

## Recitation 7

Recitation Instructor: Shivam Verma

Email: [shivamverma@nyu.edu](mailto:shivamverma@nyu.edu)

Ph: 718-362-7836

Office hours: WWH 605 (2.50 - 4.50 pm, Tuesdays)

### Brief Overview

- Review for midterm
  - Floating point system
  - Solving equations iteratively - Iteration, Newton etc.
  - LU decomposition, QR, least squares, norms
  - Special matrices - tridiagonal, banded, Householders
  - Eigen value methods - Jacobi, QR, Sturm, Inverse Iteration, Rayleigh, Gershgorin discs
- Quiz 2 review

**Note:** I've added useful resources from previous recitations, as well as some new links. Going through them would be useful revision

### 1. Floating point system [References: 1-4]

- Due to finite precision in floating point number representation, there are gaps between consecutive numbers.
- Size of these gaps depends on the size of the number and on the precision (e.g., double or single precision).
- MATLAB has the function `eps()`, which returns, for a number, the distance to the next floating point number in the same precision.

Examples:

- `eps(1)`
- `eps(single(1))`
- `eps(2^(40))`
- `eps(single(2^(40)))`

### 2. Solving equations iteratively [References: 5-11]

$$\lim_{k \rightarrow \infty} |\epsilon_{k+1}|/|\epsilon_k| = \lim_{k \rightarrow \infty} |x_{k+1} - \xi|/|x_k - \xi| = \mu$$

- If  $\mu = 0$ , converges superlinearly

- If  $\mu \in (0, 1)$ , converges linearly with asymptotic rate of convergence  $\rho = -\log_{10}\mu$
- If  $\mu = 1$ , converges sublinearly

For linearly convergent systems,  $\rho$  measures number of correct decimal digits gained in one iteration.

<u>Method</u>	<u>Step Equation</u>	<u>Rate of Convergence</u>
Iteration / Fixed Point	$x_{k+1} = g(x_k)$	Sub, linear or Super <ul style="list-style-type: none"> <li>• <math>\lim_{k \rightarrow \infty}  \epsilon_{k+1} / \epsilon_k  =  g'(\xi) </math></li> <li>• <math> g'(\xi)  \in (0, 1)</math>, then converges linearly with <math>\rho = -\log_{10} g'(\xi) </math></li> <li>• <math> g'(\xi)  &gt; 1</math>, does not converge</li> </ul>
Newton	$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}$	“Ultimately” Quadratic <ul style="list-style-type: none"> <li>• <math>\lim_{k \rightarrow \infty}  x_{k+1} - \xi / x_k - \xi ^2 =  f''(\xi) /2 f'(\xi) </math></li> </ul>
Bisection	$x_{k+1} = (a_{k+1} + b_{k+1})/2$ where <ul style="list-style-type: none"> <li>• <math>(a_{k+1}, b_{k+1}) = (a_k, x_k)</math> if <math>f(x_k)f(b_k) &gt; 0</math></li> <li>• <math>(a_{k+1}, b_{k+1}) = (x_k, b_k)</math> if <math>f(x_k)f(b_k) &lt; 0</math></li> </ul>	If $[a_0, b_0]$ chosen such that $f(a_0)f(b_0) < 0$ , then after k iterations, soln lies in interval of length $(a_0 - b_0)/2^k$ <ul style="list-style-type: none"> <li>• <math>\rho = \log_{10}2</math></li> </ul>
Secant	$x_{k+1} = x_k - f(x_k) \left( \frac{x_k - x_{k-1}}{f(x_k) - f(x_{k-1})} \right)$	Faster than linear, less than quadratic <ul style="list-style-type: none"> <li>• <math> \epsilon_{k+1}  \leq 2/3  \epsilon_k </math>, <math>\rho</math> atleast <math>\log_{10}3/2</math></li> <li>• Precisely:               <math display="block">\lim_{k \rightarrow \infty}  x_{k+1} - \xi / x_k - \xi ^q = \left(  f''(\xi) /2 f'(\xi)  \right)^{q/(1+q)}</math>               Where <math>q = \frac{1}{2} \left( 1 + \sqrt{5} \right) \approx 1.618</math> </li> </ul>

