

<< Vectors & Vector Space

> vector - object comprised of magnitude & direction

- starts at origin, given coordinate of its tip

- transpose - $(r_1, \dots, r_n)^T = \begin{pmatrix} r_1 \\ \vdots \\ r_n \end{pmatrix}$

- vector addition - $\vec{u} + \vec{v} = [u_1 + v_1, \dots, u_n + v_n]^T$

- scalar multiple - $\alpha \vec{u} = [\alpha u_1, \dots, \alpha u_n]^T$

→ Vector Space

> objects that addition & scalar multiple are defined in:

A1. If \vec{u} and \vec{v} are in V , then $\vec{u} + \vec{v}$ is defined and in V

A2. Commutativity : $\vec{u} + \vec{v} = \vec{v} + \vec{u}$

A3. Associativity : $\vec{u} + (\vec{v} + \vec{w}) = (\vec{u} + \vec{v}) + \vec{w}$

A4. Additive identity : there is a vector $\vec{0} \in V$ such that $\vec{u} + \vec{0} = \vec{u}$ for all $\vec{u} \in V$

A5. Additive inverse : for all $\vec{u} \in V$, there is a vector $-\vec{u} \in V$ such that $\vec{u} + (-\vec{u}) = \vec{0}$

M1. If α is scalar, $\alpha \vec{u}$ is defined and in V

M2. Distributive & $\alpha(\vec{u} + \vec{v}) = \alpha \vec{u} + \alpha \vec{v}$

Associative property $(\alpha + \beta) \vec{u} = \alpha \vec{u} + \beta \vec{u}$

$$\alpha \beta \vec{u} = \alpha(\beta \vec{u})$$

M3. $1 \vec{u} = \vec{u}$

$0 \vec{u} = \vec{0}$

$(-1)\vec{u} = -\vec{u}$

> field of scalar - \mathbb{R} or \mathbb{C}

→ Examples of Vector Spaces

- \mathbb{K}^n

- space of polynomials of degree 2 or less : $a_0 + a_1 x + a_2 x^2$

> subset - a collection of some/all of elements of V

> subspace - a subset that's a vector space closed under definitions

- e.g. $(c, c)^T$

2 May 2020

Linear Independence & Dependence, Span, Basis

AMATH 352

<< Linear Combination

> linear independent - if the only scalars c_1, \dots, c_m for which $c_1\vec{v}_1 + \dots + c_m\vec{v}_m = \vec{0}$
are $c_1 = \dots = c_m = 0$

> linear dependent - there exist scalars c_1, \dots, c_m , not all zero for $c_1\vec{v}_1 + \dots + c_m\vec{v}_m = \vec{0}$

> linear combination - sum of scalar multiples of vectors
(LC)

$$c_1\vec{v}_1 + \dots + c_m\vec{v}_m = \sum_{j=1}^m c_j\vec{v}_j$$

> nontrivial linear combination - LC that the scalars are not all zero

- vectors are linear dependent if there is a nontrivial LC equal to zero vector

<< Span

> span of vectors $\vec{v}_1, \dots, \vec{v}_m$ - set of all linear combinations of $\vec{v}_1, \dots, \vec{v}_m$

- If $\vec{v}_1, \dots, \vec{v}_m$ is in V , then $\text{span}(\vec{v}_1, \dots, \vec{v}_m)$ is subspace of V

- $\vec{v}_1, \dots, \vec{v}_m$ span all of V if every vector $\vec{v} \in V$ can be written as LC of $\vec{v}_1, \dots, \vec{v}_m$.

<< Basis

> basis - $\vec{v}_1, \dots, \vec{v}_m$ are linear independent and span V For \mathbb{R}^n , any set of n linear independent

- e.g. $(1, 0)^T, (0, 1)^T$ basis of \mathbb{R}^2

- $\{(1, 0)^T, (0, 1)^T, (1, 1)^T\}$ not basis of \mathbb{R}^2 : not linear independent

- $\{(1, 1)^T, (2, 2)^T\}$ not basis of \mathbb{R}^2 : do not span V

<< Standard Basis & Dimension

- infinite choices of basis for V \rightarrow all basis have same # of elements

> dimension - number of vectors in any basis of V

- e.g. \mathbb{R}^2 dimension = 2 : $(1, 0)^T, (0, 1)^T$

- $\dim(\text{polynomial degree } \leq 2) = 3 : 1, x, x^2$

> standard basis - in \mathbb{R}^n , the set of basis vectors $\vec{e}_1, \dots, \vec{e}_n$ where \vec{e}_j has 1 in position j and 0 everywhere else.

<< Coordinate

- Given a basis for V , any vector $\vec{v} \in V$ can be written uniquely as LC of the basis vectors

$$\vec{v} = \sum_{j=1}^m c_j \vec{v}_j$$

> coordinate - the scalars c_1, \dots, c_m corresponding to basis $\vec{v}_1, \dots, \vec{v}_m$

<< Dot Product

> dot product (inner product) (scalar product) of \vec{u} and \vec{v} in \mathbb{R}^n :

$$\vec{u} \cdot \vec{v} \equiv \langle \vec{u}, \vec{v} \rangle = \sum_{j=1}^n u_j v_j$$

> inner product - any rule that takes a pair of vectors \vec{u} and \vec{v} and maps them to a real number denoted $\langle \vec{u}, \vec{v} \rangle$

1. $\langle \vec{u}, \vec{v} \rangle = \langle \vec{v}, \vec{u} \rangle$
2. $\langle \vec{u}, \vec{v} + \vec{w} \rangle = \langle \vec{u}, \vec{v} \rangle + \langle \vec{u}, \vec{w} \rangle$
3. $\langle \vec{u}, \alpha \vec{v} \rangle = \alpha \langle \vec{u}, \vec{v} \rangle$ for real number α
4. $\langle \vec{u}, \vec{u} \rangle \geq 0$ with equality if and only if $\vec{u} = \vec{0}$.

> norm - $\|\vec{u}\| = \sqrt{\langle \vec{u}, \vec{u} \rangle} = \sqrt{\sum_{j=1}^n u_j^2}$

<< Orthonormality

> orthogonal - inner product of two non-zero vector is zero: $\langle \vec{u}, \vec{v} \rangle = 0$

- standard basis vector for \mathbb{R}^n are orthogonal

> orthonormal - orthogonal vectors that have norm of 1

- standard basis vector for \mathbb{R}^n are orthonormal

- $\langle \vec{e}_i, \vec{e}_j \rangle = 0, i \neq j$

- $\|\vec{e}_i\| = 1, i = 1, \dots, n$

<< Gram-Schmidt Algorithm

> construct orthonormal set $\vec{q}_1, \dots, \vec{q}_m$ from linearly independent vectors $\vec{v}_1, \dots, \vec{v}_m$

- $\tilde{q}_j = \vec{v}_j - \sum_{i=1}^{j-1} \langle \vec{v}_j, \vec{q}_i \rangle \vec{q}_i$ e.g. $\tilde{q}_1 = \vec{v}_1$ $\tilde{q}_2 = \vec{v}_2 - \langle \vec{v}_2, \tilde{q}_1 \rangle \tilde{q}_1$

- $\vec{q}_j = \frac{\tilde{q}_j}{\|\tilde{q}_j\|}$ $\vec{q}_1 = \frac{\tilde{q}_1}{\|\tilde{q}_1\|}$ $\vec{q}_2 = \frac{\tilde{q}_2}{\|\tilde{q}_2\|}$

Matrices & Linear System

> Matrix - rectangular array of elements

> Gaussian elimination - eliminating variables by adding / subtracting multiples of one equation to another

$$\begin{bmatrix} 1 & 5 & 9 \\ 2 & 6 & 10 \\ 3 & 7 & 11 \\ 4 & 8 & 12 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} x+5y+9z \\ 2x+6y+10z \\ 3x+7y+11z \\ 4x+8y+12z \end{bmatrix}$$

4x3 3x1 4x1

$$A_{(m \times n) \text{ matrix}} B_{(n \times p) \text{ matrix}} = C_{(m \times p) \text{ matrix}}, \text{ where } C_{ik} = \sum_{j=1}^n a_{ij} b_{jk}, i=1, \dots, m, k=1, \dots, p$$

$$\begin{bmatrix} 1 & 2 & 1 \\ 2 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 1 & 1 \\ 2 & 0 \end{bmatrix} = \begin{bmatrix} 4 & 3 \\ 2 & 2 \end{bmatrix}$$

$$\cdot AB \neq BA$$

• matrix addition: elementwise addition

• Matrix Multiplication :

$$(m \times p) AB = \left[A \begin{bmatrix} b_{11} \\ \vdots \\ b_{1n} \end{bmatrix}, \dots, A \begin{bmatrix} b_{p1} \\ \vdots \\ b_{pn} \end{bmatrix} \right]$$

