

Introduction

Tsung-Ming Huang

Department of Mathematics
National Taiwan Normal University

September 8, 2011



Outline

- 1 Vectors and matrices
- 2 Rank and orthogonality
- 3 Eigenvalues and Eigenvectors
- 4 Norms and eigenvalues
- 5 Backward and Forward errors



Vectors and matrices

$A \in \mathbb{F}$ with

$$A = [a_{ij}] = \begin{bmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \cdots & a_{mn} \end{bmatrix}, \quad \mathbb{F} = \mathbb{R} \text{ or } \mathbb{C}.$$

- **Product of matrices:** $C = AB$, where $c_{ij} = \sum_{k=1}^n a_{ik}b_{kj}$, $i = 1, \dots, m$, $j = 1, \dots, p$.
- **Transpose:** $C = A^T$, where $c_{ij} = a_{ji} \in \mathbb{R}$.
- **Conjugate transpose:** $C = A^*$ or $C = A^H$, where $c_{ij} = \bar{a}_{ji} \in \mathbb{C}$.
- **Differentiation:** Let $C = (c_{ij}(t))$. Then $\dot{C} = \frac{d}{dt} C = [\dot{c}_{ij}(t)]$.



- **Outer product** of $x \in \mathbb{F}^m$ and $y \in \mathbb{F}^n$:

$$xy^* = \begin{bmatrix} x_1 \bar{y}_1 & \cdots & x_1 \bar{y}_n \\ \vdots & \ddots & \vdots \\ x_m \bar{y}_1 & \cdots & x_m \bar{y}_n \end{bmatrix} \in \mathbb{F}^{m \times n}.$$

- **Inner product** of $x \in \mathbb{F}^n$ and $y \in \mathbb{F}^n$:

$$\langle y, x \rangle := x^T y = \sum_{i=1}^n x_i y_i = y^T x \in \mathbb{R},$$

$$\langle y, x \rangle := x^* y = \sum_{i=1}^n \bar{x}_i y_i = \overline{y^* x} \in \mathbb{C}.$$



- **Sherman-Morrison Formula:**

Let $A \in \mathbb{R}^{n \times n}$ be nonsingular, $u, v \in \mathbb{R}^n$. If $v^T A^{-1}u \neq -1$, then

$$(A + uv^T)^{-1} = A^{-1} - A^{-1}uv^T A^{-1} / (1 + v^T A^{-1}u). \quad (1)$$

- **Sherman-Morrison-Woodbury Formula:**

Let $A \in \mathbb{R}^{n \times n}$, be nonsingular $U, V \in \mathbb{R}^{n \times k}$. If $(I + V^T A^{-1}U)$ is invertible, then

$$(A + UV^T)^{-1} = A^{-1} - A^{-1}U(I + V^T A^{-1}U)^{-1}V^T A^{-1}.$$

Proof of (1):

$$\begin{aligned} & (A + uv^T)[A^{-1} - A^{-1}uv^T A^{-1} / (1 + v^T A^{-1}u)] \\ = & I + \frac{1}{1 + v^T A^{-1}u} [uv^T A^{-1}(1 + v^T A^{-1}u) - uv^T A^{-1} - uv^T A^{-1}uv^T A^{-1}] \\ = & I + \frac{1}{1 + v^T A^{-1}u} [u(v^T A^{-1}u)v^T A^{-1} - u(v^T A^{-1}u)v^T A^{-1}] \\ = & I. \end{aligned}$$



Example 1

$$\tilde{A} = \begin{bmatrix} 3 & -1 & 1 & 1 & 1 \\ 0 & 1 & 2 & 2 & 2 \\ 0 & 0 & 4 & 1 & 1 \\ 0 & 0 & 0 & 3 & 0 \\ 0 & -1 & 0 & 0 & 3 \end{bmatrix} = A + \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ -1 \end{bmatrix} [0 \ 1 \ 0 \ 0 \ 0],$$

where

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



Rank and orthogonality

Let $A \in \mathbb{R}^{m \times n}$. Then

- $\mathcal{R}(A) = \{y \in \mathbb{R}^m \mid y = Ax \text{ for some } x \in \mathbb{R}^n\} \subseteq \mathbb{R}^m$ is the **range space of A** .
- $\mathcal{N}(A) = \{x \in \mathbb{R}^n \mid Ax = 0\} \subseteq \mathbb{R}^n$ is the **null space of A** .
- $\text{rank}(A) = \dim[\mathcal{R}(A)]$ = the number of maximal linearly independent columns of A .
- $\text{rank}(A) = \text{rank}(A^T)$.
- $\dim(\mathcal{N}(A)) + \text{rank}(A) = n$.
- If $m = n$, then A is nonsingular $\Leftrightarrow \mathcal{N}(A) = \{0\} \Leftrightarrow \text{rank}(A) = n$.



- Let $\{x_1, \dots, x_p\}$ in \mathbb{R}^n . Then $\{x_1, \dots, x_p\}$ is said to be **orthogonal** if

$$x_i^T x_j = 0, \quad \text{for } i \neq j$$

and **orthonormal** if

$$x_i^T x_j = \delta_{ij},$$

where $\delta_{ij} = 0$ if $i \neq j$ and $\delta_{ij} = 1$ if $i = j$.

- $S^\perp = \{y \in \mathbb{R}^m \mid y^T x = 0, \text{ for } x \in S\} =$ orthogonal complement of S .
- $\mathbb{R}^m = \mathcal{R}(A) \oplus \mathcal{N}(A^T)$.
- $\mathbb{R}^n = \mathcal{R}(A^T) \oplus \mathcal{N}(A)$.
- $\mathcal{R}(A)^\perp = \mathcal{N}(A^T)$.
- $\mathcal{R}(A^T)^\perp = \mathcal{N}(A)$.