### 3. Solving linear systems [References: 12-21]

<u>Method</u>	<u>Algorithm</u>	<u>Computation cost</u>
LU decomposition	<ul style="list-style-type: none"> <li>Break up matrix into product of upper and lower triangular matrices  <math>A = LU</math></li> <li>Solve two triangular systems of equations  <math>Ax = b \Rightarrow LUx = b</math>  <math>Ly = b, Ux = y</math></li> </ul>	<ul style="list-style-type: none"> <li>Factorization:  <math>2n^3/3 - n^2/2</math></li> <li>Solving triangular sys.:  <math>n(n-1) + n^2</math></li> <li>Total  <math>2n^3/3 + 3n^2/2 \approx 2n^3/3</math></li> </ul>
Cholesky Decomp.	<ul style="list-style-type: none"> <li>If A is PSD, <math>A = LL^T</math></li> <li><math>Ax = b \Rightarrow LL^T x = b</math></li> <li>Solve: <ul style="list-style-type: none"> <li><math>Ly = b</math></li> <li><math>L^T x = y</math></li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li><math>n^3/3</math></li> <li>Half cost compared to LU</li> </ul>
Linear least squares	<ul style="list-style-type: none"> <li>Treat linear system as optimization problem:  <math>\min_x \ Ax - b\ _2^2</math></li> <li>Solve Normal Eqn.  <math>A^T Ax = A^T b</math> using Cholesky</li> </ul>	<ul style="list-style-type: none"> <li><math>O(mn^2)</math> using Cholesky/LU</li> </ul>
QR Decomp.	<ul style="list-style-type: none"> <li>Break up matrix into product of orthogonal and upper triangular matrices:  <math>A = QR</math>  <math>= [Q_1 Q_2] [R_1 0]^T = Q_1 R_1</math></li> <li>Normal Eqn. equiv. to solving a triangular system  <math>R_1 x = Q_1^T b</math></li> </ul>	<p><b>Givens rotations</b>  Use sequence of rotations in 2D subspaces:  For <math>m \approx n</math>: <math>\sim n^2/2</math> square roots, and <math>4/3n^3</math> multiplications  For <math>m \gg n</math>: <math>\sim nm</math> square roots, and <math>2mn^2</math> multiplications</p> <p><b>Householder reflections</b>  Use sequence of reflections in 2D subspaces  For <math>m \approx n</math>: <math>2/3n^3</math> multiplications  For <math>m \gg n</math>: <math>2mn^2</math> multiplications</p>

**Least Squares [See Ref. 21]**

Taking the case where  $m \geq n$ ,

- To solve  $Ax = b$ , minimize the 'residual sum of squares' or 'mean square error' or 'squared euclidean norm'
- Optimization problem:
  - $\min_x \|Ax - b\|_2^2$
  - Has a closed-form solution, known as the **normal equation**:
    - $A^T Ax = A^T b$
  - Multiple ways of solving

### Solve Normal Equation using LU, Cholesky etc.

- If  $A$  has full rank,  $A^T A$  is invertible. In general,  $A^T A$  is a symmetric positive definite. How?
  - $x^T A^T Ax = (Ax)^T (Ax) = \|Ax\|_2^2 \geq 0$
  - This property is very useful in general (see Cholesky decomposition).
- Can use the usual methods (LU, Cholesky etc.) to solve this linear system in  $O(mn^2)$ .
- Disadvantage:
  - Computing
  - May be ill-conditioned, as  $k(A^T A) = k(A)^2$

### QR decomposition

$$A = \begin{bmatrix} Q_1 & Q_2 \end{bmatrix} \begin{bmatrix} R_{11} \\ 0 \end{bmatrix} = Q_1 R_{11},$$

$$\begin{aligned} \|Ax - b\|^2 &= \|Q^T(Ax - b)\|^2 \\ &= \left\| \begin{bmatrix} R_{11} \\ 0 \end{bmatrix} x - \begin{bmatrix} Q_1^T b \\ Q_2^T b \end{bmatrix} \right\|^2 \\ &= \|R_{11}x - Q_1^T b\|^2 + \|Q_2^T b\|^2. \end{aligned}$$

- Since second term is independent of  $x$ , the minimum can be achieved when:
  - $R_{11}x = Q_1^T b$

- This is a triangular linear system. Can be solved in  $O(n^2)$
- This decomposition exists for any matrix - rectangular, non-symmetric etc.
- How can we calculate a QR decomposition?