Existence & Uniqueness of Solutions

\Rightarrow homogeneous - $\vec{b} = \vec{0}$ in $A\vec{x} = \vec{b}$, that is, $A\vec{x} = \vec{0}$

• trivial solution: $\vec{x} = \vec{0}$

• nontrivial solution : solutions where not all entries are zero

> inhomogeneous - $\vec{B} \neq \vec{0}$ in $A\vec{x} = \vec{b}$

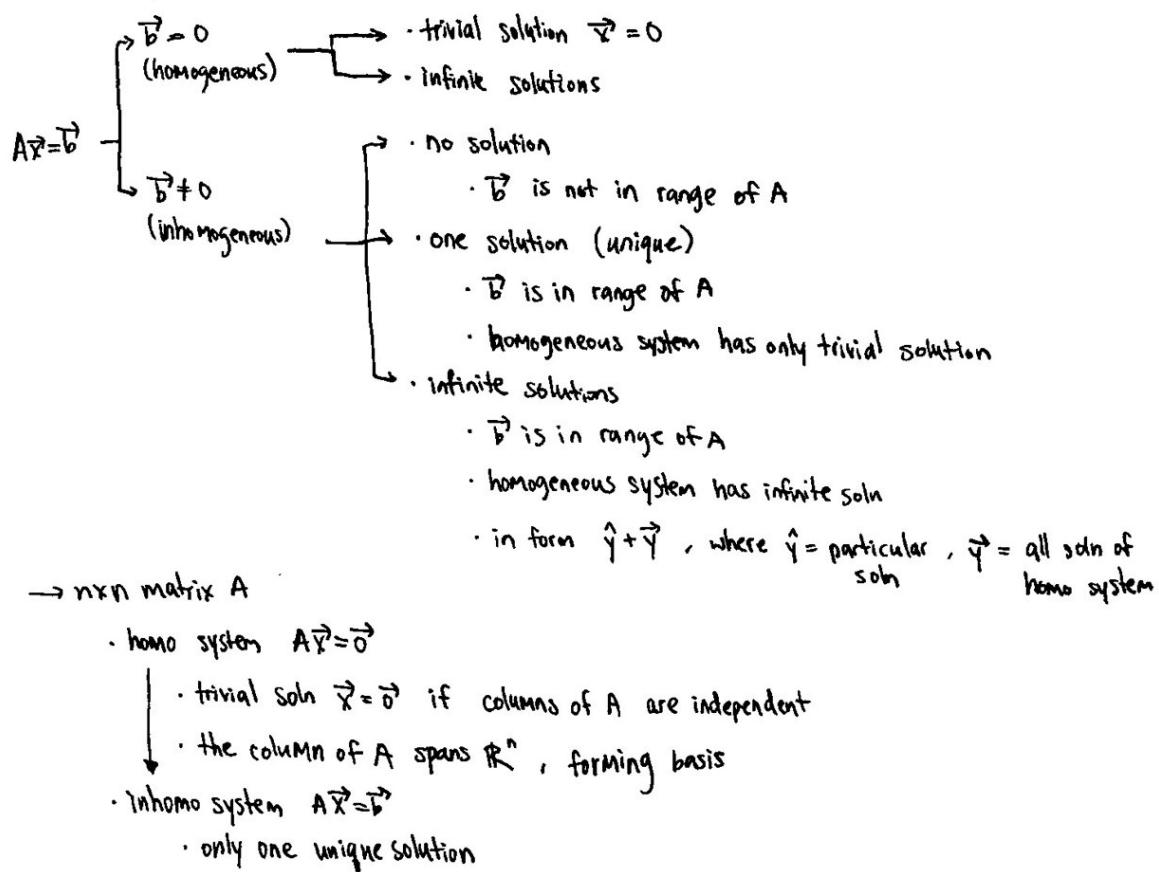
$\Rightarrow \text{range}(A) = \{ A\vec{x} : \vec{x} \in \mathbb{R}^n \} \subset \mathbb{R}^m$ (the set of all things one can get from applying A)

- all LC of columns of A

• Span of columns of A

$$\Rightarrow \text{rank}(A) = \dim(\text{range}(A))$$

<< Existence & Uniqueness of Solutions



<< Inverse of Matrix

> identity matrix I - 1 on diagonal, 0 elsewhere

$$\cdot AI = IA = A$$

$$\cdot I\vec{x} = \vec{x}$$

> invertible - for A , if there exist a matrix B such that $BA = I$ > inverse - the matrix B satisfies $BA = I$, denoted by A^{-1}

• inverse is unique.

• If A is invertible, $A\vec{x} = \vec{b}$ has unique soln, columns of A linearly independent

$$\cdot A^{-1}(A\vec{x}) = (A^{-1}A)\vec{x} = I\vec{x} = \vec{x} = A^{-1}\vec{b}$$

> nonsingular - invertible matrix

> singular - not invertible

→ Conclusion of $n \times n$ matrix

- equivalent
 - A is invertible (nonsingular) $\quad \cdot \det(A) \neq 0$
 - columns of A are linearly independent
 - $A\vec{x} = \vec{0}$ has only trivial soln $\vec{x} = \vec{0}$
 - $A\vec{x} = \vec{b}$ has unique soln
 - rows of A linearly independent

<< Transpose

- > transpose of A - matrix whose (i,j) entry is the (j,i) entry of A , denoted A^T
- $(AB)^T = B^T A^T$
- proves that row of A are linear independent if $n \times n$ A is invertible.

<< Determinant

$$\cdot \det(A) = \det \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} = a_{11}a_{22} - a_{12}a_{21}$$

$$\cdot \det(A) = \sum_{j=1}^n (-1)^{i+j} a_{ij} \det(A_{ij}) = \sum_{i=1}^n (-1)^{i+j} a_{ij} \det(A_{ij}) \quad \text{for } (n-1) \times (n-1) A_{ij} \text{matrix by deleting row } i \text{ column } j$$

↑
→ Laplace expansion

→ Cramer's Rule

- Solving $n \times n$ nonsingular A for $A\vec{x} = \vec{b}$
- $A \leftarrow^j \vec{b}$ - matrix that replaces j th column by vector \vec{b} .
- j th entry x_j of solution \vec{x} is

$$x_j = \frac{\det(A \leftarrow^j \vec{b})}{\det(A)} \quad j=1, \dots, n$$

- $\pm m \times 2^E$, $1 < m < 2$

> single precision (32 bits): 1 bit sign, 8 bit exp, 23 bit significand

→ hidden bit representation, not store the 1 before binary pt, since it's always 1.

> machine precision ϵ - gap between 1 and the next larger fltng pt #.

$$\cdot \text{single precision: } (1 + 2^{-23}) - 1 = 2^{-23} \approx 1.2 \times 10^{-7}$$

> double precision (64 bits): 1 sign, 11 exp, 52 significand

$$\cdot \epsilon = (1 + 2^{-52}) - 1 = 2^{-52} \approx 2.2 \times 10^{-16}$$

→ toy system

$$1. b_1 b_2 \quad \exp(0, 1, -1)$$

$$1.00 \times 2^0 = 1 \quad 1.01 \times 2^0 = \frac{5}{4}$$

$$1.00 \times 2^1 = 2$$

$$1.00 \times 2^{-1} = \frac{1}{2}$$

$$1.10 \times 2^0 = \frac{3}{2}$$

$$1.10 \times 2^1 = 3$$

$$1.11 \times 2^0 = \frac{7}{4}$$

$$1.11 \times 2^1 = \frac{7}{2}$$

$$1.11 \times 2^{-1} = \frac{7}{8}$$

$$\cdot \epsilon = (1 + 0.01) - 1 = 0.01 = 2^{-2} = \frac{1}{4}$$

• abs gap between # becomes large as it moves away from origin

• relative gap still ok.

• gap between smallest positive # & 0 >> that of smallest & next smallest pos. #

$$|1.00 \times 2^{-1} - 1.01 \times 2^{-1}|$$

$$|1.00 \times 2^{-1} - 0| = 2^{-1} \quad |1.00 \times 2^{-1} - 1.01 \times 2^{-1}| = \underbrace{0.01}_{6} \times 2^{-1} = \underbrace{2^{-2} \times 2^{-1}}_{6} = 2^{-3}$$

• the smallest pos. # in any system is 1×2^{-E} , NOTE: gap between two smallest # is

• the next smallest pos. # is $(1 + \epsilon) \times 2^{-E}$ the machine precision $\epsilon \times 2^{-E}$

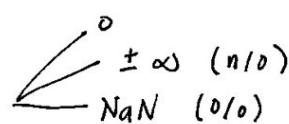
• gap between smallest & next smallest is ϵ times that of 0 & smallest.