Special matrices

$$A \in \mathbb{R}^{n \times n}$$

$$\text{Symmetric: } A^T = A$$

$$\text{skew-symmetric: } A^T = -A$$

$$\text{positive definite: } x^T A x > 0, x \neq 0$$

$$\text{non-negative definite: } x^T A x \geq 0$$

$$\text{indefinite: } (x^T A x)(y^T A y) < 0, \text{ for some } x, y$$

$$\text{orthogonal: } A^T A = I_n$$

$$\text{normal: } A^T A = A A^T$$

$$\text{positive: } a_{ij} > 0$$

$$\text{non-negative: } a_{ij} \geq 0.$$

$$A \in \mathbb{C}^{n \times n}$$

$$\text{Hermitian: } A^* = A \quad (A^H = A)$$

$$\text{skew-Hermitian: } A^* = -A$$

$$\text{positive definite: } x^* A x > 0, x \neq 0$$

$$\text{non-negative definite: } x^* A x \geq 0$$

$$\text{indefinite: } (x^* A x)(y^* A y) < 0, \text{ for some } x, y$$

$$\text{unitary: } A^* A = I_n$$

$$\text{normal: } A^* A = A A^*$$



Let $A \in \mathbb{F}^{n \times n}$. Then the matrix A is

- **diagonal** if $a_{ij} = 0$, for $i \neq j$. Denote $D = \text{diag}(d_1, \dots, d_n) \in \mathbf{D}_n$;
- **tridiagonal** if $a_{ij} = 0, |i - j| > 1$;
- **upper bi-diagonal** if $a_{ij} = 0, i > j$ or $j > i + 1$;
- **(strictly) upper triangular** if $a_{ij} = 0, i > j$ ($i \geq j$);
- **upper Hessenberg** if $a_{ij} = 0, i > j + 1$.
(Note: the lower case is the same as above.)

Sparse matrix: n^{1+r} , where $r < 1$ (usually between $0.2 \sim 0.5$). If $n = 1000$, $r = 0.9$, then $n^{1+r} = 501187$.



Eigenvalues and Eigenvectors

Definition 2

Let $A \in \mathbb{C}^{n \times n}$. Then $\lambda \in \mathbb{C}$ is called an **eigenvalue** of A , if there exists $x \neq 0$, $x \in \mathbb{C}^n$ with $Ax = \lambda x$ and x is called an **eigenvector** corresponding to λ .

Notations:

$\sigma(A) :=$ spectrum of $A =$ the set of eigenvalues of A .

$\rho(A) :=$ radius of $A = \max\{|\lambda| : \lambda \in \sigma(A)\}$.

- $\lambda \in \sigma(A) \Leftrightarrow \det(A - \lambda I) = 0$.
- $p(\lambda) = \det(\lambda I - A) =$ characteristic polynomial of A .
- $p(\lambda) = \prod_{i=1}^s (\lambda - \lambda_i)^{m(\lambda_i)}$, $\lambda_i \neq \lambda_j$ (for $i \neq j$) and $\sum_{i=1}^s m(\lambda_i) = n$.
- $m(\lambda_i) =$ algebraic multiplicity of λ_i .
- $n(\lambda_i) = n - \text{rank}(A - \lambda_i I) =$ geometric multiplicity of λ_i .
- $1 \leq n(\lambda_i) \leq m(\lambda_i)$.



If there is some i such that $n(\lambda_i) < m(\lambda_i)$, then A is called degenerated.

The following statements are equivalent:

- (1) There are n linearly independent eigenvectors;
- (2) A is diagonalizable, i.e., there is a nonsingular matrix T such that $T^{-1}AT \in \mathbf{D}_n$;
- (3) For each $\lambda \in \sigma(A)$, it holds $m(\lambda) = n(\lambda)$.

If A is degenerated, then eigenvectors plus principal vectors derive Jordan form.



Theorem 3 (Schur decomposition)

- (1) Let $A \in \mathbb{C}^{n \times n}$. There is a unitary matrix U such that U^*AU is upper triangular.
- (2) Let $A \in \mathbb{R}^{n \times n}$. There is an orthogonal matrix Q such that $Q^T A Q$ is quasi-upper triangular, i.e., an upper triangular matrix possibly with nonzero subdiagonal elements in non-consecutive positions.
- (3) A is normal if and only if there is a unitary U such that $U^*AU = D$ is diagonal.
- (4) A is Hermitian if and only if A is normal and $\sigma(A) \subseteq \mathbb{R}$.
- (5) A is symmetric if and only if there is an orthogonal U such that $U^T A U = D$ is diagonal and $\sigma(A) \subseteq \mathbb{R}$.



Norms and eigenvalues

Let X be a vectorspace over $\mathbb{F} = \mathbb{R}$ or \mathbb{C} .

Definition 4 (Vector norms)

Let N be a real-valued function defined on X ($N : X \rightarrow \mathbb{R}_+$). Then N is a (vector) norm, if

N1: $N(\alpha x) = |\alpha|N(x)$, $\alpha \in \mathbb{F}$, for $x \in X$;

N2: $N(x + y) \leq N(x) + N(y)$, for $x, y \in X$;

N3: $N(x) = 0$ if and only if $x = 0$.

The usual notation is $\|x\| = N(x)$.



Example 5

Let $X = \mathbb{C}^n$, $p \geq 1$. Then $\|x\|_p = (\sum_{i=1}^n |x_i|^p)^{1/p}$ is an l_p -norm.
Especially,

$$\|x\|_1 = \sum_{i=1}^n |x_i| \quad (l_1\text{-norm}),$$

$$\|x\|_2 = \left(\sum_{i=1}^n |x_i|^2 \right)^{1/2} \quad (\text{Euclidean-norm}),$$

$$\|x\|_\infty = \max_{1 \leq i \leq n} |x_i| \quad (\text{maximum-norm}).$$



Lemma 6

$N(x)$ is a continuous function in the components x_1, \dots, x_n of x .

Proof:

$$|N(x) - N(y)| \leq N(x - y) \leq \sum_{j=1}^n |x_j - y_j| N(e_j) \leq \|x - y\|_{\infty} \sum_{j=1}^n N(e_j).$$

□

Theorem 7 (Equivalence of norms)

Let N and M be two norms on \mathbb{C}^n . Then there exist constants $c_1, c_2 > 0$ such that

$$c_1 M(x) \leq N(x) \leq c_2 M(x), \text{ for all } x \in \mathbb{C}^n.$$

▶ Proof of Theorem 7

Remark: Theorem 7 does not hold in infinite dimensional space.