#### Givens rotations

Use sequence of rotations in 2D subspaces:

For  $m \approx n$ :  $\sim n^2/2$  square roots, and  $4/3n^3$  multiplications

For  $m \gg n$ :  $\sim nm$  square roots, and  $2mn^2$  multiplications

#### Householder reflections

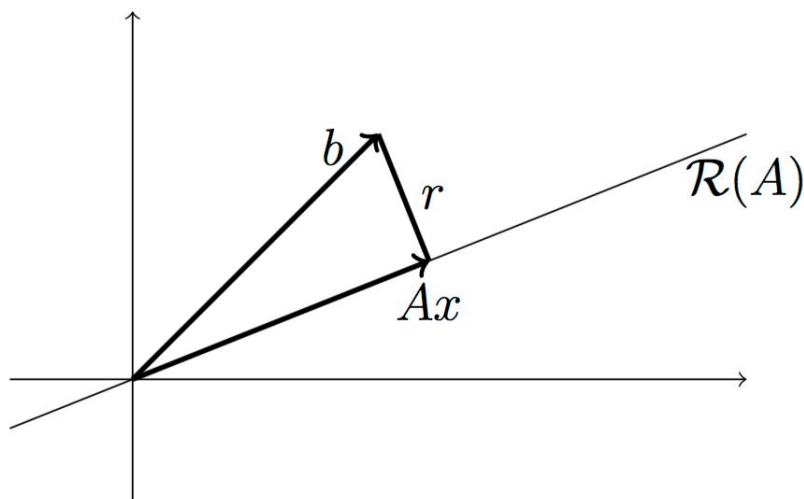
Use sequence of reflections in 2D subspaces

For  $m \approx n$ :  $2/3n^3$  multiplications

For  $m \gg n$ :  $2mn^2$  multiplications

- See textbook or Deuffhard/Hohmann for proof and discussion.
- Advantage: Better conditioned than least-squares, as  $k(R_1) = k(A)$ . How?
- $k(A^T A) = k(R_1^T Q_1^T Q_1 R_1) = k(R_1^T R_1)$

### Geometric interpretation of least squares



- $A^T(Ax - b) = 0 \Rightarrow A^T r = 0$  where  $r$  is the residual
- This means residual vector is orthogonal to any vector in the range of  $A$

- $\|Ax\|^2 + \|r\|^2 = \|b\|^2$
- Thus, least squares solves for the projection of 'b' on the range space of 'Ax', or, it solves  $Ax = b_{\text{projected}}$ , where  $b_{\text{projected}} = b \cdot \cos(\theta)$
- If  $\theta \approx \pi/2$ , then  $b \cdot \cos(\theta) \approx 0$ , and corresponding solution will be bad (model doesn't fit data!)
- In general, it may be that columns of A are nearly linearly dependent, in which case problem becomes ill-conditioned, as  $A^T A$  is not invertible.
  - One approach is called **regularization**. It involves adding a strictly positive constant to the diagonal elements to make eigenvalues non-zero.
  - $(A^T A + \lambda I)x = A^T b$
  - This is the solution of the minimization problem:
    - $\min_x \|Ax - b\|_2^2 + \lambda \|x\|_2^2$
    - This is known as L2-regularization, since the "regularization" term involves an L2-norm
  - Can you say whether we can use an L1-norm instead of the L2-norm for regularization? Is there a closed-form solution for this? why/why not?

## Norms and Condition Numbers

- An  $L_p$  norm is defined as  $\|x\|_p = \left( \sum_{i=1}^n |x_i|^p \right)^{1/p}$
- 2-norm or Euclidean norm:  $\|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$
- 1-norm:  $\|x\|_1 = \sum_{i=1}^n |x_i|$
- $\infty$ -norm:  $\|x\|_\infty = \max(|x_1|, |x_2|, \dots, |x_n|)$
- Relative condition number
  - $\sup_x (\|f\|/\|f(x)\|) / (\|x\|/\|x\|)$
  - For a matrix (image from Trefethen et al., Numerical Linear Algebra, p. 93)

$$\kappa = \sup_{\delta x} \left( \frac{\|A(x + \delta x) - Ax\|}{\|Ax\|} \right) / \left( \frac{\|\delta x\|}{\|x\|} \right) = \sup_{\delta x} \frac{\|A\delta x\|}{\|\delta x\|} / \frac{\|Ax\|}{\|x\|}$$

that is,

$$\kappa = \|A\| \frac{\|x\|}{\|Ax\|}$$

- Relative condition number more important in numerical analysis, as floating point system introduces relative errors
- Small condition number means well-conditioned. Large means ill-conditioned.

