

CS229 Section: Linear Algebra

Nandita Bhaskhar

Slides adapted from past CS229 teams

April 1, 2022

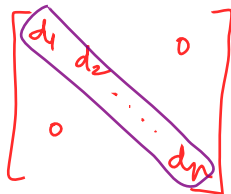
Basic Concepts and Notation

Diagonal matrices

A **diagonal matrix** is a matrix where all non-diagonal elements are 0. This is typically denoted $D = \text{diag}(d_1, d_2, \dots, d_n)$, with

$$D_{ij} = \begin{cases} d_i & i = j \\ 0 & i \neq j \end{cases}$$

Clearly, $I = \text{diag}(1, 1, \dots, 1)$.



Outline

- 1 Basic Concepts and Notation
- 2 Matrix Multiplication**
- 3 Operations and Properties
- 4 Matrix Calculus

Vector-Vector Product

- inner product* or *dot product*

vector vector

$$x^T y \in \mathbb{R} = [\underbrace{x_1 \ x_2 \ \dots \ x_n}_{(1 \times n)}] \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \sum_{i=1}^n x_i y_i$$

equal (1 x 1) scalar

need not be equal (n x 1)

- outer product*

vector matrix

$$x y^T \in \mathbb{R}^{m \times n} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} [\underbrace{y_1 \ y_2 \ \dots \ y_n}_{(1 \times n)}] = \begin{bmatrix} x_1 y_1 & x_1 y_2 & \dots & x_1 y_n \\ x_2 y_1 & x_2 y_2 & \dots & x_2 y_n \\ \vdots & \vdots & \ddots & \vdots \\ x_m y_1 & x_m y_2 & \dots & x_m y_n \end{bmatrix} \begin{matrix} (m \times n) \\ \cdot \\ \text{matrix} \end{matrix}$$

vector

Matrix-Vector Product

connect to class notes $\left[\begin{aligned} h(x^{(i)}) &= (x^{(i)})^T \theta \\ h(x) &= x \theta \end{aligned} \right]$

- If we write A by rows, then we can express Ax as,

$$y = Ax = \begin{matrix} (1 \times n) \\ \left[\begin{array}{ccc} \text{---} & a_1^T & \text{---} \\ \text{---} & a_2^T & \text{---} \\ & \vdots & \\ \text{---} & a_m^T & \text{---} \end{array} \right] \end{matrix} \begin{matrix} (n \times 1) \\ \left[\begin{array}{c} \vdots \\ x \\ \vdots \end{array} \right] \end{matrix} = \begin{matrix} (1 \times 1) \\ \left[\begin{array}{c} a_1^T x \\ a_2^T x \\ \vdots \\ a_m^T x \end{array} \right] \end{matrix} \cdot$$

vector of inner products

Matrix-Vector Product

- If we write A by columns, then we have:

$$y = Ax = \left[\begin{array}{c|c|ccc|c} a^1 & a^2 & \dots & a^n \end{array} \right] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \underbrace{\begin{bmatrix} a^1 \end{bmatrix}}_{(m \times 1)} x_1 + \underbrace{\begin{bmatrix} a^2 \end{bmatrix}}_{(m \times 1)} x_2 + \dots + \begin{bmatrix} a^n \end{bmatrix} x_n \quad (1)$$

y is a linear combination of the columns of A .



Matrix-Vector Product

It is also possible to multiply on the left by a row vector.

- If we write A by columns, then we can express $x^T A$ as,

$$y^T = x^T A = x^T \begin{array}{c} \text{\scriptsize (m \times h)} \\ \left[\begin{array}{ccc|ccc} a^1 & a^2 & \dots & a^n \\ \hline | & | & & | \\ | & | & & | \end{array} \right] = \left[\begin{array}{cccc} x^T a^1 & x^T a^2 & \dots & x^T a^n \end{array} \right] \\ \begin{array}{c} \text{\scriptsize (1 \times m)} \\ \text{\scriptsize (m \times 1)} \end{array} \end{array} \quad \begin{array}{c} \text{\scriptsize (1 \times 1)} \\ \longrightarrow n \\ \text{\scriptsize (1 \times h)} \end{array}$$

Matrix-Vector Product

It is also possible to multiply on the left by a row vector.

- expressing A in terms of rows we have:

$$y^T = x^T A = \begin{matrix} & \begin{matrix} (m \times n) \end{matrix} \\ \begin{matrix} (1 \times m) \end{matrix} & \begin{bmatrix} \text{---} & a_1^T & \text{---} \\ \text{---} & a_2^T & \text{---} \\ & \vdots & \\ \text{---} & a_m^T & \text{---} \end{bmatrix} \end{matrix} \begin{matrix} (1 \times n) \\ \\ \\ \\ \end{matrix}$$

$$= x_1 \begin{bmatrix} \text{---} & a_1^T & \text{---} \end{bmatrix} + x_2 \begin{bmatrix} \text{---} & a_2^T & \text{---} \end{bmatrix} + \dots + x_m \begin{bmatrix} \text{---} & a_m^T & \text{---} \end{bmatrix}$$

y^T is a linear combination of the rows of A .



Matrix-Matrix Multiplication (different views)

1. As a set of vector-vector products (dot product)

$$\begin{matrix}
 (m \times p) \\
 C = AB =
 \end{matrix}
 \begin{matrix}
 (m \times n) & (n \times p) \\
 \begin{bmatrix}
 - & a_1^T & - \\
 - & a_2^T & - \\
 \boxed{\dots} & \boxed{i} & \dots \\
 - & a_m^T & -
 \end{bmatrix}
 \begin{bmatrix}
 | & | & \boxed{j} & | \\
 b^1 & b^2 & & b^p \\
 | & | & & |
 \end{bmatrix}
 \end{matrix}
 =
 \begin{bmatrix}
 a_1^T b^1 & a_1^T b^2 & \dots & a_1^T b^p \\
 a_2^T b^1 & a_2^T b^2 & \dots & a_2^T b^p \\
 \vdots & \vdots & \odot_{ij} & \vdots \\
 a_m^T b^1 & a_m^T b^2 & \dots & a_m^T b^p
 \end{bmatrix}
 \cdot$$

$$\begin{matrix}
 (i \times 1) \\
 c_{ij} = a_i^T b^j \\
 (1 \times n) \quad (n \times 1)
 \end{matrix}$$

Matrix-Matrix Multiplication (different views)