Norms and eigenvalues

Definition 8 (Matrix-norms)

Let $A \in \mathbb{C}^{m \times n}$. A real value function $\|\cdot\| : \mathbb{C}^{m \times n} \rightarrow \mathbb{R}_+$ satisfying

N1: $\|\alpha A\| = |\alpha| \|A\|;$

N2: $\|A + B\| \leq \|A\| + \|B\|;$

N3: $\|A\| = 0$ if and only if $A = 0;$

N4: $\|AB\| \leq \|A\| \|B\|;$

N5: $\|Ax\|_v \leq \|A\| \|x\|_v.$

If $\|\cdot\|$ satisfies N1 to N4, then it is called a matrix norm. In addition, matrix and vector norms are compatible for some $\|\cdot\|_v$ in N5.



Example 9 (Frobenius norm)

$$\text{Let } \|A\|_F = \left\{ \sum_{i,j=1}^n |a_{i,j}|^2 \right\}^{1/2}.$$

$$\begin{aligned} \|AB\|_F &= \left(\sum_{i,j} \left| \sum_k a_{ik} b_{kj} \right|^2 \right)^{1/2} \\ &\leq \left(\sum_{i,j} \left\{ \sum_k |a_{ik}|^2 \right\} \left\{ \sum_k |b_{kj}|^2 \right\} \right)^{1/2} \quad (\text{Cauchy-Schwartz Ineq.}) \\ &= \left(\sum_i \sum_k |a_{ik}|^2 \right)^{1/2} \left(\sum_j \sum_k |b_{kj}|^2 \right)^{1/2} = \|A\|_F \|B\|_F. \end{aligned}$$

This implies that N4 holds.

$$\|Ax\|_2 = \left(\sum_i \left| \sum_j a_{ij} x_j \right|^2 \right)^{1/2} \leq \left\{ \sum_i \left(\sum_j |a_{ij}|^2 \right) \left(\sum_j |x_j|^2 \right) \right\}^{1/2} = \|A\|_F \|x\|_2.$$

This implies N5 holds. Also, N1, N2 and N3 hold obviously. ($\|I\|_F = \sqrt{n}$) \square



Example 10 (Operator norm)

Given a vector norm $\|\cdot\|$. An associated matrix norm is defined by

$$\|A\| = \sup_{x \neq 0} \frac{\|Ax\|}{\|x\|} = \max_{x \neq 0} \frac{\|Ax\|}{\|x\|} = \max_{\|x\|=1} \{\|Ax\|\}.$$

$N5$ holds immediately. On the other hand,

$$\begin{aligned}\|(AB)x\| &= \|A(Bx)\| \leq \|A\| \|Bx\| \\ &\leq \|A\| \|B\| \|x\|\end{aligned}$$

for all $x \neq 0$. This implies that

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

Thus, $N4$ holds. ($\|I\| = 1$). □



Three useful matrix norms:

$$\|A\|_1 = \sup_{x \neq 0} \frac{\|Ax\|_1}{\|x\|_1} = \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{ij}| \quad (3)$$

$$\|A\|_\infty = \sup_{x \neq 0} \frac{\|Ax\|_\infty}{\|x\|_\infty} = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}| \quad (4)$$

$$\|A\|_2 = \sup_{x \neq 0} \frac{\|Ax\|_2}{\|x\|_2} = \sqrt{\rho(A^*A)} \quad (5)$$

▶ Proof of (3)-(5)

Example 11 (Dual norm)

Let $\frac{1}{p} + \frac{1}{q} = 1$. Then $\|\cdot\|_p^* = \|\cdot\|_q$, ($p = \infty, q = 1$). (It concludes from the application of the Hölder inequality, i.e. $|y^*x| \leq \|x\|_p \|y\|_q$.)



Theorem 12

Let $A \in \mathbb{C}^{n \times n}$. Then for any operator norm $\|\cdot\|$, it holds

$$\rho(A) \leq \|A\|.$$

Moreover, for any $\epsilon > 0$, there exists an operator norm $\|\cdot\|_\epsilon$ such that

$$\|\cdot\|_\epsilon \leq \rho(A) + \epsilon.$$

▶ Proof of Theorem 12

Lemma 13

Let U and V are unitary. Then

$$\|UAV\|_F = \|A\|_F, \quad \|UAV\|_2 = \|A\|_2$$

From

$$\begin{aligned} \|UA\|_F &= \sqrt{\|Ua_1\|_2^2 + \cdots + \|Ua_n\|_2^2}, \\ \rho(A^*A) &= \rho(AA^*). \end{aligned}$$



Theorem 14 (Singular Value Decomposition (SVD))

Let $A \in \mathbb{C}^{m \times n}$. Then there exist unitary matrices

$U = [u_1, \dots, u_m] \in \mathbb{C}^{m \times m}$ and $V = [v_1, \dots, v_n] \in \mathbb{C}^{n \times n}$ such that

$$U^*AV = \text{diag}(\sigma_1, \dots, \sigma_p) = \Sigma, \quad (6)$$

where $p = \min\{m, n\}$ and $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$. (Here, σ_i denotes the i -th largest singular value of A).

▶ Proof of Theorem 14

Remark: From (6), we have $\|A\|_2 = \sqrt{\rho(A^*A)} = \sigma_1$, which is the maximal singular value of A , and

$$\|ABC\|_F = \|U\Sigma V^*BC\|_F = \|\Sigma V^*BC\|_F \leq \sigma_1 \|BC\|_F = \|A\|_2 \|BC\|_F.$$

This implies

$$\|ABC\|_F \leq \|A\|_2 \|B\|_F \|C\|_2. \quad (7)$$

In addition, by (2) and (7), we get

$$\|A\|_2 \leq \|A\|_F \leq \sqrt{n} \|A\|_2.$$



Theorem 15

Let $A \in \mathbb{C}^{n \times n}$. The statements are equivalent:

- (1) $\lim_{m \rightarrow \infty} A^m = 0$;
- (2) $\lim_{m \rightarrow \infty} A^m x = 0$ for all x ;
- (3) $\rho(A) < 1$.

Proof:

(1) \Rightarrow (2): Trivial.

(2) \Rightarrow (3): Let $\lambda \in \sigma(A)$, i.e., $Ax = \lambda x$, $x \neq 0$. This implies $A^m x = \lambda^m x \rightarrow 0$, as $\lambda^m \rightarrow 0$. Thus $|\lambda| < 1$, i.e., $\rho(A) < 1$.