IEEE Floating Point Arithmetic

• consistent representation of floating pt #

• correctly rounded arithmetic

• consistent & sensible treatment of exceptions



$\text{EXP}(b_0)$	$\text{EXP}(b_1)$	#
0	0 ... 00	$\pm (0.b_1 \dots b_{52}) \times 2^{-1022}$
1	0 ... 01	$\pm (1.b_1 \dots b_{52}) \times 2^{-1022}$
2	0 ... 10	$\pm (1.b_1 \dots b_{52}) \times 2^{-1021}$
1023	0111 ... 11	$\pm (1.b_1 \dots b_{52}) \times 2^0$
2046	111 ... 10	$\pm (1.b_1 \dots b_{52}) \times 2^{1023}$
2047	111 ... 11	$\pm \infty$ if $b_1 = b_{52} = \dots = 0$ NaN if otherwise

- comment
 - 0 / subnormal #
 - floating pt #
 - actual exp + 1023
 - Exception
- . Smallest + #: 2^{-1022}
 . next smallest: $(1 + 2^{-52}) \times 2^{-1022}$
 . gap: $2^{-52} \times 2^{-1022}$

Rounding

- Round down - largest fp# $\leq x$
- Round up - smallest fp# $\geq x$
- Round to \emptyset - { round down if $x > 0$
round up if $x < 0$
- Round to nearest

→ Absolute Rounding Error

$$> \text{abs error} = |\text{round}(x) - x|$$

$$\begin{aligned} &\cdot \text{double precision : abs error } \leq 2^{-52} \times 2^{+E} && \text{any rounding mode } \leq \epsilon \times 2^{+E} \\ &\quad \text{abs error } \leq 2^{-53} \times 2^{+E} && \text{round to nearest } \leq \frac{\epsilon}{2} \times 2^{+E} \end{aligned}$$

→ Relative Rounding Error

$$> \text{rel error} = \frac{|\text{round}(x) - x|}{|x|} = \frac{\text{abs error}}{|x|} \leq \frac{\epsilon + 2^E}{m \times 2^E} \leq \frac{\epsilon}{m} \leq \epsilon$$

$$\cdot \text{any rounding mode } \leq \epsilon \leq \epsilon$$

$$\cdot \text{round to nearest } \leq \frac{\epsilon}{2}$$

$$\cdot \text{round}(x) = x(1+\delta) \quad , \text{ where } |\delta| < \epsilon \text{ for any rounding} \\ |\delta| < \frac{\epsilon}{2} \text{ for round to nearest}$$

→ Floating Point Operations

- $a \oplus b = \text{round}(a+b) = (a+b)(1+\delta_1)$
- $a \ominus b = \text{round}(a-b) = (a-b)(1+\delta_2)$
- $a \otimes b = \text{round}(ab) = (ab)(1+\delta_3)$
- $a \oslash b = \text{round}(a/b) = (a/b)(1+\delta_4)$

$$\left. \begin{array}{l} |\delta_i| < \epsilon \text{ for any rounding} \\ |\delta_i| < \frac{\epsilon}{2} \text{ for round to nearest,} \\ i=1, \dots, 4 \end{array} \right\}$$

<< Errors in Scientific Computation

- physical prob $\xrightarrow{\text{approx.}}$ math model
 - math model $\xrightarrow{\text{approx.}}$ numeric model
 - numeric model input data measurement error
 - computer rounding error
- > abs error = $|\hat{y} - y|$
- > rel error = $\frac{|\hat{y} - y|}{|y|}$, y is real value, \hat{y} is computed value

<< Conditioning of Problem

> conditioning - how sensitive the answer is to small changes in input

$$> \text{abs condition number } C(x) : |\hat{y} - y| = C(x) |\hat{x} - x|$$

$$> \text{rel condition number } K(x) : \left| \frac{\hat{y} - y}{y} \right| = K(x) \left| \frac{\hat{x} - x}{x} \right|$$

$$\cdot C(x) = \left| \frac{\hat{y} - y}{\hat{x} - x} \right| = \left| \frac{f(\hat{x}) - f(x)}{\hat{x} - x} \right| = |f'(x)|$$

$$\cdot K(x) = \left| \frac{\hat{y} - y}{\hat{x} - x} \right| \left| \frac{x}{\hat{x} - x} \right| = \left| \frac{\hat{y} - y}{\hat{x} - x} \cdot \frac{x}{\hat{y} - y} \right| = \left| \frac{x f'(\hat{x})}{f(x)} \right|$$

<< Stability of Algorithm

> stable - algorithm that achieves level of accuracy defined by the conditioning of prob

> unstable - algorithm that gets unnecessarily inaccurate results due to roundoff

> forward error - output error

> backward error - input error

<< Matrix Multiplication

$$\cdot A = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \dots & a_{mn} \end{bmatrix}$$

$$\cdot \vec{b} = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix}$$

$$\cdot A^T = \begin{bmatrix} a_{11} & \dots & a_{m1} \\ \vdots & \ddots & \vdots \\ a_{1n} & \dots & a_{nn} \end{bmatrix}$$

$$\cdot A\vec{b} = \begin{bmatrix} a_{11}b_{11} + \dots + a_{1n}b_{n1} \\ \vdots \\ a_{m1}b_{11} + \dots + a_{mn}b_{n1} \end{bmatrix}$$

> Symmetric ~ $A = A^T$

<< LU Decomposition

$$> A \equiv LU$$

$$L_2 L_1 A = U$$

$$\cdot A\vec{x} = \vec{b}$$

$$A = (L_2 L_1)^{-1} U$$

$$LU\vec{x} = \vec{b}$$

$$L = (L_1 L_2)^{-1}$$

$$> U\vec{x} = \vec{y}$$

$$\cdot L\vec{y} = \vec{b}$$

<< Gaussian Elimination

- If A is $n \times n$ non-singular, for n -vector \vec{b} , $A\vec{x} = \vec{b}$ has unique soln \vec{x}
- replacement - add to one row a multiple of another row
- interchange - interchange two rows
- scaling - Multiply all entries by a nonzero constant
- after Gaussian elimination \rightarrow upper triangular, solve by back sub.

→ Gaussian Elimination as Factorization

$$L_2 L_1 A = U \quad U: \text{upper } \Delta \text{ matrix}$$

$$\text{let } L = (L_2 L_1)^{-1} \quad L: \text{lower } \Delta \text{ matrix}$$

$$A = LU$$

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

$$LU\vec{x} = \vec{b}$$

$$\text{let } U\vec{x} = \vec{y}$$

$$L\vec{y} = \vec{b}$$

→ Algorithm of Gaussian Elimination w/o partial pivoting

for $j = 1 : n-1$ % columns (for the pivots)

for $i = j+1 : n$ % rows below j

$$\text{mult} = A(i,j) / A(j,j);$$

$\% A(1,:) = A(1,:)$ - mult * $A(j,:)$; // more than necessary

$$A(1,:)=A(1,:)-\text{mult} * A(j,:);$$

$$b(i) = b(i) - \text{mult} * b(j);$$

end

and

• may fail - the (j,j) entry must be non-zero

<< Operation Counts

• $\oplus \ominus \otimes \odot$

• incrementing int loop index not counted

• Gaussian elimination operating on all elements

$$\sum_{i=1}^{n-1} \sum_{j=j+1}^n (2n+3) = (2n+3) \frac{n(n-1)}{2} > n^3 + O(n^2)$$

• Gaussian elimination operating efficiently on nonzero elements

$$\begin{aligned} \sum_{j=1}^{n-1} \sum_{i=j+1}^n [2(n-j)+5] &= 2 \sum_{k=1}^{n-1} k^2 + 5 \sum_{k=1}^{n-1} k = 2 \frac{(n-1)n(2n-1)}{6} + 5 \frac{(n-1)n}{2} \\ &= \frac{2}{3}n^3 + O(n^2) \end{aligned}$$

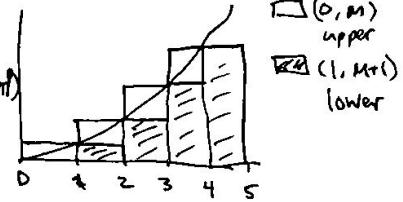
<< Operation Count (Cont.)