2. As a sum of outer products

$$C = AB = \begin{matrix} \text{\textcircled{m \times n}} & \text{\textcircled{m \times p}} & & \text{\textcircled{p \times n}} \\ \left[\begin{array}{ccc|c} a^1 & a^2 & \dots & a^p \\ \hline & & & \\ \hline & & & \\ \hline & & & \end{array} \right] & \begin{matrix} \left[\begin{array}{c|c|c} - & b_1^T & - \\ - & b_2^T & - \\ & \vdots & \\ - & b_p^T & - \end{array} \right] & \end{matrix} & = & \sum_{i=1}^p \text{\textcircled{a^i b_i^T}}. \end{matrix}$$

(m x n) *(m x p)* *(p x n)*

(m x 1) → p *(1 x n)*

Matrix-Matrix Multiplication (different views)

3. As a set of matrix-vector products.

$$C = AB = A \left[\begin{array}{c|c|c|c} b^1 & b^2 & \dots & b^n \\ \hline \end{array} \right] = \left[\begin{array}{c|c|c|c} Ab^1 & Ab^2 & \dots & Ab^n \\ \hline \end{array} \right]. \quad (2)$$

Here the i th column of C is given by the matrix-vector product with the vector on the right, $c_i = Ab_i$. These matrix-vector products can in turn be interpreted using both viewpoints given in the previous subsection.

Matrix-Matrix Multiplication (different views)

4. As a set of vector-matrix products.

$$C = AB = \begin{bmatrix} - & a_1^T & - \\ - & a_2^T & - \\ & \vdots & \\ - & a_m^T & - \end{bmatrix} B = \begin{bmatrix} - & a_1^T B & - \\ - & a_2^T B & - \\ & \vdots & \\ - & a_m^T B & - \end{bmatrix}.$$

j^{th} row of C : $a_j^T B$

Matrix-Matrix Multiplication (properties)

$\neq BAC$ (order is $\cdot \cdot \cdot$)



• Associative: $(AB)C = A(BC)$.

• Distributive: $A(B + C) = AB + AC$.

} prove!

• In general, not commutative; that is, it can be the case that $AB \neq BA$. (For example, if $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{n \times q}$, the matrix product BA does not even exist if m and q are not equal!)

not necessary

counter example!

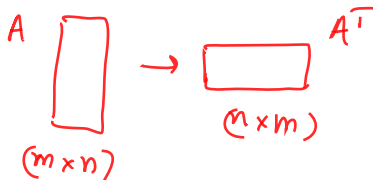
Outline

- 1 Basic Concepts and Notation
- 2 Matrix Multiplication
- 3 Operations and Properties**
- 4 Matrix Calculus

The Transpose

The *transpose* of a matrix results from "flipping" the rows and columns. Given a matrix $A \in \mathbb{R}^{m \times n}$, its transpose, written $A^T \in \mathbb{R}^{n \times m}$, is the $n \times m$ matrix whose entries are given by

$$(A^T)_{ij} = A_{ji}.$$



The following properties of transposes are easily verified:


- $(A^T)^T = A$ → flip twice
- $(AB)^T = \underline{B^T} \underline{A^T}$ order changes!
- $(A + B)^T = A^T + B^T$

sum → flip ≡ flip → sum

} prove!

Trace

The trace of a square matrix $A \in \mathbb{R}^{n \times n}$, denoted $\text{tr}A$, is the sum of diagonal elements in the matrix:

$$\text{tr}A = \sum_{i=1}^n A_{ii}.$$


The trace has the following properties:

• For $A \in \mathbb{R}^{n \times n}$, $\text{tr}A = \text{tr}A^T$. *→ flipping doesn't change diagonals*

• For $A, B \in \mathbb{R}^{n \times n}$, $\text{tr}(A + B) = \text{tr}A + \text{tr}B$. *verify!*

• For $A \in \mathbb{R}^{n \times n}$, $t \in \mathbb{R}$, $\text{tr}(tA) = t \text{tr}A$.

• For A, B such that AB is square, $\text{tr}AB = \text{tr}BA$. *diag. are special!*

• For A, B, C such that ABC is square, $\text{tr}ABC = \text{tr}BCA = \text{tr}CAB$, and so on for the product of more matrices. *≠ tr ACB (order is ~~not~~.)* *Prove!*

Norms

A norm of a vector $\|x\|$ is informally a measure of the “length” of the vector.

More formally, a norm is any function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ that satisfies 4 properties:

1. For all $x \in \mathbb{R}^n$, $f(x) \geq 0$ (non-negativity).
 2. $f(x) = 0$ if and only if $x = 0$ (definiteness).
 3. For all $x \in \mathbb{R}^n$, $t \in \mathbb{R}$, $f(tx) = |t|f(x)$ (homogeneity). *scaling comp. scales length*
 4. For all $x, y \in \mathbb{R}^n$, $f(x + y) \leq f(x) + f(y)$ (triangle inequality). *(geometry)*
- } length !*

Examples of Norms

The commonly-used Euclidean or l_2 norm,

$$\|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2}.$$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

The l_1 norm,

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

The l_∞ norm,

$$\|x\|_\infty = \max_j |x_j|.$$

Examples of Norms

In fact, all three norms presented so far are examples of the family of ℓ_p norms, which are parameterized by a real number $p \geq 1$, and defined as

$$\|x\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p}$$

Matrix Norms

Norms can also be defined for matrices, such as the Frobenius norm,

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n A_{ij}^2} = \sqrt{\text{tr}(A^T A)}.$$

Many other norms exist, but they are beyond the scope of this review.

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ a_{im} & \dots & \dots & a_{mn} \end{bmatrix}$$

↪ useful,
will see this
again &
again

Linear Independence

A set of vectors $\{x_1, x_2, \dots, x_n\} \subset \mathbb{R}^m$ is said to be **(linearly) dependent** if ^{any} one vector belonging to the set can be represented as a linear combination of the remaining vectors; that is, if

$$x_n = \sum_{i=1}^{n-1} \alpha_i x_i$$

Scalar coeff.

for some scalar values $\alpha_1, \dots, \alpha_{n-1} \in \mathbb{R}$; otherwise, the vectors are **(linearly) independent**.

Linear Independence

A set of vectors $\{x_1, x_2, \dots, x_n\} \subset \mathbb{R}^m$ is said to be (*linearly*) *dependent* if one vector belonging to the set *can* be represented as a linear combination of the remaining vectors; that is, if

$$x_n = \sum_{i=1}^{n-1} \alpha_i x_i$$

for some scalar values $\alpha_1, \dots, \alpha_{n-1} \in \mathbb{R}$; otherwise, the vectors are (*linearly*) *independent*.

Example:

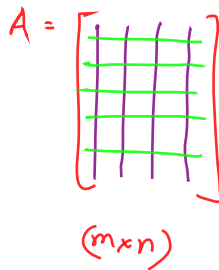
$$x_1 = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \quad x_2 = \begin{bmatrix} 4 \\ 1 \\ 5 \end{bmatrix} \quad x_3 = \begin{bmatrix} 2 \\ -3 \\ -1 \end{bmatrix}$$

are linearly dependent because $x_3 = -2x_1 + x_2$.