- Condition number of A:

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

#### 4. Special matrices [References: 18-19]

- Banded matrices
  - $a_{ij} = 0$  for all i and j such that  $|i - j| < k$  for a k-band matrix
  - Eg. tri-diagonal matrix
- Positive-definite matrices
  - $x^T A x > 0$  for all non-zero x in real space
- What is a Householder matrix? Where is it used (hint: QR algorithm)? Why is it useful?

#### 5. Eigenvalue problems [References: 22-32]

<u>Method</u>	<u>Idea</u>	<u>Algorithm</u>
<b>Jacobi</b>	<ul style="list-style-type: none"> <li>• Use orthogonal transformations (pre- and post- multiply) to convert matrix to diagonal form.</li> </ul>	$R(\varphi) = \begin{pmatrix} \cos \varphi & \sin \varphi \\ -\sin \varphi & \cos \varphi \end{pmatrix}$ $\varphi = \frac{1}{2} \tan^{-1} \frac{2a_{pq}}{a_{qq} - a_{pp}}$
<b>Sturm Sequence</b>	<ul style="list-style-type: none"> <li>• # of consecutive sign agreements in sequence <math>p(\lambda)</math> = # of eig. Values <math>&gt; \lambda</math></li> <li>• Take interval using Gershgorin theorem, use bisection method to find any eig. value</li> </ul>	$T = \begin{pmatrix} a_1 & b_2 & & & & \\ b_2 & a_2 & b_3 & & & \\ & b_3 & a_3 & b_4 & & \\ & & \dots & \dots & \dots & \\ & & & \dots & \dots & \dots \\ & & & & \dots & \dots \\ & & & & & b_{n-1} & a_{n-1} & b_n \\ & & & & & & b_n & a_n \end{pmatrix}$ $p_r(\lambda) = (a_r - \lambda)p_{r-1}(\lambda) - b_r^2 p_{r-2}(\lambda), \quad r = 2, 3, \dots, n$
<b>Inverse Iteration</b>	<ul style="list-style-type: none"> <li>• Take an estimate of eigenvalue</li> <li>• Iterate to find corresponding eigenvector</li> </ul>	$(A - \vartheta I)w^{(k)} = v^{(k)},$ $v^{(k+1)} = c_k w^{(k)},$ <p>where <math>c_k = 1/\ w^{(k)}\ _2</math>, the sequence <math>\{v^{(k)}\}</math> converges to the normalized eigenvector <math>\bar{v}</math> for the eigenvalue <math>\lambda</math> closest to <math>\vartheta</math>.</p>

<b>QR</b>	<ul style="list-style-type: none"> <li>Take symmetric tridiagonal matrix, and convert to upper-diagonal</li> </ul>	<ul style="list-style-type: none"> <li><math>A_k = Q_k R_k</math></li> <li><math>A_{k+1} = R_k Q_k</math></li> <li>Can add shift parameter</li> </ul>
<b>Rayleigh coefficient</b>	<ul style="list-style-type: none"> <li>If we have a fairly close approximation of the eigenvector, <math>R(x)</math> gives us a good approximate of the corresponding eigenvalue</li> <li>Th. 5.12:  <math>\lambda_{min} \leq R(x) \leq \lambda_{max}</math></li> </ul>	$R(x) = \frac{x^T A x}{x^T x}$ where A is symmetric. <ul style="list-style-type: none"> <li>If <math>x</math> is eigenvector, <math>R(x)</math> is corresponding eigenvalue</li> <li>Otherwise, if</li> </ul> $x = \sum_{j=1}^n \alpha_j x^{(j)},$ $R(x) = \frac{\sum_{j=1}^n \lambda_j \alpha_j^2}{\sum_{j=1}^n \alpha_j^2}$