→ order equation

$$\sum_{n=1}^m n^p = \int_0^M x^p dx \leq 1^p + 2^p + \dots + m^p \leq \int_1^{M+1} x^p dx \quad (M \in \mathbb{Z}, p \in \mathbb{Z}^+)$$

upper bound
 |
 $= \frac{m^{p+1}}{p+1}$
 |
 $\frac{(m+1)^{p+1}-1}{p+1} = \frac{m^{p+1}}{p+1} + O(m^p)$

$$\sum_{n=1}^m n^p = \frac{m^{p+1}}{p+1} + O(m^p)$$



<< LU Factorization

- use LU to solve systems w/ same coeff. matrix but diff. right hand side vector \vec{b}
- w/o extra storage

→ LU Factorization w/o pivoting

```

for j=1:n-1 % columns
    for i = j+1:n % rows
        mult = A(i,j) / A(j,j)
        A(i, j+1:n) = A(i, j+1:n) - mult * A(j, j+1:n); % upper Δ matrix U
        A(i,j) = mult; % lower Δ matrix L
    end
end

```

→ Solving $L\vec{y} = \vec{b}$ for \vec{y}

$$\begin{bmatrix} l_{11} & 0 & \dots & 0 \\ l_{21} & l_{22} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ l_{n1} & l_{n2} & \dots & l_{nn} \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix} \Rightarrow \begin{aligned} y_1 &= b_1 / l_{11} && (\text{fwd solve}) \\ y_2 &= (b_2 - l_{21}y_1) / l_{22} \\ y_i &= (b_i - \sum_{j=1}^{i-1} l_{ij}y_j) / l_{ii} \end{aligned}$$

→ MATLAB code : L has unit diagonal & lower Δ .

```

function y = lsolve(L,b)
    n = length(b);
    for i = 1:n
        y(i) = b(i);
        for j = 1:i-1
            y(i) = y(i) - L(i,j) * y(j);
        end
    end
end

```

no division by diagonal term since have unit diagonal

→ Operation Count

$$\sum_{i=1}^n \sum_{j=1}^{i-1} 2 = n^2 + O(n)$$

↳ Partial Pivoting

- avoids (ij) entry being zero
- > pivoting - process of interchanging rows
- > pivot element - nonzero entry A_{kj} ($k > j$)
- If no nonzero entries or below main diagonal in column j , matrix is singular
 - no soln
 - infinite soln

- > partial pivoting: search for largest entry in the column in abs value & uses as pivot element
- multiplier $c = 1$
 - $O(n^2)$ finding largest entry for each column
 - additional data movement cost

Theorem Every $n \times n$ nonsingular matrix A can be factored in the form $A = PLU$, where P is a permutation matrix, L is a unit lower Δ matrix, U is ~~upper~~ an upper Δ matrix

→ MATLAB Code

```

for j = 1 : n-1 % columns
    [pivot, k] = max(abs(A(j:n, j))); % largest value as pivot
    if pivot == 0 % pivot = 0 → singular
        disp("Matrix is singular")
        break
    end
    temp = A(j, :);
    A(j, :) = A(k+j-1, :); % exchange rows if A & b
    A(k+j-1, :) = temp;
    b(j) = b(k+j-1); % b(k+j-1) = b(j)
    for i = j+1:n % make A(i,j)=0
        mult = A(i, j) / A(j, j);
        A(i, j:n) = A(i, j:n) - mult * A(j, j:n);
        b(i) = b(i) - mult * b(j);
    end
end

```

23 Apr 2020

Banded Matrix

AMATH 352

<< Banded Matrix

- > symmetric matrix - if A satisfies $A = A^T$
- > positive definite - if symmetric matrix has all eigenvalues positive
 - symmetric positive definite matrix does not need pivot
- > Cholesky decomposition - $A = LL^T$
- > Strict Diagonal Dominance - $|a_{ii}| > \sum_{j \neq i} |a_{ij}|$ for all i , the abs of diagonal entry is greater than the sum of abs of off diagonal entries
- > Banded - $a_{ij} = 0$ if $|i-j| > m$, where $m \ll n$ is the half bandwidth
 - half bandwidth m
 - full bandwidth $2m+1$
 - e.g. tridiagonal matrix - diagonal & one above & below, 0 elsewhere.

$$\begin{array}{l} \text{tridiagonal matrix} \\ \text{in Gaussian elimination:} \\ \text{3(n-1) ops} \end{array} \left[\begin{array}{cccc} a_1 & b_1 & 0 & 0 \\ b_1 & a_2 & b_2 & 0 \\ 0 & b_2 & a_3 & b_3 \\ 0 & 0 & b_3 & a_4 & b_4 \\ 0 & 0 & 0 & b_4 & a_5 \end{array} \right] \sim \left[\begin{array}{ccccc} a_1 & b_1 & 0 & 0 & 0 \\ 0 & \tilde{a}_2 & b_2 & 0 & 0 \\ 0 & 0 & \tilde{a}_3 & b_3 & 0 \\ 0 & 0 & 0 & \tilde{a}_4 & b_4 \\ 0 & 0 & 0 & 0 & \tilde{a}_5 \end{array} \right]$$

$$\tilde{a}_2 = a_2 - \frac{b_1}{a_1} b_1$$

(3 ops) each ~~row~~ row 2

- GE of banded matrix is $2m^2n$

<< Implementation Considerations of High Performance

- Gaussian elimination { fetch data: $3(n-j+1)$
arithmetic: $2(n-j+1)$

BLAS	BLAS 1 - vector operations, least efficient	$\vec{y} \leftarrow \alpha \vec{x} + \vec{y}$
	• $3n$ memory access • $2n$ floating pt operation	} ratio: $\frac{2}{3} \frac{\text{flop}}{\text{ma}} \leftarrow \text{memory access}$
	BLAS 2 - matrix vector operation, intermediate efficiency, $\vec{y} \leftarrow A\vec{x} + \vec{y}$	
	• $n^2 + 3n$ memory access • $2n^2$ floating pt operation	} ratio: $2 \frac{\text{flop}}{\text{ma}}$
	BLAS 3 - matrix-matrix operation, highest efficiency, $C \leftarrow AB + C$	
	• $4n^2$ memory access • $2n^2$ floating pt op.	} ratio: $\frac{n}{2} \frac{\text{flop}}{\text{ma}}$

<< Cholesky Algorithm

- For A that's symmetric positive definite (SPD)

$$\begin{aligned} A &= LU = L\tilde{U}^T \\ &= LD\tilde{U} \\ &= LDL^T \\ &= \underbrace{LD^{\frac{1}{2}}}_{L} \underbrace{D^{\frac{1}{2}}}_{D} \underbrace{L^T}_{\tilde{U}} \\ &= L \tilde{U}^T \end{aligned}$$

where
 L is lower Δ matrix with 1 on diagonal
 U is upper Δ matrix

L is lower Δ matrix with nonzeros on diagonal
 D is diagonal matrix of U
 \tilde{U} is upper Δ matrix with 1 on diagonal
 $\cdot \tilde{U} = L^T$
 $\cdot DU = \tilde{U}$

[Ex] $A = \begin{bmatrix} 3 & 1 \\ 1 & 5 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ \frac{1}{3} & 1 \end{bmatrix} \begin{bmatrix} 3 & 1 \\ 0 & \frac{14}{3} \end{bmatrix} = \begin{bmatrix} 3 & 0 \\ \frac{1}{3} & \frac{14}{3} \end{bmatrix} \begin{bmatrix} 1 & \frac{1}{3} \\ 0 & \frac{14}{3} \end{bmatrix}$

$$\begin{bmatrix} 3 & 1 \\ 1 & 5 \end{bmatrix} \sim \begin{bmatrix} 3 & 1 \\ 0 & \frac{14}{3} \end{bmatrix}$$

$A \uparrow \quad U \uparrow$

$$L_i = \begin{bmatrix} 1 & 0 \\ \frac{1}{3} & 1 \end{bmatrix} \quad L_i^{-1} = \begin{bmatrix} 1 & 0 \\ \frac{1}{3} & 1 \end{bmatrix}$$