Rank of a Matrix

- The column rank of a matrix $A \in \mathbb{R}^{m \times n}$ is the largest number of columns of A that constitute a linearly independent set.

$$\text{col. rank} \leq \# \text{ cols} = n$$



Rank of a Matrix

- The *column rank* of a matrix $A \in \mathbb{R}^{m \times n}$ is the largest number of columns of A that constitute a linearly independent set.
- The *row rank* is the largest number of rows of A that constitute a linearly independent set.

$$\text{row rank} \leq \# \text{ rows} = m$$

$$A = \begin{bmatrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{bmatrix} \quad (m \times n)$$

Rank of a Matrix

- The **column rank** of a matrix $A \in \mathbb{R}^{m \times n}$ is the largest number of columns of A that constitute a linearly independent set.
- The **row rank** is the largest number of rows of A that constitute a linearly independent set.
- For any matrix $A \in \mathbb{R}^{m \times n}$, it turns out that the column rank of A is equal to the row rank of A (prove it yourself!), and so both quantities are referred to collectively as the **rank** of A , denoted as $\text{rank}(A)$.

$$\begin{array}{c} \text{col. rank} \\ \text{of } A \end{array} = \begin{array}{c} \text{row rank} \\ \text{of } A \end{array} = \text{rank}(A)$$

(prove!)

Properties of the Rank

- For $A \in \mathbb{R}^{m \times n}$, $\text{rank}(A) \leq \min(m, n)$. If $\text{rank}(A) = \min(m, n)$ then A is said to be **full rank**.

Properties of the Rank

- For $A \in \mathbb{R}^{m \times n}$, $\text{rank}(A) \leq \min(m, n)$. If $\text{rank}(A) = \min(m, n)$, then A is said to be **full rank**.
- For $A \in \mathbb{R}^{m \times n}$, $\text{rank}(A) = \text{rank}(A^T)$. *flipping doesn't change linear independence!*

Properties of the Rank

- For $A \in \mathbb{R}^{m \times n}$, $\text{rank}(A) \leq \min(m, n)$. If $\text{rank}(A) = \min(m, n)$, then A is said to be **full rank**.
- For $A \in \mathbb{R}^{m \times n}$, $\text{rank}(A) = \text{rank}(A^T)$.
- For $A \in \mathbb{R}^{m \times p}$, $B \in \mathbb{R}^{p \times n}$, $\text{rank}(AB) \leq \min(\text{rank}(A), \text{rank}(B))$.



Properties of the Rank

- For $A \in \mathbb{R}^{m \times n}$, $\text{rank}(A) \leq \min(m, n)$. If $\text{rank}(A) = \min(m, n)$, then A is said to be **full rank**.
- For $A \in \mathbb{R}^{m \times n}$, $\text{rank}(A) = \text{rank}(A^T)$.
- For $A \in \mathbb{R}^{m \times p}$, $B \in \mathbb{R}^{p \times n}$, $\text{rank}(AB) \leq \min(\text{rank}(A), \text{rank}(B))$.
- For $A, B \in \mathbb{R}^{m \times n}$, $\text{rank}(A + B) \leq \text{rank}(A) + \text{rank}(B)$.

The Inverse of a Square Matrix

- The inverse of a square matrix $A \in \mathbb{R}^{n \times n}$ is denoted A^{-1} , and is the unique matrix such that

$$A^{-1}A = I = AA^{-1}.$$

Does this always exist?

The Inverse of a Square Matrix

- The *inverse* of a square matrix $A \in \mathbb{R}^{n \times n}$ is denoted A^{-1} , and is the unique matrix such that

$$A^{-1}A = I = AA^{-1}.$$

- We say that A is *invertible* or *non-singular* if A^{-1} exists and *non-invertible* or *singular* otherwise.

When does it exist ??

The Inverse of a Square Matrix

- The *inverse* of a square matrix $A \in \mathbb{R}^{n \times n}$ is denoted A^{-1} , and is the unique matrix such that

$$A^{-1}A = I = AA^{-1}.$$

- We say that A is *invertible* or *non-singular* if A^{-1} exists and *non-invertible* or *singular* otherwise.
- In order for a square matrix A to have an inverse A^{-1} , then A must be full rank.

*if $A \rightarrow$ full rank : invertible
square non-singular*

The Inverse of a Square Matrix

- The *inverse* of a square matrix $A \in \mathbb{R}^{n \times n}$ is denoted A^{-1} , and is the unique matrix such that

$$A^{-1}A = I = AA^{-1}.$$

- We say that A is *invertible* or *non-singular* if A^{-1} exists and *non-invertible* or *singular* otherwise.
- In order for a square matrix A to have an inverse A^{-1} , then A must be full rank.
- Properties (Assuming $A, B \in \mathbb{R}^{n \times n}$ are non-singular):
 - ✓ $(A^{-1})^{-1} = A$
 - ▶ $(AB)^{-1} = \underline{B^{-1}} \underline{A^{-1}}$ *order!*
 - ▶ $(A^{-1})^T = (\underline{A^T})^{-1}$. For this reason this matrix is often denoted A^{-T} .

} prove!

Orthogonal Matrices



- Two vectors $x, y \in \mathbb{R}^n$ are **orthogonal** if $x^T y = 0$.
- A vector $x \in \mathbb{R}^n$ is **normalized** if $\|x\|_2 = 1$.
- A square matrix $U \in \mathbb{R}^{n \times n}$ is **orthogonal** if all its columns are orthogonal to each other and are normalized (the columns are then referred to as being orthonormal).

$$U = \begin{bmatrix} | & | & \dots & | \\ u_1 & u_2 & \dots & u_n \\ | & | & \dots & | \end{bmatrix}$$

$(n \times n)$

cols: orthogonal to each other
normalized

$$u_i^T u_j = \boxed{0} \quad u_i^T u_i = \boxed{1}$$

$(j \neq i)$

Orthogonal Matrices

- Two vectors $x, y \in \mathbb{R}^n$ are **orthogonal** if $x^T y = 0$.
- A vector $x \in \mathbb{R}^n$ is **normalized** if $\|x\|_2 = 1$.
- A square matrix $U \in \mathbb{R}^{n \times n}$ is **orthogonal** if all its columns are orthogonal to each other and are normalized (the columns are then referred to as being **orthonormal**).

