall \epsilon > 0 \exists N \text{ s.t. } \forall i, j \geq N \ \|v^{(i)} - v^{(j)}\| < \epsilon$$

Normann Series

Let $A \in \mathbb{R}^{n \times n}$ be such that $\|A\| < 1$, then

i) $I - A$ is invertible

$$\text{ii) } (I - A)^{-1} = \sum_{k=0}^{\infty} A^k$$

Proof: Assume for contradiction that $I - A$ is singular

$$\Rightarrow \exists z \in \mathbb{R}^n \text{ s.t. } (I - A)z = 0, \text{ take w.l.o.g. } z \text{ s.t. } \|z\| = 1$$

then:

$$1 = \|z\| = \|Az\| < \underbrace{\|A\|}_{< 1} \|z\| < \|z\| = 1. \text{ Contradiction!}$$

Now we need to show that

(28)

$$\sum_{k=0}^m A^k \rightarrow (I-A)^{-1} \quad (\text{ie. partial sums converge})$$

Notice:

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

↙ telescoping series.

Now using properties of induced matrix norm: ($\|AB\| \leq \|A\| \|B\|$)

$$\|A^{m+1}\| \leq \|A\|^m \rightarrow 0$$

Thus:

$$\left\| (I-A) \sum_{k=0}^m A^k - I \right\| \leq \|A\|^{m+1} \rightarrow 0 \text{ as } m \rightarrow \infty.$$

Q.E.D.

Note: This theorem is a very intuitive generalization of what you already know about geometric series.

$$\frac{1}{1-x} = 1+x+x^2+\dots, \text{ when } |x| < 1.$$

Here is the kind of result we can show with Neumann Series:

$$\|(I-A)^{-1}\| \leq \sum_{k=0}^{\infty} \|A^k\| \leq \sum_{k=0}^{\infty} \|A\|^k = \frac{1}{1-\|A\|}$$

Note: This theorem is very useful for the theory, however it is seldom used in practice. (There are faster iterative methods for solving $Ax = b$).

Here is an easy but handy generalization of

Neuman Series:

Theorem

If A and B are matrices s.t. $\|I-AB\| < 1$ then
 A, B are invertible and

$$A^{-1} = B \sum_{k=0}^{\infty} (I-AB)^k$$

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

Proof: By Neumann Series:

$I - (I-AB) = AB$ is invertible and:

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

$$\Rightarrow A^{-1} = B B^{-1} A^{-1} = B \sum_{k=0}^{\infty} (I-AB)^k$$

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

Iterative refinement

Let $x^{(0)}$ be an approx. sol to

$$Ax = b.$$

then:

$$x = x^{(0)} + A^{-1}(b - Ax^{(0)})$$

We often call

$$r^{(0)} = b - Ax^{(0)} = \text{residual vector}$$

$$e^{(0)} = A^{-1}r^{(0)} = \text{error vector.}$$

Algorithm

$x^{(0)}$ given

for $k=0, 1, \dots$

$$\left\{ \begin{array}{l} r^{(k)} = b - Ax^{(k)} \\ Ae^{(k)} = r^{(k)} \\ x^{(k+1)} = x^{(k)} + e^{(k)} \end{array} \right.$$

← usually inexact solve

Solution can be improved even if the solves are inexact for example:

- LU factorization is stopped prematurely (\equiv incomplete LU factor)
- $Ae^{(k)} = r^{(k)}$ is solved with an iterative method
- A is not known exactly.

To be more precise about what we mean by "solving approximately"

Let $B \equiv$ "approximate" inverse of A

$$\approx A^{-1} \text{ (in some sense we shall see)}$$

The iterates are:

$$x^{(0)} = Bb$$

$$x^{(k+1)} = x^{(k)} + B(b - Ax^{(k)}), \quad k \geq 0.$$

Looking at a few iterates we see a pattern emerge.

$$x^{(0)} = Bb$$

$$x^{(1)} = x^{(0)} + B(b - Ax^{(0)})$$

$$= Bb + B(b - ABb)$$

$$= Bb + B(I - AB)b$$

$$x^{(2)} = x^{(1)} + B(b - Ax^{(1)})$$

$$= Bb + B(I - AB)b + B(b - A(Bb + B(I - AB)b))$$

$$= Bb + B(I - AB)b + B((I - AB)b - AB(I - AB)b)$$

$$= Bb + B(I - AB)b + B(I - AB)^2 b$$

$$x^{(m)} = B \sum_{k=0}^m (I - AB)^k b$$

Theorem (Iterative refinements)

Iterates from iterative refinement are:

$$x^{(m)} = B \sum_{k=0}^m (I - AB)^k b \quad (m \geq 0)$$

and because of the generalized Neuman series:

$$x^{(m)} \rightarrow x \quad \text{as } m \rightarrow \infty \quad \text{when } \|I - AB\| < 1.$$

proof: can be done by induction on m .

Solving equations with iterative methods

Idea: Instead of having a direct method which finds solution in finitely many steps, use an iterative method to find successive approx that converge to sol.

Motivational example

$$\begin{bmatrix} 7 & -6 \\ -8 & 9 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 3 \\ -4 \end{bmatrix}$$

Jacobi method: solve i -th equation for i -th unknown:

$$\begin{cases} x_1^{(k)} = \frac{3}{7} + \frac{6}{7} x_2^{(k-1)} \\ x_2^{(k)} = \frac{-4}{9} + \frac{8}{9} x_1^{(k-1)} \end{cases}$$

Gauss Seidel method: use improved value of $x_1^{(k)}$ in second equation.

$$\begin{cases} x_1^{(k)} = \frac{3}{7} + \frac{6}{7} x_2^{(k-1)} \\ x_2^{(k)} = \frac{-4}{9} + \frac{8}{9} x_1^{(k)} \end{cases} \quad \text{converges faster!}$$

Note: convergence for these methods depends on initial guess $(x_1^{(0)}, x_2^{(0)})$

Matrix splitting

Consider the system $Ax = b$.

Introduce a non-singular matrix Q and rewrite system as:

$$Qx = (Q - A)x + b$$

Get iteration:

$$Qx^{(k)} = (Q - A)x^{(k-1)} + b, \quad k \geq 1, \quad x^{(0)} \text{ given}$$

Q should have the desirable properties:

- ① $x^{(k)} \rightarrow x$ where x is a sol. of $Ax=b$
- ② Solving systems with Q is cheap.

If $x^{(k)} \rightarrow x$ then the limit x must satisfy:

$$Qx = (Q-A)x + b \Rightarrow Ax = b$$

Assuming A & Q are non-singular:

$$x^{(k)} = (I - Q^{-1}A)x^{(k-1)} + Q^{-1}b$$

$$x = (I - Q^{-1}A)x + Q^{-1}b$$

$$x^{(k)} - x = (I - Q^{-1}A)(x^{(k-1)} - x)$$

$$\Rightarrow \|x^{(k)} - x\| \leq \|I - Q^{-1}A\| \|x^{(k-1)} - x\|$$

⋮

$$\leq \|I - Q^{-1}A\|^k \|x^{(0)} - x\|$$

Thus when $\|I - Q^{-1}A\| < 1$ we have $x^{(k)} \rightarrow x$, where $Ax = b$.

Theorem If $\|I - Q^{-1}A\| < 1$ for some induced matrix norm

then seq.

$$Qx^{(k)} = (Q-A)x^{(k-1)} + b$$

converges to sol. of $Ax = b$ for any starting vector $x^{(0)}$.