(3) \Rightarrow (1): There is a norm $\|\cdot\|$ with $\|A\| < 1$ (by Theorem 12). Therefore, $\|A^m\| \leq \|A\|^m \rightarrow 0$, i.e., $A^m \rightarrow 0$. □



Theorem 16

$$\rho(A) = \lim_{k \rightarrow \infty} \|A^k\|^{1/k}.$$

Proof: Since

$$\rho(A)^k = \rho(A^k) \leq \|A^k\| \Rightarrow \rho(A) \leq \|A^k\|^{1/k},$$

for $k = 1, 2, \dots$. If $\epsilon > 0$, then $\tilde{A} = [\rho(A) + \epsilon]^{-1}A$ has spectral radius < 1 and $\|\tilde{A}^k\| \rightarrow 0$ as $k \rightarrow \infty$. There is an $N = N(\epsilon, A)$ such that $\|\tilde{A}^k\| < 1$ for all $k \geq N$. Thus,

$$\|A^k\| \leq [\rho(A) + \epsilon]^k, \text{ for all } k \geq N$$

or

$$\|A^k\|^{1/k} \leq \rho(A) + \epsilon, \text{ for all } k \geq N.$$

Since $\rho(A) \leq \|A^k\|^{1/k}$, and k, ϵ are arbitrary, $\lim_{k \rightarrow \infty} \|A^k\|^{1/k}$ exists and equals $\rho(A)$.



Theorem 17

Let $A \in \mathbb{C}^{n \times n}$, and $\rho(A) < 1$. Then $(I - A)^{-1}$ exists and

$$(I - A)^{-1} = I + A + A^2 + \dots$$

Proof: Since $\rho(A) < 1$, the eigenvalues of $(I - A)$ are nonzero. Therefore, by Theorem 15, $(I - A)^{-1}$ exists and

$$(I - A)(I + A + A^2 + \dots + A^m) = I - A^{m+1} \rightarrow I.$$

□

Corollary 18

If $\|A\| < 1$, then $(I - A)^{-1}$ exists and

$$\|(I - A)^{-1}\| \leq \frac{1}{1 - \|A\|}.$$

Proof: Since $\rho(A) \leq \|A\| < 1$ (by Theorem 12),

$$\|(I - A)^{-1}\| = \left\| \sum_{i=0}^{\infty} A^i \right\| \leq \sum_{i=0}^{\infty} \|A\|^i = (1 - \|A\|)^{-1}.$$



□

Theorem 19 (without proof)

For $A \in \mathbb{F}^{n \times n}$ the following statements are equivalent:

- (1) There is a multiplicative norm p with $p(A^k) \leq 1, k = 1, 2, \dots$
- (2) For each multiplicative norm p the power $p(A^k)$ are uniformly bounded, i.e., there exists a $M(p) < \infty$ such that $p(A^k) \leq M(p), k = 0, 1, 2, \dots$
- (3) $\rho(A) \leq 1$ and all eigenvalue λ with $|\lambda| = 1$ are not degenerated. (i.e., $m(\lambda) = n(\lambda)$.)

(See Householder: *The theory of matrix*, pp.45-47.)



In the following, we prove some important inequalities of vector norms and matrix norms.

$$1 \leq \frac{\|x\|_p}{\|x\|_q} \leq n^{(q-p)/pq}, \quad (p \leq q). \quad (8)$$

▶ Proof of (8)

$$1 \leq \frac{\|x\|_p}{\|x\|_\infty} \leq n^{\frac{1}{p}}. \quad (9)$$

▶ Proof of (9)

$$\max_{1 \leq j \leq n} \|a_j\|_p \leq \|A\|_p \leq n^{(p-1)/p} \max_{1 \leq j \leq n} \|a_j\|_p, \quad (10)$$

where $A = [a_1, \dots, a_n] \in \mathbb{R}^{m \times n}$.

▶ Proof of (10)



$$\max_{i,j} |a_{ij}| \leq \|A\|_p \leq n^{(p-1)/p} m^{1/p} \max_{i,j} |a_{ij}|, \quad (11)$$

where $A \in \mathbb{R}^{m \times n}$.

Proof of (11): By (9) and (10) immediately. □

$$m^{(1-p)/p} \|A\|_1 \leq \|A\|_p \leq n^{(p-1)/p} \|A\|_1. \quad (12)$$

Proof of (12): By (10) and (8) immediately. □



Hölder inequality:

$$|x^T y| \leq \|x\|_p \|y\|_q, \text{ where } \frac{1}{p} + \frac{1}{q} = 1. \quad (13)$$

Proof of (13): Let $\alpha_i = \frac{x_i}{\|x\|_p}$, $\beta_i = \frac{y_i}{\|y\|_q}$. Then

$$(\alpha_i^p)^{1/p} (\beta_i^q)^{1/q} \leq \frac{1}{p} \alpha_i^p + \frac{1}{q} \beta_i^q. \quad (\text{Jensen Inequality})$$

Since $\|\alpha\|_p = 1$, $\|\beta\|_q = 1$, it follows that

$$\sum_{i=1}^n \alpha_i \beta_i \leq \frac{1}{p} + \frac{1}{q} = 1.$$

Then we have $|x^T y| \leq \|x\|_p \|y\|_q$. □



$$\max\{|x^T y| : \|x\|_p = 1\} = \|y\|_q. \quad (14)$$

Proof of (14): Take $x_i = y_i^{q-1} / \|y\|_q^{q/p}$. Then we have

$$\|x\|_p^p = \frac{\sum |y_i|^q}{\|y\|_q^{q/p}} = \frac{\|y\|_q^q}{\|y\|_q^{q/p}} = 1. \quad (\because (q-1)p = 1)$$

It follows

$$\left| \sum_{i=1}^n x_i^T y_i \right| = \frac{\sum |y_i|^q}{\|y\|_q^{q/p}} = \frac{\|y\|_q^q}{\|y\|_q^{q/p}} = \|y\|_q.$$

□

Remark: $\exists \hat{z}$ with $\|\hat{z}\|_p = 1$ s.t. $\|y\|_q = \hat{z}^T y$. Let $z = \hat{z} / \|y\|_q$. Then we have $\exists z$ s.t. $z^T y = 1$ with $\|z\|_p = \frac{1}{\|y\|_q}$.



$$\|A\|_p = \|A^T\|_q \quad (15)$$

▶ Proof of (15)

$$n^{-\frac{1}{p}} \|A\|_\infty \leq \|A\|_p \leq m^{\frac{1}{p}} \|A\|_\infty. \quad (16)$$

▶ Proof of (16)

$$\|A\|_2 \leq \sqrt{\|A\|_p \|A\|_q}, \quad \left(\frac{1}{p} + \frac{1}{q} = 1\right). \quad (17)$$

▶ Proof of (17)

$$n^{(p-q)/pq} \|A\|_q \leq \|A\|_p \leq m^{(q-p)/pq} \|A\|_q, \quad (18)$$

where $A \in \mathbb{R}^{m \times n}$ and $q \geq p \geq 1$.