- **Properties:**

- ▶ The inverse of an orthogonal matrix is its transpose.

$$U^T U = I = U U^T.$$

Orthogonal Matrices

- Two vectors $x, y \in \mathbb{R}^n$ are **orthogonal** if $x^T y = 0$.
- A vector $x \in \mathbb{R}^n$ is **normalized** if $\|x\|_2 = 1$.
- A square matrix $U \in \mathbb{R}^{n \times n}$ is **orthogonal** if all its columns are orthogonal to each other and are normalized (the columns are then referred to as being **orthonormal**).

- **Properties:**

- ▶ The inverse of an orthogonal matrix is its transpose.

$$U^T U = I = U U^T.$$

- ▶ Operating on a vector with an orthogonal matrix will not change its Euclidean norm, i.e.,

$$\|Ux\|_2 = \|x\|_2 \quad \text{no scaling in } l_2 \text{ norm}$$

for any $x \in \mathbb{R}^n$, $U \in \mathbb{R}^{n \times n}$ orthogonal.

Span and Projection

- The span of a set of vectors $\{x_1, x_2, \dots, x_n\}$ is the set of all vectors that can be expressed as a linear combination of $\{x_1, \dots, x_n\}$. That is,

$$\text{span}(\{x_1, \dots, x_n\}) = \left\{ v : v = \sum_{i=1}^n \alpha_i x_i, \alpha_i \in \mathbb{R} \right\}.$$



Span and Projection

- The **span** of a set of vectors $\{x_1, x_2, \dots, x_n\}$ is the set of all vectors that can be expressed as a linear combination of $\{x_1, \dots, x_n\}$. That is,

$$\text{span}(\{x_1, \dots, x_n\}) = \left\{ v : v = \sum_{i=1}^n \alpha_i x_i, \alpha_i \in \mathbb{R} \right\}.$$

- The **projection** of a vector $y \in \mathbb{R}^m$ onto the **span** of $\{x_1, \dots, x_n\}$ is the vector $v \in \text{span}(\{x_1, \dots, x_n\})$, such that v is as close as possible to y , as measured by the Euclidean norm $\|v - y\|_2$.

$$\text{Proj}(y; \{x_1, \dots, x_n\}) = \underset{v \in \text{span}(\{x_1, \dots, x_n\})}{\text{argmin}} \|y - v\|_2.$$

Range

- The range or the column space of a matrix $A \in \mathbb{R}^{m \times n}$, denoted $\mathcal{R}(A)$, is the span of the columns of A . In other words,

$$\mathcal{R}(A) = \{v \in \mathbb{R}^m : v = Ax, x \in \mathbb{R}^n\}.$$

Range

- The **range** or the column space of a matrix $A \in \mathbb{R}^{m \times n}$, denoted $\mathcal{R}(A)$, is the the span of the columns of A . In other words,

$$\mathcal{R}(A) = \{v \in \mathbb{R}^m : v = Ax, x \in \mathbb{R}^n\}.$$

- Assuming A is full rank and $n < m$, the projection of a vector $y \in \mathbb{R}^m$ onto the range of A is given by,

$$\text{Proj}(y; A) = \underset{v \in \mathcal{R}(A)}{\text{argmin}} \|v - y\|_2.$$

Null space

The nullspace of a matrix $A \in \mathbb{R}^{m \times n}$, denoted $\mathcal{N}(A)$ is the set of all vectors that equal 0 when multiplied by A , i.e.,

$$\mathcal{N}(A) = \{x \in \mathbb{R}^n : \underline{Ax} = 0\}.$$

remember orthogonal vectors!

(Rank-nullity theorem)

The Determinant

matrix *scalar*

The determinant of a square matrix $A \in \mathbb{R}^{n \times n}$, is a function $\det : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}$, and is denoted $|A|$ or $\det A$.

Given a matrix

$$\begin{bmatrix} - & a_1^T & - \\ - & a_2^T & - \\ & \vdots & \\ - & a_n^T & - \end{bmatrix},$$

geometrically:

consider the set of points $S \subset \mathbb{R}^n$ as follows:

*restricted
span*

$$S = \{v \in \mathbb{R}^n : v = \sum_{i=1}^n \alpha_i a_i \text{ where } 0 \leq \alpha_i \leq 1, i = 1, \dots, n\}.$$

The absolute value of the determinant of A is a measure of the “volume” of the set S .

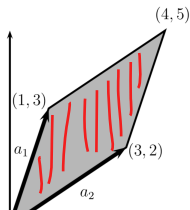
The Determinant: Intuition

For example, consider the 2×2 matrix,

$$A = \begin{bmatrix} 1 & 3 \\ 3 & 2 \end{bmatrix} \quad (3)$$

Here, the rows of the matrix are

$$a_1 = \begin{bmatrix} 1 \\ 3 \end{bmatrix} \quad a_2 = \begin{bmatrix} 3 \\ 2 \end{bmatrix}$$



parallelogram

The Determinant: Properties

Algebraically, the determinant satisfies the following three properties:

1. The determinant of the identity is 1, $|I| = 1$. (Geometrically, the volume of a unit hypercube is 1).

The Determinant: Properties

Algebraically, the determinant satisfies the following three properties:

1. The determinant of the identity is 1, $|I| = 1$. (Geometrically, the volume of a unit hypercube is 1).
2. Given a matrix $A \in \mathbb{R}^{n \times n}$, if we multiply a single row in A by a scalar $t \in \mathbb{R}$, then the determinant of the new matrix is $t|A|$, (Geometrically, multiplying one of the sides of the set S by a factor t causes the volume to increase by a factor t .)

The Determinant: Properties

Algebraically, the determinant satisfies the following three properties:

1. The determinant of the identity is 1, $|I| = 1$. (Geometrically, the volume of a unit hypercube is 1).
2. Given a matrix $A \in \mathbb{R}^{n \times n}$, if we multiply a single row in A by a scalar $t \in \mathbb{R}$, then the determinant of the new matrix is $t|A|$, (Geometrically, multiplying one of the sides of the set S by a factor t causes the volume to increase by a factor t .)
3. If we exchange any two rows a_i^T and a_j^T of A , then the determinant of the new matrix is $-|A|$, for example

The Determinant: Properties

Algebraically, the determinant satisfies the following three properties:

1. The determinant of the identity is 1, $|I| = 1$. (Geometrically, the volume of a unit hypercube is 1).
2. Given a matrix $A \in \mathbb{R}^{n \times n}$, if we multiply a single row in A by a scalar $t \in \mathbb{R}$, then the determinant of the new matrix is $t|A|$, (Geometrically, multiplying one of the sides of the set S by a factor t causes the volume to increase by a factor t .)
3. If we exchange any two rows a_i^T and a_j^T of A , then the determinant of the new matrix is $-|A|$, for example

The Determinant: Properties

Algebraically, the determinant satisfies the following three properties:

1. The determinant of the identity is 1, $|I| = 1$. (Geometrically, the volume of a unit hypercube is 1).
2. Given a matrix $A \in \mathbb{R}^{n \times n}$, if we multiply a single row in A by a scalar $t \in \mathbb{R}$, then the determinant of the new matrix is $t|A|$, (Geometrically, multiplying one of the sides of the set S by a factor t causes the volume to increase by a factor t .)
3. If we exchange any two rows a_i^T and a_j^T of A , then the determinant of the new matrix is $-|A|$, for example

In case you are wondering, it is not immediately obvious that a function satisfying the above three properties exists. In fact, though, such a function does exist, and is unique (which we will not prove here).