$L \uparrow$

$$D = \begin{bmatrix} 3 & 0 \\ 0 & \frac{14}{3} \end{bmatrix} \quad \tilde{U} = L^T = \begin{bmatrix} 1 & \frac{1}{3} \\ 0 & 1 \end{bmatrix}$$

$$A = LD\tilde{U} = LDL^T = \begin{bmatrix} 1 & 0 \\ \frac{1}{3} & 1 \end{bmatrix} \begin{bmatrix} 3 & 0 \\ 0 & \frac{14}{3} \end{bmatrix} \begin{bmatrix} 1 & \frac{1}{3} \\ 0 & 1 \end{bmatrix}$$

$$= (L D^{\frac{1}{2}}) (D^{\frac{1}{2}} L^T) = \begin{bmatrix} 1 & 0 \\ \frac{1}{3} & 1 \end{bmatrix} \begin{bmatrix} \sqrt{3} & 0 \\ 0 & \sqrt{\frac{14}{3}} \end{bmatrix} \begin{bmatrix} \sqrt{3} & 0 \\ 0 & \sqrt{\frac{14}{3}} \end{bmatrix} \begin{bmatrix} 1 & \frac{1}{3} \\ 0 & 1 \end{bmatrix}$$

$$= L \tilde{U}^T = \begin{bmatrix} \sqrt{3} & 0 \\ \frac{\sqrt{14}}{3} & \sqrt{\frac{14}{3}} \end{bmatrix} \begin{bmatrix} \sqrt{3} & \frac{1}{3} \\ 0 & \sqrt{\frac{14}{3}} \end{bmatrix}$$

LU decomposition
 • keep track of mult in L
 • take inverse of L_i as

Cholesky Algorithm
 • find diagonal matrix
 • calculate $\tilde{U} = L^T$
 • follow formula
 $L D^{\frac{1}{2}} D^{\frac{1}{2}} L^T$

<< Methods for Solving $A\vec{x} = \vec{b}$

→ Gaussian Elimination (& LU Decomposition)

- PLU factorization: $\frac{2}{3}n^3$ flops
- solve system: $2n^2$ flops

→ Inversion of Matrix

- append identity matrix to right of A and reduce A to identity matrix,
 the original identity matrix becomes A^{-1} : $[A | I] \sim [I | A^{-1}]$
- matrix inversion: $\frac{8}{3}n^3$ flops ← reduce A to upper Δ , apply multiplication to \vec{e}_1, \vec{e}_2
- solve system: $2n^2$ flops

$$\frac{2}{3}n^3 + n(2n^2) = \frac{8}{3}n^3$$

<< Methods of Solving $A\vec{x} = \vec{b}$ (Cont.)

→ Cramer's Rule

- good for small, not for big

- j th component of solution \vec{x} : replace column j of A by \vec{b} , calc $\frac{\det(A_j(\vec{b}))}{\det(A)}$

- determinant computation: $> n!$

→ Summary

Method	Operation Count
Gaussian Elimination	$\frac{2}{3}n^3$
Inversion of Matrix	$\frac{4}{3}n^3$
Cramer's Rule	$> n!$

<< Conditioning of Linear Systems

→ Norms

> A norm for vectors is a function $\|\cdot\|$ satisfying

(i) $\|\vec{v}\| \geq 0$, where $\|\vec{v}\|=0$ if and only if $\vec{v}=0$.

(ii) $\|\alpha\vec{v}\| = |\alpha| \|\vec{v}\|$ for any scalar α

(iii) $\|\vec{v} + \vec{w}\| \leq \|\vec{v}\| + \|\vec{w}\|$ (triangular inequality)

- . Types of Norms
- 2-norm (Euclidean norm) - $\|\vec{v}\|_2 = \sqrt{\sum_{i=1}^n |v_i|^2}$; $\|\vec{v}\|_2 = \langle \vec{v}, \vec{v} \rangle^{\frac{1}{2}}$, where $\langle \vec{x}, \vec{y} \rangle = \sum_{i=1}^n x_i y_i$ (2-norm from inner product)
 - ∞ -norm - $\|\vec{v}\|_\infty = \max_{i=1, \dots, n} |v_i|$
 - 1-norm - $\|\vec{v}\|_1 = \sum_{i=1}^n |v_i|$
 - p -norm $(p \geq 1)$ - $\|\vec{v}\|_p = \left(\sum_{i=1}^n |v_i|^p \right)^{\frac{1}{p}}$

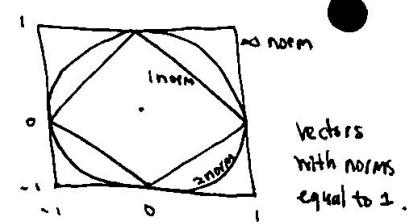
> A matrix norm is function $\|\cdot\|$ for $m \times n$ matrix A, B

(i) $\|A\| \geq 0$, where $\|A\|=0$ if and only if $A=0$.

(ii) $\|\alpha A\| = |\alpha| \|A\|$ for any scalar α

(iii) $\|A+B\| \leq \|A\| + \|B\|$ (triangular inequality)

(iv) $\|AB\| \leq \|A\| \cdot \|B\|$ (submultiplicative)
($m \times n$) ($n \times p$)



<< Norms

→ Matrix Norm

induced matrix norm: $\|A\| = \max_{\|\vec{v}\|=1} \|A\vec{v}\| = \max_{\|\vec{v}\|\neq 0} \frac{\|A\vec{v}\|}{\|\vec{v}\|} \sim \|A \frac{\vec{v}}{\|\vec{v}\|}\|$

$\cdot \|A\| \geq \frac{\|A\vec{v}\|}{\|\vec{v}\|} \rightarrow \|A\vec{v}\| \leq \|A\| \cdot \|\vec{v}\| \quad \begin{matrix} \text{(vector norm} \\ \text{compatible with} \\ \text{matrix norm}) \end{matrix}$

→ Matrix Norm induced by 1-norm

• maximum absolute column sum

$\cdot \|A\|_1 = \max_{j=1,\dots,n} \sum_{i=1}^n |a_{ij}| \quad \text{for } m \times n \text{ matrix induced by } n\text{-vector}$

→ Matrix Norm induced by ∞ -norm

• maximum absolute row sum

$\cdot \|A\|_\infty = \max_{i=1,\dots,n} \sum_{j=1}^n |a_{ij}| \quad \text{for } m \times n \text{ matrix induced by } n\text{-vector}$

→ Matrix Norm induced by 2-norm

$\cdot \|A\|_2 = \sqrt{\text{largest eigenvalue of } A^T A}$

• square root of largest eigenvalue of $A^T A$.

<< Conditioning of Linear Systems

→ Error in \vec{b} • $A\vec{x} = \vec{b}$ is the exact system \Leftrightarrow • $A\vec{x} = \vec{b}$, where \vec{b} is rounded, giving $\hat{\vec{x}}$.

$$\vec{x} - \hat{\vec{x}} = \vec{b} - \hat{\vec{b}}$$

$$\vec{x} - \hat{\vec{x}} = A^{-1}(\vec{b} - \hat{\vec{b}})$$

$$\begin{aligned} \|\vec{x} - \hat{\vec{x}}\| &\leq \|A^{-1}\| \cdot \|\vec{b} - \hat{\vec{b}}\|, \quad C(x) = \|A^{-1}\| \\ \frac{\|\vec{x} - \hat{\vec{x}}\|}{\|\vec{x}\|} &\leq \|A^{-1}\| \cdot \frac{\|\vec{b} - \hat{\vec{b}}\|}{\|\vec{b}\|} \cdot \frac{\|\vec{b}\|}{\|\vec{x}\|} = \|A^{-1}\| \cdot \frac{\|\vec{b} - \hat{\vec{b}}\|}{\|\vec{b}\|} \cdot \frac{\|A\vec{x}\|}{\|\vec{x}\|} \leq \|A^{-1}\| \cdot \|A\| \cdot \frac{\|\vec{b} - \hat{\vec{b}}\|}{\|\vec{b}\|} \end{aligned}$$

$$K(x) = \|A^{-1}\| \|A\|$$

> Condition number of nonsingular matrix - $K(x) = \|A^{-1}\| \|A\|$ → Error in A and \vec{b} • Let $A\vec{x} = \vec{b}$ be exact, $A+E$ be another nonsingular matrix, \vec{b}' another n -vector, $\hat{\vec{x}}$ satisfies $(A+E)\vec{b}'$

then $\frac{\|\vec{x} - \hat{\vec{x}}\|}{\|\vec{x}\|} \leq \|(A+E)^{-1}\| \cdot \|A\| \left(\frac{\|\vec{b}' - \hat{\vec{b}}\|}{\|\vec{b}'\|} + \frac{\|E\|}{\|A\|} \right)$

if $\|E\|$ is small enough so $\|A^{-1}\| \cdot \|E\| < 1$, then

$$\begin{aligned} K(A) &= \|(A+E)^{-1}\| \|A\| \\ &\approx \|A^{-1}\| \|A\| \end{aligned}$$

$$\frac{\|\vec{x} - \hat{\vec{x}}\|}{\|\vec{x}\|} \leq \frac{K(A)}{1 - K(A)\|E\|/\|A\|} \cdot \left(\frac{\|\vec{b}' - \hat{\vec{b}}\|}{\|\vec{b}'\|} + \frac{\|E\|}{\|A\|} \right)$$