▶ Proof of (18)



Backward error and Forward error

Let $x = F(a)$. We define backward and forward errors in Figure 1. In Figure 1, $\hat{x} + \Delta x = F(a + \Delta a)$ is called a mixed forward-backward error, where $|\Delta x| \leq \varepsilon|x|$, $|\Delta a| \leq \eta|a|$.

Definition 20

- (i) An algorithm is **backward stable**, if for all a , it produces a computed \hat{x} with a small backward error, i.e., $\hat{x} = F(a + \Delta a)$ with Δa small.
- (ii) An algorithm is **numerical stable**, if it is stable in the mixed forward-backward error sense, i.e., $\hat{x} + \Delta x = F(a + \Delta a)$ with both Δa and Δx small.
- (iii) If a method which produces answers with forward errors of similar magnitude to those produced by a backward stable method, is called a **forward stable**.



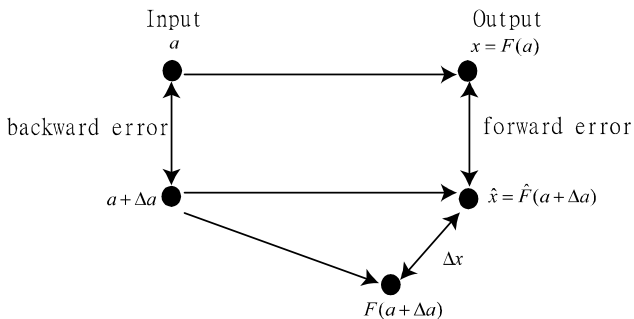


Figure: Relationship between backward and forward errors.

Remark:

- (i) Backward stable \Rightarrow forward stable, no vice versa!
- (ii) Forward error \leq condition number \times backward error



Consider

$$\hat{x} - x = F(a + \Delta a) - F(a) = F'(a)\Delta a + \frac{F''(a + \theta\Delta a)}{2}(\Delta a)^2, \quad \theta \in (0, 1).$$

Then we have

$$\frac{\hat{x} - x}{x} = \left(\frac{aF'(a)}{F(a)} \right) \frac{\Delta a}{a} + O((\Delta a)^2).$$

The quantity $C(a) = \left| \frac{aF'(a)}{F(a)} \right|$ is called the condition number of F . If x or F is a vector, then the condition number is defined in a similar way using norms and it measures the maximum relative change, which is attained for some, but not all Δa .

Backward error: $\left\{ \begin{array}{l} \text{\AA} priori error estimate ! \\ \text{\AA} posteriori error estimate ! \end{array} \right.$



Lemma 21

$$\begin{cases} Ax = b \\ (A + \Delta A)\hat{x} = b + \Delta b \end{cases}$$

with $\|\Delta A\| \leq \delta \|A\|$ and $\|\Delta b\| \leq \delta \|b\|$. If $\delta \kappa(A) = r < 1$ then $A + \Delta A$ is nonsingular and $\frac{\|\hat{x}\|}{\|x\|} \leq \frac{1+r}{1-r}$.

Proof: Since $\|A^{-1}\Delta A\| < \delta \|A^{-1}\| \|A\| = r < 1$, it follows that $A + \Delta A$ is nonsingular. From $(I + A^{-1}\Delta A)\hat{x} = x + A^{-1}\Delta b$, we have

$$\begin{aligned} \|\hat{x}\| &\leq \|(I + A^{-1}\Delta A)^{-1}\| (\|x\| + \delta \|A^{-1}\| \|b\|) \\ &\leq \frac{1}{1-r} (\|x\| + \delta \|A^{-1}\| \|b\|) \\ &= \frac{1}{1-r} \left(\|x\| + r \frac{\|b\|}{\|A\|} \right) \end{aligned}$$

$$\Rightarrow \|\hat{x}\| \leq \frac{1}{1-r} (\|x\| + r \|x\|). \quad (\because \|b\| = \|Ax\| \leq \|A\| \|x\|)$$



Normwise Forward Error Bound

Theorem 22

If the condition of Lemma 21 hold, then

$$\frac{\|x - \hat{x}\|}{\|x\|} \leq \frac{2\delta}{1-r} \kappa(A).$$

Proof: Since $\hat{x} - x = A^{-1}\Delta b - A^{-1}\Delta A\hat{x}$, we have

$$\|\hat{x} - x\| \leq \delta \|A^{-1}\| \|b\| + \delta \|A^{-1}\| \|A\| \|\hat{x}\|.$$

So, by Lemma 21, we have

$$\begin{aligned} \frac{\|\hat{x} - x\|}{\|x\|} &\leq \delta \kappa(A) \frac{\|b\|}{\|A\| \|x\|} + \delta \kappa(A) \frac{\|\hat{x}\|}{\|x\|} \\ &\leq \delta \kappa(A) \left(1 + \frac{1+r}{1-r}\right) = \frac{2\delta}{1-r} \kappa(A). \end{aligned}$$



Componentwise Forward Error Bounds

Theorem 23

Let $Ax = b$ and $(A + \Delta A)\hat{x} = b + \Delta b$. Let $|\Delta A| \leq \delta |A|$ and $|\Delta b| \leq \delta |b|$. If $\delta \kappa_\infty(A) = r < 1$ then $(A + \Delta A)$ is nonsingular and

$$\frac{\|\hat{x} - x\|_\infty}{\|x\|_\infty} \leq \frac{2\delta}{1-r} \| |A^{-1}| |A| \|_\infty.$$

Proof: Since $\|\Delta A\|_\infty \leq \delta \|A\|_\infty$ and $\|\Delta b\|_\infty \leq \delta \|b\|_\infty$, the conditions of Lemma 21 are satisfied in ∞ -norm. Then $A + \Delta A$ is nonsingular and $\frac{\|\hat{x}\|_\infty}{\|x\|_\infty} \leq \frac{1+r}{1-r}$.