## Gershgorin's theorems

Theorem 1:

*Every eigenvalue of matrix  $A_{nn}$  satisfies:*

$$|\lambda - A_{ii}| \leq \sum_{j \neq i} |A_{ij}| \quad i \in \{1, 2, \dots, n\}$$

*Every eigenvalue of a matrix A must lie in a Gershgorin disc corresponding to the columns of A.*

Theorem 2:

*A Subset G of the Gershgorin discs is called a disjoint group of discs if no disc in the group G intersects a disc which is not in G. If a disjoint group G contains r nonconcentric discs, then there are r eigenvalues.*



## Quiz 2 discussion

### **Q1:**

Find the Householder transformation matrix which maps the column vector  $(1, 1, 1, 1, 1, 1)^T$  into a vector of the form  $(1, 1, 1, *, 0, 0)^T$  and determined the value of the fourth element of the second vector marked by \*.

### **Soln.**

If  $x$  and  $y$  have the same length ( $\text{norm}(x) = \text{norm}(y)$ ),

$Hx = y$  if  $H = I - 2uu^T$  where  $u = (x - y)/\|x - y\|$  (definition of Householder transformation).

Therefore, in our case,  $x = [1, 1, 1, 1, 1, 1]^T, y = [1, 1, 1, *, 0, 0]^T$ . It must be that  $\|x\| = \|y\|$ , therefore,  $* = \pm\sqrt{3}$ , and  $u$  is given by the above formula.

For theory and example, see <http://web.csulb.edu/~tgao/math423/s93.pdf>.

## Helpful links

1. [https://en.wikipedia.org/wiki/IEEE\\_floating\\_point](https://en.wikipedia.org/wiki/IEEE_floating_point)
2. Numerical Computing with IEEE Floating Point Arithmetic, Michael L. Overton (NYU)
3. Sec. 2.5, Numerical Mathematics, Alfio Quarteroni et al.
4. Sec. 2.1, Numerical Analysis in Modern Scientific Computing, Peter Deufhard and Andreas Hohmann.
5. [https://www.math.ust.hk/~mamu/courses/231/Slides/ch02\\_2b.pdf](https://www.math.ust.hk/~mamu/courses/231/Slides/ch02_2b.pdf)
6. [A casual \(but interesting\) introduction to Newton's method](#)
7. [A good lecture on rates of convergence](#)
8. [A proof of Newton's method's quadratic convergence](#)
9. [Newton's method, complex numbers and pretty fractals](#)
10. [Newton's method on functions with multiple same roots](#)
11. [Secant method and why it's order of convergence is the golden ratio](#)
12. [More about cost of computation](#)
13. [Wiki page for LU decomp.](#)
14. [MIT OCW example for LU](#)
15. [LU decomp. explained very well here](#)
16. [Detailed notes on LU decomp. and linear systems](#)
17. [Cholesky decomposition explained](#)
18. [About householder transformation - 1](#)
19. [About householder transformation - 2](#)
20. *Numerical Mathematics (Quarteroni et al)* is a very good resource for LU/Cholesky/etc. It's also freely available for NYU students via Springer!
21. <http://www.cs.cornell.edu/~bindel/class/cs3220-s12/notes/lec11.pdf>

22. [Nice reference on Gerschgorin's theorems](#)
23. [Sturm's theorem](#)
24. [Proof of Sturm's sequence property](#)
25. [Using Sturm's theorem for finding eigenvalues \(sec. 4.6.2\)](#)
26. [QR algorithm explained quite well](#)
27. [QR algorithm proof](#)
28. [Rayleigh coefficient and Inverse iteration](#)
29. *Numerical Linear Algebra, Trefethen & Bau* also has a nice discussion on condition numbers (see Lec. 12).
30. [https://en.wikipedia.org/wiki/Wilkinson%27s\\_polynomial](https://en.wikipedia.org/wiki/Wilkinson%27s_polynomial)
31. <http://blogs.mathworks.com/cleve/2013/03/04/wilkinsons-polynomials/>
32. [http://college.cengage.com/mathematics/larson/elementary\\_linear/5e/students/ch08-10/chap\\_10\\_3.pdf](http://college.cengage.com/mathematics/larson/elementary_linear/5e/students/ch08-10/chap_10_3.pdf)