↔ Stability of Gaussian Elimination

- without pivoting - unstable

- with pivoting - almost always stable

- backward stable - when solve $A\vec{x} = \vec{b}$, produce \hat{x} satisfies $(A+E)\hat{x} = \vec{b}'$,

$$\text{where } \frac{\|E\|}{\|A\|} = O(\varepsilon), \quad \frac{\|\vec{b}' - \vec{b}\|}{\|\vec{b}\|} = O(\varepsilon)$$

: Suppose $PA = LU$, and GE w/ partial pivoting produces $\tilde{L}\tilde{U}$, then $\tilde{L}\tilde{U} = A+E$,

$$\text{where } \frac{\|E\|}{\|\tilde{L}\|\|\tilde{U}\|} = O(\varepsilon)$$

- there are exceptions, as proposed by Wilkinson.

$$\left[\begin{array}{cccc|c} 1 & 0 & 0 & 1 \\ -1 & 1 & 0 & 1 \\ -1 & -1 & 1 & 1 \end{array} \right] \rightarrow \left[\begin{array}{cccc|c} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 2 \\ 0 & -1 & 1 & 2 \end{array} \right] \rightarrow \left[\begin{array}{cccc|c} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 2 \\ 0 & 0 & 1 & 4 \end{array} \right] \rightarrow \left[\begin{array}{cccc|c} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 2 \\ 0 & 0 & 1 & 4 \end{array} \right]$$

if $n \times n$,
then last entry is 2^{n-1}

<< Overview of Least Squares Problem

- used for overdetermined system
- plotting close fit of data pts
- minimizing square of 2-norm of residual : $\min \| \vec{b} - A\vec{x} \|_2^2 = \sum_{i=1}^m (b_i - \sum_{j=1}^n a_{ij} x_j)^2$
- for $A\vec{x} = \vec{b}$, A is $m \times n$, where $m > n$, b is m -vector, x is n -vector

$$\begin{matrix} A & \vec{x} \\ (m \times n) & (n \times 1) \end{matrix} = \begin{matrix} \vec{b} \\ (m \times 1) \end{matrix}, \text{ find closest } \vec{x}$$

<< Normal Equations: approach I

- differentiate squared error to find min : (a is element of A)

$$\begin{aligned} \frac{\partial}{\partial x_k} \| \vec{b} - A\vec{x} \|_2^2 &= \frac{\partial}{\partial x_k} \sum_{i=1}^m (b_i - \sum_{j=1}^n a_{ij} x_j)^2 \\ &= \sum_{i=1}^m 2(b_i - \sum_{j=1}^n a_{ij} x_j)(-a_{ik}) = 0, \quad k=1, \dots, n \end{aligned}$$

so $\sum_{i=1}^m a_{ik} b_i = \sum_{i=1}^m \left(a_{ik} \left(\sum_{j=1}^n a_{ij} x_j \right) \right) = \sum_{i=1}^m a_{ik} (A\vec{x})_i$

since $\sum_{i=1}^m a_{ik} (A\vec{x})_i = \sum_{i=1}^m (A^T)_{ki} (A\vec{x})_i = \sum_{i=1}^m (A^T)_{ki} b_i$

that is,

$$A^T A \vec{x} = A^T \vec{b} \quad \text{normal equation}$$

Ex1 For $A\vec{x} = \vec{b}$: $\begin{bmatrix} 1 & 1 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 2 \\ 8 \end{bmatrix}, \quad A^T = \begin{bmatrix} 1 & 2 \\ 1 & 1 \end{bmatrix}$ ← overdetermined

by $A^T A \vec{x} = A^T \vec{b}$, $\begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ 8 \end{bmatrix}$

$$\begin{bmatrix} 6 & 2 \\ 2 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 10 \\ 6 \end{bmatrix} \quad \text{← NOT overdetermined}$$

$$\begin{cases} x_1 = \frac{9}{7} \\ x_2 = \frac{8}{7} \end{cases}$$

- drawback - 2-norm condition # of $A^T A$ is square of A

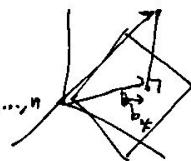
$$\begin{aligned} \| A^T A \|_2^2 &= \sqrt{\max \operatorname{eig}(A^T A)^T (A^T A))} \quad \text{that is, } \lambda_2(A^T A) = \| A \|_2^2 \\ &= \sqrt{\max \operatorname{eig}((A^T A)(A^T A))} \\ &= \sqrt{\max \operatorname{eig}((A^T A)^2)} \\ &= \max \operatorname{eig}(A^T A) \\ &= \| A \|_2^2 \end{aligned}$$

<< QR Decomposition: approach II

- goal - find \vec{x} for $A\vec{x} = \vec{b}_*$, where \vec{b}_* is closest to \vec{b} in 2 norm in range(A)
- equiv goal - minimize $\|\vec{b} - A\vec{x}\|_2$, since $\|\vec{b} - \vec{b}_*\|_2 \leq \|\vec{b} - A\vec{x}\|_2$ for all \vec{x}
- the closest vector to a given vector from a subspace is the orthogonal projection of that vector onto that subspace.

→ Gram-Schmidt Algorithm

$$\vec{g}_j = \frac{\vec{v}_j}{\|\vec{v}_j\|}, \text{ where } \vec{g}_j = \vec{v}_j - \sum_{i=1}^{j-1} \langle \vec{v}_j, \vec{g}_i \rangle \vec{g}_i, \text{ given } \vec{v}_j \text{ for } j=1, \dots, n$$



→ QR Decomposition (Reduced) Solving Least Squares Problem

- Given $A\vec{x} = \vec{b}$, where A is $m \times n$, where $m > n$
- Decompose A to $A = QR$, where Q has orthonormal columns, and R is upper Δ .
- Algorithm

QR Decomposition

1. Apply Gram-Schmidt Algorithm to A. The orthonormal vectors are columns of Q.
2. Find $R = Q^T A$. $\left\{ \begin{array}{l} \text{proof: } A = QR \\ Q^T A = Q^T QR \leftarrow (Q^T Q = I \text{ for orthonormal matrix}) \\ Q^T A = R \end{array} \right.$

Use QR Decomp. Solve system

3. Solve $R\vec{x} = Q^T \vec{b}$ for \vec{x} $\left\{ \begin{array}{l} \text{proof: } A\vec{x} = \vec{b} \\ QR\vec{x} = \vec{b} \\ Q^T QR\vec{x} = Q^T \vec{b} \leftarrow \\ R\vec{x} = Q^T \vec{b} \end{array} \right.$

$$A = Q R$$

$$m \times n \quad m \times n \quad n \times n$$

→ Full QR Decomposition

- first complete reduced QR Factorization
- complete the orthonormal columns to an orthonormal basis for \mathbb{R}^m

$$A = Q \begin{bmatrix} R \\ 0 \end{bmatrix}$$

$$m \times n \quad m \times m \quad m \times n$$

→ MATLAB Command

- $[Q, R] = qr(A, 0)$ for reduced QR
- $[Q, R] = qr(A)$ for full QR

<< Linear Fit

- Given a set of data (x_i, y_i) for $i = 1:m$
- fit linear equation $y = mx + b$ for all data pts

$$y_1 \approx mx_1 + b$$

$$y_2 \approx mx_2 + b$$

 \vdots

$$y_m \approx mx_m + b$$

$$\Rightarrow \begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_m & 1 \end{bmatrix} \begin{bmatrix} m \\ b \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

that minimizes squared error $\sum_{i=1}^m (y_i - (mx_i + b))^2$

- Overdetermined system can be approx. solved by { Normal equation
QR Decomposition }

<< Polynomial Fit

- Given m data pts, can fit $n-1 < m$ degree polynomial
- Find coefficients c_0, c_1, \dots, c_{n-1} that

$$y_i \approx c_0 + c_1 x_i + c_2 x_i^2 + c_3 x_i^3 + \dots + c_{n-1} x_i^{n-1}, \quad i = 1:m$$

so that squared error $\sum_{i=1}^m (y_i - \sum_{j=0}^{n-1} c_j x_i^j)^2$ is minimized.

- Solve by QR / normal eq:

$$\begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^{n-1} \\ 1 & x_2 & x_2^2 & \dots & x_2^{n-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_m & x_m^2 & \dots & x_m^{n-1} \end{bmatrix} \begin{bmatrix} c_0 \\ c_1 \\ \vdots \\ c_{n-1} \end{bmatrix} \approx \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

 $m \times n$
 $n \times 1$
 $m \times 1$

<< Power Method

→ Power Method

- > for computing the eigenvalue of the largest abs value and corresponding normalized eigenvector. ($|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n|$)

Given nonzero vector \vec{w} , set $\vec{y}^{(0)} = \frac{\vec{w}}{\|\vec{w}\|}$.