Since $\hat{x} - x = A^{-1}\Delta b - A^{-1}\Delta A\hat{x}$, we have

$$\begin{aligned} |\hat{x} - x| &\leq |A^{-1}| |\Delta b| + |A^{-1}| |\Delta A| |\hat{x}| \\ &\leq \delta |A^{-1}| |b| + \delta |A^{-1}| |A| |\hat{x}| \leq \delta |A^{-1}| |A| (|x| + |\hat{x}|). \end{aligned}$$



Taking ∞ -norm, we get

$$\begin{aligned}\|\hat{x} - x\|_{\infty} &\leq \delta \| |A^{-1}| |A| \|_{\infty} \left(\|x\|_{\infty} + \frac{1+r}{1-r} \|x\|_{\infty} \right) \\ &= \frac{2\delta}{1-r} \underbrace{\| |A^{-1}| |A| \|_{\infty}}_{\text{Skeel condition number}}.\end{aligned}$$



Condition Number by First Order Approximation

$$(A + \epsilon F)x(\epsilon) = b + \epsilon f, \quad x(0) = x$$

$$\dot{x}(0) = A^{-1}(f - Fx)$$

$$x(\epsilon) = x + \epsilon \dot{x}(0) + o(\epsilon^2)$$

$$\frac{\|x(\epsilon) - x\|}{\|x\|} \leq \epsilon \|A^{-1}\| \left\{ \frac{\|f\|}{\|x\|} + \|F\| \right\} + o(\epsilon^2)$$

$$\text{Condition number } \kappa(A) := \|A\| \|A^{-1}\|$$

$$\|b\| \leq \|A\| \|x\|,$$

$$\frac{\|x(\epsilon) - x\|}{\|x\|} \leq \kappa(A)(\rho_A + \rho_b) + o(\epsilon^2).$$

$$\rho_A = \epsilon \frac{\|F\|}{\|A\|}, \quad \rho_b = \epsilon \frac{\|f\|}{\|b\|}, \quad \kappa_2(A) = \frac{\sigma_1(A)}{\sigma_n(A)}.$$



Normwise Backward Error Bound

Theorem 24

Let \hat{x} be the computed solution of $Ax = b$. Then the normwise backward error bound

$$\eta(\hat{x}) := \min \{ \epsilon \mid (A + \Delta A)\hat{x} = b + \Delta b, \quad \|\Delta A\| \leq \epsilon \|A\|, \quad \|\Delta b\| \leq \epsilon \|b\| \}$$

is given by

$$\eta(\hat{x}) = \frac{\|r\|}{\|A\| \|\hat{x}\| + \|b\|}, \quad (19)$$

where $r = b - A\hat{x}$ is the residual.



Proof: The right hand side of (19) is a upper bound of $\eta(\hat{x})$. This upper bound is attained for the perturbation (by construction)

$$\Delta A_{\min} = \frac{\|A\| \|\hat{x}\| r z^T}{\|A\| \|\hat{x}\| + \|b\|}, \quad \Delta b_{\min} = -\frac{\|b\|}{\|A\| \|\hat{x}\| + \|b\|} r,$$

where z is the dual vector of \hat{x} , i.e. $z^T \hat{x} = 1$ and $\|z\|_* = \frac{1}{\|\hat{x}\|}$.

Check:

$$\|\Delta A_{\min}\| = \eta(\hat{x}) \|A\|,$$

or

$$\|\Delta A_{\min}\| = \frac{\|A\| \|\hat{x}\| \|r z^T\|}{\|A\| \|\hat{x}\| + \|b\|} = \left(\frac{\|r\|}{\|A\| \|\hat{x}\| + \|b\|} \right) \|A\|,$$

i.e. claim

$$\|r z^T\| = \frac{\|r\|}{\|\hat{x}\|}.$$

Since

$$\|r z^T\| = \max_{\|u\|=1} \left\| (r z^T) u \right\| = \|r\| \max_{\|u\|=1} \left| z^T u \right| = \|r\| \|z\|_* = \|r\| \frac{1}{\|\hat{x}\|},$$

we have done. Similarly, $\|\Delta b_{\min}\| = \eta(\hat{x}) \|b\|$.



Componentwise Backward Error Bound

Theorem 25

The componentwise backward error bound

$$\omega(\hat{x}) := \min \{ \epsilon | (A + \Delta A)\hat{x} = b + \Delta b, \quad |\Delta A| \leq \epsilon |A|, \quad |\Delta b| \leq \epsilon |b| \}$$

is given by

$$\omega(\hat{x}) = \max_i \frac{|r|_i}{(A|\hat{x}| + b)_i}, \quad (20)$$

where $r = b - A\hat{x}$. (note: $\xi/0 = 0$ if $\xi = 0$; $\xi/0 = \infty$ if $\xi \neq 0$.)

Proof: The right hand side of (20) is an upper bound for $\omega(\hat{x})$. This bound is attained for the perturbation

$$\Delta A = D_1 A D_2, \quad \Delta b = -D_1 b,$$

where

$$D_1 = \text{diag}(r_i / (A|\hat{x}| + b)_i) \text{ and } D_2 = \text{diag}(\text{sign}(\hat{x}_i)).$$



Determinants and Nearness to Singularity

$$B_n = \begin{bmatrix} 1 & -1 & \cdots & -1 \\ & 1 & \ddots & \vdots \\ & & 1 & -1 \\ 0 & & & 1 \end{bmatrix}, \quad B_n^{-1} = \begin{bmatrix} 1 & 1 & \cdots & 2^{n-2} \\ & & \ddots & \vdots \\ & & & 1 \\ 0 & & & 1 \end{bmatrix},$$

$$\det(B_n) = 1, \quad \kappa_\infty(B_n) = n2^{n-1}, \quad \sigma_n(B_n) \approx 10^{-8} (n = 30).$$

$$D_n = \begin{bmatrix} 10^{-1} & & 0 \\ & \ddots & \\ 0 & & 10^{-1} \end{bmatrix},$$

$$\det(D_n) = 10^{-n}, \quad \kappa_p(D_n) = 1, \quad \sigma_n(D_n) = 10^{-1}.$$



Appendix

Proof of Theorem 7: Without loss of generality (W.L.O.G.) we can assume that $M(x) = \|x\|_\infty$ and N is arbitrary. We claim

$$c_1 \|x\|_\infty \leq N(x) \leq c_2 \|x\|_\infty$$

or

$$c_1 \leq N(z) \leq c_2, \text{ for } z \in S = \{z \in \mathbf{C}^n \mid \|z\|_\infty = 1\}.$$