The Determinant: Properties

- For $A \in \mathbb{R}^{n \times n}$, $|A| = |A^T|$.
- For $A, B \in \mathbb{R}^{n \times n}$, $|AB| = |A||B|$. *not full rank!*
- For $A \in \mathbb{R}^{n \times n}$, $|A| = 0$ if and only if A is singular (i.e., non-invertible). (If A is singular then it does not have full rank, and hence its columns are linearly dependent. In this case, the set S corresponds to a “flat sheet” within the n -dimensional space and hence has zero volume.)
- For $A \in \mathbb{R}^{n \times n}$ and A non-singular, $|A^{-1}| = 1/|A|$.

The Determinant: Formula !! 😱

Let $A \in \mathbb{R}^{n \times n}$, $A_{\setminus i, \setminus j} \in \mathbb{R}^{(n-1) \times (n-1)}$ be the *matrix* that results from deleting the i th row and j th column from A .

The general (recursive) formula for the determinant is

$$\begin{aligned} |A| &= \sum_{i=1}^n (-1)^{i+j} a_{ij} |A_{\setminus i, \setminus j}| \quad (\text{for any } j \in 1, \dots, n) \\ &= \sum_{j=1}^n (-1)^{i+j} a_{ij} |A_{\setminus i, \setminus j}| \quad (\text{for any } i \in 1, \dots, n) \end{aligned}$$

with the initial case that $|A| = a_{11}$ for $A \in \mathbb{R}^{1 \times 1}$. If we were to expand this formula completely for $A \in \mathbb{R}^{n \times n}$, there would be a total of $n!$ (n factorial) different terms. For this reason, we hardly ever explicitly write the complete equation of the determinant for matrices bigger than 3×3 .

The Determinant: Examples

However, the equations for determinants of matrices up to size 3×3 are fairly common, and it is good to know them:

$$\begin{aligned}
 |[a_{11}]| &= a_{11} \\
 \left| \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \right| &= a_{11}a_{22} - a_{12}a_{21} \\
 \left| \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \right| &= a_{11}a_{22}a_{33} + a_{12}a_{23}a_{31} + a_{13}a_{21}a_{32} \\
 &\quad - a_{11}a_{23}a_{32} - a_{12}a_{21}a_{33} - a_{13}a_{22}a_{31}
 \end{aligned}$$

Quadratic Forms



Given a square matrix $A \in \mathbb{R}^{n \times n}$ and a vector $x \in \mathbb{R}^n$, the scalar value $x^T A x$ is called a *quadratic form*. Written explicitly, we see that

$$x^T A x = \sum_{i=1}^n x_i (A x)_i = \sum_{i=1}^n x_i \left(\sum_{j=1}^n A_{ij} x_j \right) = \sum_{i=1}^n \sum_{j=1}^n A_{ij} x_i x_j .$$

↗ every element of A
 ↗ every element of x
 ↘ every element of x

⇓

$$(x^T A x)^T = x^T A x$$

Quadratic Forms

Given a square matrix $A \in \mathbb{R}^{n \times n}$ and a vector $x \in \mathbb{R}^n$, the scalar value $x^T Ax$ is called a *quadratic form*. Written explicitly, we see that

$$x^T Ax = \sum_{i=1}^n x_i (Ax)_i = \sum_{i=1}^n x_i \left(\sum_{j=1}^n A_{ij} x_j \right) = \sum_{i=1}^n \sum_{j=1}^n A_{ij} x_i x_j .$$

We often implicitly assume that the matrices appearing in a quadratic form are symmetric.

$$x^T Ax = (x^T Ax)^T = x^T A^T x = x^T \left(\frac{1}{2} A + \frac{1}{2} A^T \right) x,$$

scalar

only

parts contribute to Q. form

Positive Semidefinite Matrices

A symmetric matrix $A \in \mathbb{S}^n$ is:

- **positive definite** (PD), denoted $A \succ 0$ if for all non-zero vectors $x \in \mathbb{R}^n$, $x^T Ax > 0$.
- **positive semidefinite** (PSD), denoted $A \succeq 0$ if for all vectors $x^T Ax \geq 0$.
- **negative definite** (ND), denoted $A \prec 0$ if for all non-zero $x \in \mathbb{R}^n$, $x^T Ax < 0$.
- **negative semidefinite** (NSD), denoted $A \preceq 0$ if for all $x \in \mathbb{R}^n$, $x^T Ax \leq 0$.
- **indefinite**, if it is neither positive semidefinite nor negative semidefinite — i.e., if there exists $x_1, x_2 \in \mathbb{R}^n$ such that $x_1^T Ax_1 > 0$ and $x_2^T Ax_2 < 0$.

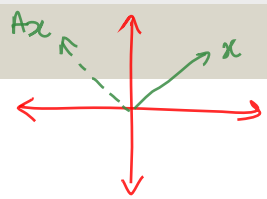
Positive Semidefinite Matrices

- One important property of positive definite and negative definite matrices is that they are always full rank, and hence, invertible. $\det \neq 0$
- Given any matrix $A \in \mathbb{R}^{m \times n}$ (not necessarily symmetric or even square), the matrix $G = A^T A$ (sometimes called a **Gram matrix**) is always positive semidefinite. Further, if $m \geq n$ and A is full rank, then $G = A^T A$ is positive definite.

$$\text{any } A \rightarrow \text{Gram}(A) = A^T A$$

Eigenvalues and Eigenvectors

$Ax \rightarrow$ transformation of x using A



Given a square matrix $A \in \mathbb{R}^{n \times n}$, we say that $\lambda \in \mathbb{C}$ is an **eigenvalue** of A and $x \in \mathbb{C}^n$ is the corresponding **eigenvector** if

$$Ax = \lambda x, \quad x \neq 0.$$

Intuitively, this definition means that multiplying A by the vector x results in a new vector that points in the same direction as x , but scaled by a factor λ .