For $k=1, 2, \dots$,

$$\text{Compute } \tilde{\vec{y}}^{(k)} = A\vec{y}^{(k-1)}$$

$$\text{Set } \lambda^{(k)} = \langle \tilde{\vec{y}}^{(k)}, \vec{y}^{(k-1)} \rangle$$

$$\text{Form } \vec{y}^{(k)} = \frac{\tilde{\vec{y}}^{(k)}}{\|\tilde{\vec{y}}^{(k)}\|}$$

$$[\text{Note that } \lambda^{(k)} = \langle A\vec{y}^{(k-1)}, \vec{y}^{(k-1)} \rangle]$$

- only works if the largest abs(eigenvalue) is strictly greater than others

- rate of convergence depend on $|\frac{\lambda_2}{\lambda_1}|$, the smaller the better

→ Power Method with Shift

- > Computing the eigenvalue farthest from shift s and the corresponding normalized eigenv.

Given nonzero vector \vec{w} , set $\vec{y}^{(0)} = \frac{\vec{w}}{\|\vec{w}\|}$

$$(A-sI)\vec{v} = (\lambda-s)\vec{v}$$

For $k=1, 2, \dots$

$$\text{Compute } \tilde{\vec{y}}^{(k)} = (A-sI)\vec{y}^{(k-1)}$$

$$\text{Set } \lambda^{(k)} = \langle \tilde{\vec{y}}^{(k)}, \vec{y}^{(k-1)} \rangle + s$$

$$\text{Form } \vec{y}^{(k)} = \frac{\tilde{\vec{y}}^{(k)}}{\|\tilde{\vec{y}}^{(k)}\|}$$

<< Deflation

- > deflation - process of modifying initial vector or other vectors generated by the algorithm so that they are orthogonal to the already computed eigenvectors.

- Given: A is symmetric with eigenvalues $\lambda_1, \dots, \lambda_n$ satisfying $|\lambda_1| > |\lambda_2| > |\lambda_3| \geq \dots \geq |\lambda_n|$. λ_1 and \vec{v}_1 found by power method.

$$\text{Let } \vec{w} = \sum_{j=1}^n c_j \vec{v}_j, \quad \hat{\vec{w}} = \vec{w} - \langle \vec{w}, \vec{v}_1 \rangle \vec{v}_1 = \sum_{j=2}^n c_j \vec{v}_j$$

$$A^k \hat{\vec{w}} = \sum_{j=2}^n c_j A^k \vec{v}_j = \sum_{j=2}^n c_j \lambda_j^k \vec{v}_j = \lambda_2^k \left[c_2 \vec{v}_2 + \underbrace{\sum_{j=3}^n c_j \left(\frac{\lambda_2}{\lambda_j} \right)^k \vec{v}_j}_{\text{Converge to } \vec{v}_2} \right]$$

$$\text{Let } \tilde{\vec{y}}^{(k)} = \vec{y}^{(k)} - \langle \vec{y}^{(k)}, \vec{v}_1 \rangle \vec{v}_1 \text{ to prevent going back to } \vec{v}_1$$

• Can use to find all eigenpairs

$$\lim_{k \rightarrow \infty} \underbrace{\text{Converge to } \lambda_2}_{\text{to } \vec{v}_2} = 0.$$

<< Inverse Iteration

- Given A is invertible
- eigenvectors of $(A - sI)^{-1}$ are same as A ; the eigenvalues are $(\lambda_i - s)^{-1} \forall i=1,\dots,n$

$$A\vec{v}_j = \lambda_j \vec{v}_j$$

$$(A - sI)\vec{v}_j = \lambda_j \vec{v}_j - s\vec{v}_j = (\lambda_j - s)\vec{v}_j$$

$$(A - sI)^{-1}(A - sI)\vec{v}_j = (A - sI)^{-1}(\lambda_j - s)\vec{v}_j$$

$$\frac{1}{\lambda_j - s} \vec{v}_j = (A - sI)^{-1}\vec{v}_j$$

> computing the eigenvalue of A that is closest to shift s and the corresponding normalized ~~vector~~ eigenvector.

- Given a shift s and a nonzero vector \vec{w} .

$$\text{Set } \vec{y}^{(0)} = \frac{\vec{w}}{\|\vec{w}\|}$$

For $k=1, 2, \dots$

$$\text{Solve } (A - sI)\vec{y}^{(k)} = \vec{y}^{(k-1)} \text{ for } \vec{y}^{(k)}$$

$$\text{Set } \lambda^{(k)} = \frac{1}{\langle \vec{y}^{(k)}, \vec{y}^{(k-1)} \rangle} + s \quad [\text{Note } \frac{1}{\lambda^{(k)} - s} = \langle (A - sI)^{-1}\vec{y}^{(k)}, \vec{y}^{(k-1)} \rangle]$$

$$\text{Form } \vec{y}^{(k)} = \frac{\vec{y}^{(k)}}{\|\vec{y}^{(k)}\|}$$

- The max abs eigenvalue of $(A - sI)^{-1}$ is the largest $\frac{1}{\lambda_j - s}$, which is the smallest $\lambda_j - s$.
- inverse iteration has faster convergence
- inverse iteration can converge to any eigenvalue. } inverse of small eigenvalue in the middle can be very large

<< Intro to Eigenvalues & Eigenvectors

- > eigenvector - \vec{v} that satisfies $A\vec{v} = \lambda\vec{v}$ } A is $n \times n$
- > eigenvalue - λ that satisfies $A\vec{v} = \lambda\vec{v}$ } \rightarrow can be real / complex : $\lambda \in \mathbb{C}$
- > characteristic polynomial - $\det(A - \lambda I) = 0$
 - polynomial rootfinding can be ill-conditioned
- > eigenspace of λ - $\{\vec{v}, \text{ eigenvectors of } \lambda\}$
 - if \vec{v} is an eigenvector of λ , then $c\vec{v}$ is also. ($c \neq 0$)
 - $A\vec{v} = \lambda\vec{v} \Rightarrow A(c\vec{v}) = \lambda(c\vec{v})$
 - if $\vec{v}_1, \dots, \vec{v}_m$ are eigenv. of λ , then linear comb. of them are also.
 - geometric multiplicity of λ - dimension of eigenspace
 - eigenspace of λ == nullspace $(A - \lambda I)$
 - use $(A - \lambda I)^{\frac{1}{2}} = 0$ to find eigenv.

<< Diagonalization

- > diagonalizable - $n \times n$ matrix A that has n linear independent eigenv.

$$A(\vec{v}_1, \vec{v}_2, \dots, \vec{v}_n) = (\lambda_1 \vec{v}_1, \lambda_2 \vec{v}_2, \dots, \lambda_n \vec{v}_n) = (\vec{v}_1, \vec{v}_2, \dots, \vec{v}_n) \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{bmatrix}$$

$$\boxed{AV = V\Lambda} \quad \left\{ \begin{array}{l} V = (\vec{v}_1, \vec{v}_2, \dots, \vec{v}_n) \\ \Lambda = \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{bmatrix} \end{array} \right.$$

$$\boxed{A = V\Lambda V^{-1}}$$

$$\boxed{\Lambda = V^{-1}AV}$$

<< Similarity Transformation

- > similarity transformation - $A \rightarrow W^{-1}AW$ (W is nonsingular)

Theorem 1 If A is $n \times n$ matrix with n distinct eigenvalues, then A is diagonalizable

Theorem 2 If A is real & symmetric ($A = A^T$), then the eigenvalues of A are real and A is diagonalizable via an orthogonal similar transformation : $A = Q\Lambda Q^T$, where $Q^T = Q^{-1}$, Λ is $[\lambda_1, \dots, \lambda_n]$

- not all $n \times n$ matrix have n independent eigenvectors

- some may have double root at an eigenvalue

Theorem 3Jordan
Canonical
FormEvery $n \times n$ matrix A is similar to one of the form

$$J = \begin{bmatrix} J_1 & & \\ & \ddots & \\ & & J_m \end{bmatrix},$$

where block J_i has form

$$J_i = \begin{bmatrix} \lambda_i & 1 & & \\ 0 & \ddots & & \\ & & \ddots & \\ & & & \lambda_i \end{bmatrix}$$

The number of linear independent eigenvectors is the number of blocks m .The matrix is diagonalizable if and only if $m=n$.The geometric multiplicity of an eigenvalue λ_i is the number of Jordan blocks w/ λ_i .The algebraic multiplicity of λ_i (its degree as a root of characteristic polynomial) is the sum of the orders of all Jordan blocks with eigenvalue λ_i .