From Lemma 6, N is continuous on S (closed and bounded). By maximum and minimum principle, there are $c_1, c_2 \geq 0$ and $z_1, z_2 \in S$ such that

$$c_1 = N(z_1) \leq N(z) \leq N(z_2) = c_2.$$

If $c_1 = 0$, then $N(z_1) = 0$. Thus, $z_1 = 0$. This contradicts that $z_1 \in S$. □

▶ Return



Proof of (3):

$$\|Ax\|_1 = \sum_i \left| \sum_j a_{ij}x_j \right| \leq \sum_i \sum_j |a_{ij}| |x_j| = \sum_j |x_j| \sum_i |a_{ij}|.$$

Let

$$\mathcal{C} := \sum_i |a_{ik}| = \max_j \sum_i |a_{ij}|.$$

Then $\|Ax\|_1 \leq \mathcal{C} \|x\|_1$, thus $\|A\|_1 \leq \mathcal{C}$. On the other hand, $\|e_k\|_1 = 1$ and $\|Ae_k\|_1 = \sum_{i=1}^n |a_{ik}| = \mathcal{C}$. □



Proof of (4):

$$\begin{aligned}\|Ax\|_\infty &= \max_i \left| \sum_j a_{ij}x_j \right| \leq \max_i \sum_j |a_{ij}x_j| \\ &\leq \max_i \sum_j |a_{ij}| \|x\|_\infty \equiv \sum_j |a_{kj}| \|x\|_\infty \equiv \hat{C} \|x\|_\infty.\end{aligned}$$

This implies, $\|A\|_\infty \leq \hat{C}$. If $A = 0$, then there is nothing to prove. Assume $A \neq 0$. Thus, the k -th row of A is nonzero. Define $z = [z_i] \in \mathbb{C}^n$ by

$$\begin{cases} z_i = \frac{\bar{a}_{ki}}{|a_{ki}|} & \text{if } a_{ki} \neq 0, \\ z_i = 1 & \text{if } a_{ki} = 0. \end{cases}$$

Then $\|z\|_\infty = 1$ and $a_{kj}z_j = |a_{kj}|$, for $j = 1, \dots, n$. It follows

$$\|A\|_\infty \geq \|Az\|_\infty = \max_i \left| \sum_j a_{ij}z_j \right| \geq \left| \sum_j a_{kj}z_j \right| = \sum_{j=1}^n |a_{kj}| \equiv \hat{C}.$$

Then, $\|A\|_\infty \geq \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}| \equiv \hat{C}$.



Proof of (5): Let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$ be the eigenvalues of A^*A . There are mutually orthonormal vectors v_j , $j = 1, \dots, n$ such that $(A^*A)v_j = \lambda_j v_j$. Let $x = \sum_j \alpha_j v_j$. Since $\|Ax\|_2^2 = (Ax, Ax) = (x, A^*Ax)$,

$$\|Ax\|_2^2 = \left(\sum_j \alpha_j v_j, \sum_j \alpha_j \lambda_j v_j \right) = \sum_j \lambda_j |\alpha_j|^2 \leq \lambda_1 \|x\|_2^2.$$

Therefore, $\|A\|_2^2 \leq \lambda_1$. Equality follows by choosing $x = v_1$ and $\|Av_1\|_2^2 = (v_1, \lambda_1 v_1) = \lambda_1$. So, we have $\|A\|_2 = \sqrt{\rho(A^*A)}$. □

▶ Return



Proof of Theorem 12: Let $|\lambda| = \rho(A) \equiv \rho$ and x be the associated eigenvector with $\|x\| = 1$. Then,

$$\rho(A) = |\lambda| = \|\lambda x\| = \|Ax\| \leq \|A\| \|x\| = \|A\|.$$

Claim: $\|\cdot\|_\epsilon \leq \rho(A) + \epsilon$. There is a unitary U such that $A = U^*RU$, where R is upper triangular.

Let $D_t = \text{diag}(t, t^2, \dots, t^n)$. For $t > 0$ large enough, the sum of all absolute values of the off-diagonal elements of $D_tRD_t^{-1}$ is less than ϵ . So, it holds $\|D_tRD_t^{-1}\|_1 \leq \rho(A) + \epsilon$ for large $t(\epsilon) > 0$. Define $\|\cdot\|_\epsilon$ for any B by

$$\begin{aligned} \|B\|_\epsilon &= \|D_tUBU^*D_t^{-1}\|_1 \\ &= \|(UD_t^{-1})^{-1}B(UD_t^{-1})\|_1. \end{aligned}$$

This implies,

$$\|A\|_\epsilon = \|D_tRD_t^{-1}\| \leq \rho(A) + \epsilon.$$



Return

Proof of Theorem 14: There are $x \in \mathbb{C}^n$, $y \in \mathbb{C}^m$ with $\|x\|_2 = \|y\|_2 = 1$ such that $Ax = \sigma y$, where $\sigma = \|A\|_2$ ($\|A\|_2 = \sup_{\|x\|_2=1} \|Ax\|_2$). Let $V = [x, V_1] \in \mathbb{C}^{n \times n}$, and $U = [y, U_1] \in \mathbb{C}^{m \times m}$ be unitary. Then

$$A_1 \equiv U^* A V = \begin{bmatrix} \sigma & w^* \\ 0 & B \end{bmatrix}.$$

Since $\left\| A_1 \begin{pmatrix} \sigma \\ w \end{pmatrix} \right\|_2^2 \geq (\sigma^2 + w^* w)^2$, it follows

$$\|A_1\|_2^2 \geq \sigma^2 + w^* w \quad \text{from} \quad \frac{\left\| A_1 \begin{pmatrix} \sigma \\ w \end{pmatrix} \right\|_2^2}{\left\| \begin{pmatrix} \sigma \\ w \end{pmatrix} \right\|_2^2} \geq \sigma^2 + w^* w.$$

But $\sigma^2 = \|A\|_2^2 = \|A_1\|_2^2$, it implies $w = 0$. Hence, the theorem holds by induction.