Eigenvalues and Eigenvectors

$$Ax = \lambda Ix$$

We can rewrite the equation above to state that (λ, x) is an eigenvalue-eigenvector pair of A if,

$$(\lambda I - A)x = 0, \quad x \neq 0.$$

But $(\lambda I - A)x = 0$ has a non-zero solution to x if and only if $(\lambda I - A)$ has a non-empty nullspace, which is only the case if $(\lambda I - A)$ is singular, i.e.,

To find λ \Rightarrow $|(\lambda I - A)| = 0.$

We can now use the previous definition of the determinant to expand this expression $|(\lambda I - A)|$ into a (very large) polynomial in λ , where λ will have degree n . It's often called the characteristic polynomial of the matrix A .

Properties of eigenvalues and eigenvectors

- The trace of a A is equal to the sum of its eigenvalues,

$$\text{tr}A = \sum_{i=1}^n \lambda_i.$$

Properties of eigenvalues and eigenvectors

- The trace of a A is equal to the sum of its eigenvalues,

$$\operatorname{tr}A = \sum_{i=1}^n \lambda_i.$$

- The determinant of A is equal to the product of its eigenvalues,

$$|A| = \prod_{i=1}^n \lambda_i.$$

Properties of eigenvalues and eigenvectors

- The trace of a A is equal to the sum of its eigenvalues,

$$\operatorname{tr}A = \sum_{i=1}^n \lambda_i.$$

- The determinant of A is equal to the product of its eigenvalues,

$$|A| = \prod_{i=1}^n \lambda_i.$$

- The rank of A is equal to the number of non-zero eigenvalues of A .

Properties of eigenvalues and eigenvectors

- The trace of a A is equal to the sum of its eigenvalues,

$$\operatorname{tr}A = \sum_{i=1}^n \lambda_i.$$

- The determinant of A is equal to the product of its eigenvalues,

$$|A| = \prod_{i=1}^n \lambda_i.$$

- The rank of A is equal to the number of non-zero eigenvalues of A .
- Suppose A is non-singular with eigenvalue λ and an associated eigenvector x . Then $1/\lambda$ is an eigenvalue of A^{-1} with an associated eigenvector x , i.e., $A^{-1}x = (1/\lambda)x$.

Properties of eigenvalues and eigenvectors

- The trace of a A is equal to the sum of its eigenvalues,

$$\operatorname{tr}A = \sum_{i=1}^n \lambda_i.$$

- The determinant of A is equal to the product of its eigenvalues,

$$|A| = \prod_{i=1}^n \lambda_i.$$

- The rank of A is equal to the number of non-zero eigenvalues of A .
- Suppose A is non-singular with eigenvalue λ and an associated eigenvector x . Then $1/\lambda$ is an eigenvalue of A^{-1} with an associated eigenvector x , i.e., $A^{-1}x = (1/\lambda)x$.
- The eigenvalues of a diagonal matrix $D = \operatorname{diag}(d_1, \dots, d_n)$ are just the diagonal entries d_1, \dots, d_n .

Eigenvalues and Eigenvectors of Symmetric Matrices

Throughout this section, let's assume that A is a symmetric real matrix (i.e., $A^T = A$). We have the following properties:

1. All eigenvalues of A are real numbers. We denote them by $\lambda_1, \dots, \lambda_n$.
2. There exists a set of eigenvectors u_1, \dots, u_n such that (i) for all i , u_i is an eigenvector with eigenvalue λ_i and (ii) u_1, \dots, u_n are unit vectors and orthogonal to each other.

orthonormal
vectors

u_1, u_2, \dots, u_n for A (symmetric) with $\lambda_1, \lambda_2, \dots, \lambda_n$
eigenvalues

New Representation for Symmetric Matrices

- Let U be the orthonormal matrix that contains u_i 's as columns:

$$U = \begin{bmatrix} | & | & \cdots & | \\ u_1 & u_2 & \cdots & u_n \\ | & | & & | \end{bmatrix}$$

$(n \times n)$

$$A : n \times n$$

New Representation for Symmetric Matrices

- Let U be the orthonormal matrix that contains u_i 's as columns:

$$U = \begin{bmatrix} | & | & \cdots & | \\ u_1 & u_2 & \cdots & u_n \\ | & | & \cdots & | \end{bmatrix}$$

- Let $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ be the diagonal matrix that contains $\lambda_1, \dots, \lambda_n$.

$$AU = \begin{bmatrix} | & | & \cdots & | \\ \boxed{Au_1} & Au_2 & \cdots & Au_n \\ | & | & \cdots & | \end{bmatrix} = \begin{bmatrix} | & | & \cdots & | \\ \boxed{\lambda_1 u_1} & \lambda_2 u_2 & \cdots & \lambda_n u_n \\ | & | & \cdots & | \end{bmatrix} = U \text{diag}(\lambda_1, \dots, \lambda_n) = U\Lambda$$

$$\boxed{AU = U\Lambda}$$

New Representation for Symmetric Matrices

- Let U be the orthonormal matrix that contains u_i 's as columns:

$$U = \begin{bmatrix} | & | & \cdots & | \\ u_1 & u_2 & \cdots & u_n \\ | & | & & | \end{bmatrix}$$

- Let $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ be the diagonal matrix that contains $\lambda_1, \dots, \lambda_n$.

$$AU = \begin{bmatrix} | & | & \cdots & | \\ Au_1 & Au_2 & \cdots & Au_n \\ | & | & & | \end{bmatrix} = \begin{bmatrix} | & | & \cdots & | \\ \lambda_1 u_1 & \lambda_2 u_2 & \cdots & \lambda_n u_n \\ | & | & & | \end{bmatrix} = U \text{diag}(\lambda_1, \dots, \lambda_n) = U\Lambda$$

- Recalling that orthonormal matrix U satisfies that $UU^T = I$, we can diagonalize matrix A :

$$A = AUU^T = U\Lambda U^T \quad (4)$$

Background: representing vector w.r.t. another basis

- Any orthonormal matrix $U = \begin{bmatrix} | & | & & | \\ u_1 & u_2 & \cdots & u_n \\ | & | & & | \end{bmatrix}$ defines a new basis of \mathbb{R}^n . *changing the coordinate axes!*
- For any vector $x \in \mathbb{R}^n$ can be represented as a linear combination of u_1, \dots, u_n with coefficient $\hat{x}_1, \dots, \hat{x}_n$:

$$x = \hat{x}_1 u_1 + \cdots + \hat{x}_n u_n = U \hat{x}$$

$$x = U \hat{x}$$

- Indeed, such \hat{x} uniquely exists

$$x = U \hat{x} \Leftrightarrow U^T x = \hat{x}$$

$$\Rightarrow \hat{x} = U^T x$$

In other words, the vector $\hat{x} = U^T x$ can serve as another representation of the vector x w.r.t the basis defined by U .