A tiny change in matrix can make huge change in Jordan form

Theorem 4Schur
FormEvery square matrix A can be written in the form $A = QTQ^*$, where Q is a unitary matrix, and T is upper triangular.**Theorem 5**

Gershgorin

Let A be an $n \times n$ matrix with entries a_{ij} and let r_i denote the sum of the abs values of off-diag entries in row i : $r_i = \sum_{j=1, j \neq i}^n |a_{ij}|$.Let D_i denote disk in complex plane centered at a_{ii} and has radius r_i : $D_i = \{z \in \mathbb{C} : |z - a_{ii}| \leq r_i\}$ Then all eigenvalues of A lie in the union $\cup_{i=1}^n D_i$ of the Gershgorin disks. If m of these disks are connected and disjoint from the others, then exactly m eigenvalues of A lie in this connected component.

This works similarly for column disks.

29 May 2020

Singular Value Decomposition

AMATH 352

<< Intro to SVD

$$A = U \Sigma V^T$$

 $[U, \Sigma, V] = \text{svd}(A)$ - full SVD $[U, \Sigma, V] = \text{svd}(A, 0)$ - reduced SVD• A - any $m \times n$ matrix• U - matrix of left singular vectors, orthogonal vectors

$AV = U\Sigma$

• Σ - matrix with singular values at diagonal• V - matrix of right singular vectors, orthogonal vectors $\rightarrow m \times n, n > m, \text{full SVD}$

$$\begin{array}{c} A \\ \text{---} \\ m \times n \end{array} = \begin{array}{c} U \\ \text{---} \\ m \times m \end{array} \begin{array}{c} \Sigma \\ \text{---} \\ m \times n \end{array} \begin{array}{c} V^T \\ \text{---} \\ n \times n \end{array}$$

 $\rightarrow m < n, \text{full SVD}$

$$\begin{array}{c} A \\ \text{---} \\ m \times n \end{array} = \begin{array}{c} U \\ \text{---} \\ m \times m \end{array} \begin{array}{c} \Sigma \\ \text{---} \\ m \times n \end{array} \begin{array}{c} V^T \\ \text{---} \\ n \times n \end{array}$$

 $\rightarrow m > n, \text{reduced SVD}$

$$\begin{array}{c} A \\ \text{---} \\ m \times n \end{array} = \begin{array}{c} \hat{U} \\ \text{---} \\ m \times n \end{array} \begin{array}{c} \hat{\Sigma} \\ \text{---} \\ n \times n \end{array} \begin{array}{c} V^T \\ \text{---} \\ n \times n \end{array}$$

 $\rightarrow m < n, \text{reduced SVD}$

$$\begin{array}{c} A \\ \text{---} \\ m \times n \end{array} = \begin{array}{c} U \\ \text{---} \\ m \times m \end{array} \begin{array}{c} \hat{\Sigma} \\ \text{---} \\ m \times m \end{array} \begin{array}{c} \hat{V}^T \\ \text{---} \\ m \times n \end{array}$$

 $\rightarrow \cancel{m < n} m = n$

$$\begin{array}{c} A \\ \text{---} \\ m \times n \end{array} = \begin{array}{c} U \\ \text{---} \\ m \times m \end{array} \begin{array}{c} \Sigma \\ \text{---} \\ m \times n \end{array} \begin{array}{c} V^T \\ \text{---} \\ n \times n \end{array}$$

• full and reduced SVD are the same.

~~ SVD Theorem

Let A be $m \times n$ matrix,

$$A = U \Sigma V^T$$

where U is $m \times n$ matrix with orthonormal columns

Σ is $m \times n$ matrix with nonzeros only on its diag, where entries are real & nonneg

V is $n \times n$ matrix with orthonormal columns

is the full singular value decomposition (SVD)

- > right singular vectors - columns of V
- > singular values - diagonal entries of Σ
- > left singular vectors - columns of U
- Right singular vectors corresponding to distinct singular values are determined up to multiplication by scalars of modulus 1 (± 1).
- Left singular vectors can be determined if right singular vectors are determined.

~~ SVD facts

- Known $A = U \Sigma V^T$, $A^T A = V \Sigma^2 V^T$ is eigendecomposition of $A^T A$ if A is square

$$A^T A = (U \Sigma V^T)^T (U \Sigma V^T) = (V \Sigma^T U^T)(U \Sigma V^T) = V \Sigma^2 V^T$$

Since $V^T V = I$, $\Sigma^T = \Sigma$

- Singular value of $A \rightarrow$ nonnegative $\sqrt{\text{eig}(A^T A)}$
- right singular vector of $A \rightarrow \text{eigvec}(A^T A)$
- $A A^T = U \Sigma^2 U^T$ is eigendecomposition of $A A^T$

$$A A^T = (U \Sigma V^T)(U \Sigma V^T)^T = (U \Sigma V^T)(V \Sigma^T U^T) = U \Sigma^2 U^T$$

Since $V^T V = I$, $\Sigma^T = \Sigma$

- singular value of $A \rightarrow$ nonnegative $\sqrt{\text{eigval}(A A^T)}$
- left singular vector of $A \rightarrow \text{eigvec}(A A^T)$
- $A V = U \Sigma$
- $A^T U = V \Sigma$

<< Algorithm of Singular Value Decomposition

- Given mxn matrix A
- Get $A = U \Sigma V^T$
- 1. Find $A^T A$ and AA^T
- 2. Calculate eigenvalue & eigenvector of $A^T A$
 - singular val of A = $\sqrt{\text{eigval}(A^T A)}$
 - right singular vec of A = $\text{eigvec}(A^T A)$
- 3. Use $AV = U\Sigma$, known A, V, Σ already, solve U.

→ Easy Matrices

1. Diagonal Matrices

- $A = [d_1 \dots d_n]$, $d_j \geq 0$, $\forall j=1, \dots, n$

$$\Sigma = A, U = V = I$$

- $A = [d_1 \dots d_n] [\pm 1 \dots \pm 1]$ if some d_j is negative

$$\Sigma = \text{diag}(|d_1|, \dots, |d_n|), U = I, I = \text{diag}(\pm 1, \dots, \pm 1)$$

2. Real Symmetric Matrices

$$A = Q \Lambda Q^T, Q^T = Q : \Sigma = |\Lambda|$$

3. Orthogonal Matrices

- square matrix with orthonormal columns

- $\Sigma = I$ since $A^T A = AA^T = I$

- $\begin{cases} V = A^T, U = I & (AV = AA^T = I = U\Sigma) \\ V = I, U = A & (AV = A = U\Sigma) \end{cases}$

<< Solving Linear Systems

→ Square Nonsingular Matrix

$$A\vec{x} = \vec{b} \quad \text{and} \quad A = U\Sigma V^T$$

$$U\Sigma V^T \vec{x} = \vec{b}$$

$$\vec{x} = V\Sigma^{-1}U^T \vec{b} \quad \text{since } U^T = U^{-1}, V^T = V^{-1}$$

- not usually used: Gaussian elimination faster

→ Overdetermined System ($m > n$)

$$A\vec{x} \approx \vec{b}$$

$$\vec{x} = V \sum_{j=1}^n \sigma_j \vec{v}_j \vec{u}_j^T \vec{b}$$

- not usually used: QR decomposition faster

<< Matrix Approximation

- If $A = U\Sigma V^T$,

$$\text{then } A = \sum_{j=1}^{\min(m,n)} \sigma_j \vec{u}_j \vec{v}_j^T$$

- Each $\sigma_j \vec{u}_j \vec{v}_j^T$ is rank 1.

- σ_j is (j,j) entry of Σ

- \vec{u}_j is j th column of U

- \vec{v}_j is j th column of V

- Assume that A satisfy above approx, where $\sigma_1 > \sigma_2 \geq \dots \geq \sigma_{\min(m,n)}$.

For any $v < \min(m,n)$, define

$$A_v = \sum_{j=1}^v \sigma_j \vec{u}_j \vec{v}_j^T$$

Then A_v is the closest rank v matrix to A in both 2-norm and the Frobenius norm:

$$\|A - A_v\|_2 = \min_{\substack{B \in \mathbb{R}^{m,n} \\ \text{rank}(B) \leq v}} \|A - B\|_2 = \sigma_{v+1}$$

$$\|A - A_v\|_F = \min_{\substack{B \in \mathbb{R}^{m,n} \\ \text{rank}(B) \leq v}} \|A - B\|_F = \sqrt{\sigma_{v+1}^2 + \dots + \sigma_{\min(m,n)}^2}$$