Return

Proof of (8): Claim $\|x\|_q \leq \|x\|_p$, ($p \leq q$): It holds

$$\|x\|_q = \left\| \left\| \|x\|_p \frac{x}{\|x\|_p} \right\|_q \right\| = \|x\|_p \left\| \frac{x}{\|x\|_p} \right\|_q \leq \mathcal{C}_{p,q} \|x\|_p,$$

where

$$\mathcal{C}_{p,q} = \max_{\|e\|_p=1} \|e\|_q, \quad e = (e_1, \dots, e_n)^T.$$

We now show that $\mathcal{C}_{p,q} \leq 1$. From $p \leq q$, we have

$$\|e\|_q^q = \sum_{i=1}^n |e_i|^q \leq \sum_{i=1}^n |e_i|^p = 1 \quad (\text{by } |e_i| \leq 1).$$

Hence, $\mathcal{C}_{p,q} \leq 1$, thus $\|x\|_q \leq \|x\|_p$.



To prove the second inequality: Let $\alpha = q/p > 1$. Then the Jensen inequality holds for the convex function $\varphi(x) \equiv x^\alpha$:

$$\int_{\Omega} |f|^q dx = \int_{\Omega} (|f|^p)^{q/p} dx \geq \left(\int_{\Omega} |f|^p dx \right)^{q/p}$$

with $|\Omega| = 1$. Consider the discrete measure $\sum_{i=1}^n \frac{1}{n} = 1$ and $f(i) = |x_i|$. It follows that

$$\sum_{i=1}^n |x_i|^q \frac{1}{n} \geq \left(\sum_{i=1}^n |x_i|^p \frac{1}{n} \right)^{q/p}.$$

Hence, we have

$$n^{-\frac{1}{q}} \|x\|_q \geq n^{-\frac{1}{p}} \|x\|_p.$$

Thus,

$$n^{(q-p)/pq} \|x\|_q \geq \|x\|_p.$$



Return

Proof of (9): Let $q \rightarrow \infty$ and $\lim_{q \rightarrow \infty} \|x\|_q = \|x\|_\infty$:

$$\|x\|_\infty = |x_k| = (|x_k|^q)^{\frac{1}{q}} \leq \left(\sum_{i=1}^n |x_i|^q \right)^{\frac{1}{q}} = \|x\|_q.$$

On the other hand,

$$\|x\|_q = \left(\sum_{i=1}^n |x_i|^q \right)^{\frac{1}{q}} \leq (n \|x\|_\infty^q)^{\frac{1}{q}} \leq n^{\frac{1}{q}} \|x\|_\infty.$$

It follows that $\lim_{q \rightarrow \infty} \|x\|_q = \|x\|_\infty$. □

▶ Return



To prove the second inequality: Let $\alpha = q/p > 1$. Then the Jensen inequality holds for the convex function $\varphi(x) \equiv x^\alpha$:

$$\int_{\Omega} |f|^q dx = \int_{\Omega} (|f|^p)^{q/p} dx \geq \left(\int_{\Omega} |f|^p dx \right)^{q/p}$$

with $|\Omega| = 1$.

Consider the discrete measure $\sum_{i=1}^n \frac{1}{n} = 1$ and $f(i) = |x_i|$. It follows that

$$\sum_{i=1}^n |x_i|^q \frac{1}{n} \geq \left(\sum_{i=1}^n |x_i|^p \frac{1}{n} \right)^{q/p}.$$

Hence, we have

$$n^{-\frac{1}{q}} \|x\|_q \geq n^{-\frac{1}{p}} \|x\|_p.$$

Thus,

$$n^{(q-p)/pq} \|x\|_q \geq \|x\|_p.$$



Return

Proof of (10): The first inequality holds obviously. Now, for the second inequality, we have

$$\begin{aligned}\|Ay\|_p &\leq \sum_{j=1}^n |y_j| \|a_j\|_p \\ &\leq \sum_{j=1}^n |y_j| \max_j \|a_j\|_p \\ &= \|y\|_1 \max_j \|a_j\|_p \\ &\leq n^{(p-1)/p} \max_j \|a_j\|_p. \quad (\text{by (8)})\end{aligned}$$

□

[Return](#)

Proof of (15): It holds

$$\begin{aligned}\max_{\|x\|_p=1} \|Ax\|_p &= \max_{\|x\|_p=1} \max_{\|y\|_q=1} |(Ax)^T y| \\ &= \max_{\|y\|_q=1} \max_{\|x\|_p=1} |x^T (A^T y)| \\ &= \max_{\|y\|_q=1} \|A^T y\|_q \\ &= \|A^T\|_q.\end{aligned}$$

□

Proof of (16): By (12) and (15), we get

$$\begin{aligned}m^{\frac{1}{p}} \|A\|_\infty &= m^{\frac{1}{p}} \|A^T\|_1 = m^{1-\frac{1}{q}} \|A^T\|_1 \\ &= m^{(q-1)/q} \|A^T\|_1 \geq \|A^T\|_q = \|A\|_p.\end{aligned}$$

□



Return

Proof of (17): It holds

$$\|A\|_p \|A\|_q = \|A^T\|_q \|A\|_q \geq \|A^T A\|_q \geq \|A^T A\|_2.$$

The last inequality holds by the following statement: Let S be a symmetric matrix. Then $\|S\|_2 \leq \|S\|$, for any matrix operator norm $\|\cdot\|$. Since $|\lambda| \leq \|S\|$,

$$\|S\|_2 = \sqrt{\rho(S^*S)} = \sqrt{\rho(S^2)} = \max_{\lambda \in \sigma(S)} |\lambda| = |\lambda_{\max}|.$$

This implies, $\|S\|_2 \leq \|S\|$.



Proof of (18): By (8), we get

$$\|A\|_p = \max_{\|x\|_p=1} \|Ax\|_p \leq \max_{\|x\|_q \leq 1} m^{(q-p)/pq} \|Ax\|_q = m^{(q-p)/pq} \|A\|_q.$$



Return