“Diagonalizing” matrix-vector multiplication

- Left-multiplying matrix A can be viewed as left-multiplying a diagonal matrix w.r.t the basis of the eigenvectors.
 - ▶ Suppose x is a vector and \hat{x} is its representation w.r.t to the basis of U . $\hat{x} = U^T x$
 - ▶ Let $z = Ax$ be the matrix-vector product.
 - ▶ the representation z w.r.t the basis of U :

$$\hat{z} = U^T z = U^T Ax = \underbrace{U^T U}_{I} \underbrace{\Lambda U^T}_{\hat{x}} x = \Lambda \hat{x} = \begin{bmatrix} \lambda_1 \hat{x}_1 \\ \lambda_2 \hat{x}_2 \\ \vdots \\ \lambda_n \hat{x}_n \end{bmatrix}$$

complicated matrix-vector product \rightarrow scaling!

- We see that left-multiplying matrix A in the original space is equivalent to left-multiplying the diagonal matrix Λ w.r.t the new basis, which is merely scaling each coordinate by the corresponding eigenvalue.

“Diagonalizing” matrix-vector multiplication

Taking a power of a matrix

Under the new basis, multiplying a matrix multiple times becomes much simpler as well. For example, suppose $q = A^3 x$.

$$\hat{q} = U^T q = U^T A^3 x = U^T U \Lambda U^T U \Lambda U^T U \Lambda U^T x = \Lambda^3 \hat{x} = \begin{bmatrix} \lambda_1^3 \hat{x}_1 \\ \lambda_2^3 \hat{x}_2 \\ \vdots \\ \lambda_n^3 \hat{x}_n \end{bmatrix}$$

$$A^{200} x = (\Lambda^{200}) x$$

“Diagonalizing” quadratic form

As a direct corollary, the quadratic form $x^T Ax$ can also be simplified under the new basis

$$x^T Ax = x^T U \Lambda U^T x = \hat{x}^T \Lambda \hat{x} = \sum_{i=1}^n \lambda_i \hat{x}_i^2$$

(Recall that with the old representation, $x^T Ax = \sum_{i=1, j=1}^n x_i x_j A_{ij}$ involves a sum of n^2 terms instead of n terms in the equation above.)

The definiteness of the matrix A depends entirely on the sign of its eigenvalues

1. If all $\lambda_i > 0$ then the matrix A is positive definite because $x^T Ax = \sum_{i=1}^n \lambda_i \hat{x}_i^2 > 0$ for any $\hat{x} \neq 0$.¹
2. If all $\lambda_i \geq 0$, it is positive semidefinite because $x^T Ax = \sum_{i=1}^n \lambda_i \hat{x}_i^2 \geq 0$ for all \hat{x} .
3. Likewise, if all $\lambda_i < 0$ or $\lambda_i \leq 0$, then A is negative definite or negative semidefinite respectively.
4. Finally, if A has both positive and negative eigenvalues, say $\lambda_i > 0$ and $\lambda_j < 0$, then it is indefinite. This is because if we let \hat{x} satisfy $\hat{x}_i = 1$ and $\hat{x}_k = 0, \forall k \neq i$, then $x^T Ax = \sum_{i=1}^n \lambda_i \hat{x}_i^2 > 0$. Similarly we can let \hat{x} satisfy $\hat{x}_j = 1$ and $\hat{x}_k = 0, \forall k \neq j$, then $x^T Ax = \sum_{i=1}^n \lambda_i \hat{x}_i^2 < 0$.

¹Note that $\hat{x} \neq 0 \Leftrightarrow x \neq 0$.

Outline

- 1 Basic Concepts and Notation
- 2 Matrix Multiplication
- 3 Operations and Properties
- 4 Matrix Calculus**

Matrix Calculus

The Gradient

Note that the size of $\nabla_A f(A)$ is always the same as the size of A . So if, in particular, A is just a vector $x \in \mathbb{R}^n$,

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\nabla_x f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{bmatrix} \cdot \left[\nabla_x \right]_i = \frac{\partial f}{\partial x_i}$$

$(n \times 1)$ (under $\nabla_x f(x)$)

$(n \times 1)$ (under the vector in the product)

The Gradient

Note that the size of $\nabla_A f(A)$ is always the same as the size of A . So if, in particular, A is just a vector $x \in \mathbb{R}^n$,

$$\nabla_x f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{bmatrix}.$$

It follows directly from the equivalent properties of partial derivatives that:

- $\nabla_x(f(x) + g(x)) = \nabla_x f(x) + \nabla_x g(x)$.
 - For $t \in \mathbb{R}$, $\nabla_x(t f(x)) = t \nabla_x f(x)$.
- } Prove!

Analogous to 2nd derivative

The Hessian

(NOT matrix) vector scalar

Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a function that takes a vector in \mathbb{R}^n and returns a real number.

Then the Hessian matrix with respect to x , written $\nabla_x^2 f(x)$ or simply as H is the $n \times n$ matrix of partial derivatives,

$x : n \times 1$

$H_x : n \times n$

$f : \mathbb{R}^n \rightarrow \mathbb{R}$

$$\nabla_x^2 f(x) \in \mathbb{R}^{n \times n} = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1^2} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_2^2} & \cdots & \frac{\partial^2 f(x)}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_n \partial x_1} & \frac{\partial^2 f(x)}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_n^2} \end{bmatrix}.$$

$H_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j}$

In other words, $\nabla_x^2 f(x) \in \mathbb{R}^{n \times n}$, with

$$\boxed{(\nabla_x^2 f(x))_{ij} = \frac{\partial^2 f(x)}{\partial x_i \partial x_j}}.$$

The Hessian

Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a function that takes a vector in \mathbb{R}^n and returns a real number. Then the **Hessian** matrix with respect to x , written $\nabla_x^2 f(x)$ or simply as H is the $n \times n$ matrix of partial derivatives,

$$\nabla_x^2 f(x) \in \mathbb{R}^{n \times n} = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1^2} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_2^2} & \cdots & \frac{\partial^2 f(x)}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_n \partial x_1} & \frac{\partial^2 f(x)}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_n^2} \end{bmatrix}.$$

Jacobian is different
(if time permits)

Note that the Hessian is always symmetric, since

$$\frac{\partial^2 f(x)}{\partial x_i \partial x_j} = \frac{\partial^2 f(x)}{\partial x_j \partial x_i}.$$

Gradients of Linear Functions

$$f(x) = \sum b_i x_i$$

$$b: n \times 1$$

$$x: n \times 1$$

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$

For $x \in \mathbb{R}^n$, let $f(x) = b^T x$ for some known vector $b \in \mathbb{R}^n$. Then

$$f(x) = \sum_{i=1}^n b_i x_i$$

$$\left(\nabla_x f \right)_i = \frac{\partial f}{\partial x_i} = b_i$$

so

$$\frac{\partial f(x)}{\partial x_k} = \frac{\partial}{\partial x_k} \sum_{i=1}^n b_i x_i = b_k.$$

From this we can easily see that $\nabla_x b^T x = b$. This should be compared to the analogous situation in single variable calculus, where $\partial/(\partial x) ax = a$.

Gradients of Quadratic Function

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$

Now consider the quadratic function $f(x) = x^T Ax$ for $A \in \mathbb{S}^n$. Remember that

$$f(x) = \sum_{i=1}^n \sum_{j=1}^n A_{ij} x_i x_j. \quad \left(\nabla_x f\right)_i = \frac{\partial f}{\partial x_i}$$

To take the partial derivative, we'll consider the terms including x_k and x_k^2 factors separately:

$$\frac{\partial f(x)}{\partial x_k} = \frac{\partial}{\partial x_k} \sum_{i=1}^n \sum_{j=1}^n A_{ij} x_i x_j$$

Gradients of Quadratic Function

Now consider the quadratic function $f(x) = x^T Ax$ for $A \in \mathbb{S}^n$. Remember that

$$f(x) = \sum_{i=1}^n \sum_{j=1}^n A_{ij} x_i x_j.$$

To take the partial derivative, we'll consider the terms including x_k and x_k^2 factors separately:

$$\begin{aligned} \frac{\partial f(x)}{\partial x_k} &= \frac{\partial}{\partial x_k} \sum_{i=1}^n \sum_{j=1}^n A_{ij} x_i x_j \\ &= \frac{\partial}{\partial x_k} \left[\sum_{i \neq k} \sum_{j \neq k} A_{ij} x_i x_j + \sum_{i \neq k} A_{ik} x_i x_k + \sum_{j \neq k} A_{kj} x_k x_j + A_{kk} x_k^2 \right] \end{aligned}$$

Gradients of Quadratic Function

Now consider the quadratic function $f(x) = x^T Ax$ for $A \in \mathbb{S}^n$. Remember that

$$f(x) = \sum_{i=1}^n \sum_{j=1}^n A_{ij} x_i x_j.$$

To take the partial derivative, we'll consider the terms including x_k and x_k^2 factors separately:

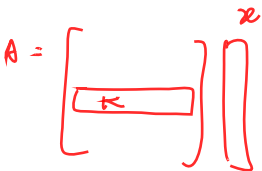
$$\begin{aligned} \frac{\partial f(x)}{\partial x_k} &= \frac{\partial}{\partial x_k} \sum_{i=1}^n \sum_{j=1}^n A_{ij} x_i x_j \\ &= \frac{\partial}{\partial x_k} \left[\cancel{\sum_{i \neq k} \sum_{j \neq k} A_{ij} x_i x_j} + \sum_{i \neq k} \underbrace{A_{ik} x_i x_k}_{\text{red wavy}} + \sum_{j \neq k} \underbrace{A_{kj} x_k x_j}_{\text{red wavy}} + \underbrace{A_{kk} x_k^2}_{\text{red wavy}} \right] \\ &= \sum_{i \neq k} A_{ik} x_i + \sum_{j \neq k} A_{kj} x_j + 2A_{kk} x_k \end{aligned}$$

Gradients of Quadratic Function

Now consider the quadratic function $f(x) = x^T Ax$ for $A \in \mathbb{S}^n$. Remember that

$$f(x) = \sum_{i=1}^n \sum_{j=1}^n A_{ij} x_i x_j.$$

To take the partial derivative, we'll consider the terms including x_k and x_k^2 factors separately:



$A = \left[\begin{array}{c} \dots \\ \boxed{k} \\ \dots \end{array} \right]$ $x = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$

$$\frac{\partial f(x)}{\partial x_k} = \frac{\partial}{\partial x_k} \sum_{i=1}^n \sum_{j=1}^n A_{ij} x_i x_j$$

$$= \sum_{i \neq k} A_{ik} x_i + \sum_{j \neq k} A_{kj} x_j + 2A_{kk} x_k$$

$$= \sum_{i=1}^n A_{ik} x_i + \sum_{j=1}^n A_{kj} x_j = 2 \sum_{i=1}^n A_{ki} x_i$$

$\therefore \nabla_x f = 2Ax$

$\therefore A$ is symmetric

Hessian of Quadratic Functions

Finally, let's look at the Hessian of the quadratic function $f(x) = x^T A x$

In this case,

$$\frac{\partial^2 f(x)}{\partial x_k \partial x_\ell} = \frac{\partial}{\partial x_k} \left[\frac{\partial f(x)}{\partial x_\ell} \right] = \frac{\partial}{\partial x_k} \left[2 \sum_{i=1}^n A_{\ell i} x_i \right] = 2A_{\ell k} = 2A_{k\ell}.$$

Therefore, it should be clear that $\nabla_x^2 x^T A x = 2A$, which should be entirely expected (and again analogous to the single-variable fact that $\frac{\partial^2}{(\partial x^2)} ax^2 = 2a$).

Recap

- $\nabla_x \underline{b^T x} = \underline{b} = (b^T)^T$
- $\nabla_x^2 b^T x = 0$
- $\nabla_x x^T Ax = 2Ax$ (if A symmetric)
- $\nabla_x^2 x^T Ax = 2A$ (if A symmetric)

Matrix Calculus Example: Least Squares

- Given a full rank matrix $A \in \mathbb{R}^{m \times n}$, and a vector $b \in \mathbb{R}^m$ such that $b \notin \mathcal{R}(A)$, we want to find a vector x such that Ax is as close as possible to b , as measured by the square of the Euclidean norm $\|Ax - b\|_2^2$. (ℓ_2)

$$\operatorname{arg\,min}_x \|Ax - b\|_2^2$$

Matrix Calculus Example: Least Squares

- Given a full rank matrix $A \in \mathbb{R}^{m \times n}$, and a vector $b \in \mathbb{R}^m$ such that $b \notin \mathcal{R}(A)$, we want to find a vector x such that Ax is as close as possible to b , as measured by the square of the Euclidean norm $\|Ax - b\|_2^2$.
- Using the fact that $\|x\|_2^2 = x^T x$, we have

$$\|Ax - b\|_2^2 = \underbrace{(Ax - b)}^T \underbrace{(Ax - b)} = \underbrace{x^T A^T A x} - 2b^T A x + b^T b$$

