

Optimization Techniques for Machine Learning

AMLZC326 · #01 Analytic Geometry

Anshid Aboobacker

MOTIVATION

- Mathematics is the language of Machine Learning.
- You have already gone through
 - ▶ Linear Algebra
 - ▶ Calculus
 - ▶ Discrete Maths
 - ▶ Probability & Statistics
- This course will take you through **Optimization**

LEARNING OBJECTIVES

Before optimizing functions, we must understand the space itself.
In this session we will:

- Recapitulation of vector spaces and bases
- Work with norms and inner products
- Interpret angles and orthogonality
- Recognize orthonormal bases and matrices

TABLE OF CONTENTS

- 1 Vector Spaces & Norms
- 2 Inner Products
- 3 Orthogonality & Gram-Schmidt

VECTOR SPACES

A vector space $(V, +, \cdot)$ over real numbers is a set V with two operations: Vector addition $+ : V \times V \rightarrow V$, and Scalar multiplication $\cdot : \mathbb{R} \times V \rightarrow V$, which satisfies the following:

- 1 $\mathbf{x} + \mathbf{y} = \mathbf{y} + \mathbf{x}$ (Commutativity)
- 2 $(\mathbf{x} + \mathbf{y}) + \mathbf{z} = \mathbf{x} + (\mathbf{y} + \mathbf{z})$ (Associativity)
- 3 $\mathbf{x} + \mathbf{0} = \mathbf{x}$ (Existence of zero element)
- 4 $\mathbf{x} + (-\mathbf{x}) = \mathbf{0}$ (Existence of additive inverse)
- 5 $a \cdot (\mathbf{x} + \mathbf{y}) = a \cdot \mathbf{x} + a \cdot \mathbf{y}$
- 6 $(a + b) \cdot \mathbf{x} = a \cdot \mathbf{x} + b \cdot \mathbf{x}$
- 7 $(ab) \cdot \mathbf{x} = a \cdot (b \cdot \mathbf{x})$
- 8 $1 \cdot \mathbf{x} = \mathbf{x}$

EXAMPLES

- \mathbb{R}^n
- Polynomials
 $P_n = \{p(x) \mid p(x) \text{ is a polynomial of degree } \leq n\}$
- Matrices $\mathbb{R}^{m \times n}$

BASIS OF A VECTOR SPACE

A set $B = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ is a basis of V if:

- B is a *spanning* set: Every $\mathbf{x} \in V$ can be written as

$$\mathbf{x} = a_1 \mathbf{v}_1 + \dots + a_n \mathbf{v}_n$$

- B is a *linearly independent* set:

$$a_1 \mathbf{v}_1 + \dots + a_n \mathbf{v}_n = \mathbf{0} \Rightarrow a_1 = \dots = a_n = 0$$

Example: Standard basis for \mathbb{R}^3

$$[1, 0, 0]^T, [0, 1, 0]^T, [0, 0, 1]^T$$

BASIS IS NOT UNIQUE

In \mathbb{R}^2 , both of the following sets are bases:

- $B_1 = \{[1, 0]^T, [0, 1]^T\}$
- $B_2 = \{[1, 1]^T, [1, -1]^T\}$

DIMENSION

The dimension of V , written $\dim(V)$, is the number of vectors in any basis of V .

Examples:

- $\dim(\mathbb{R}^n) = n$
- $\dim(P_2) = 3$

NORMS

- A norm on a vector space is a function $\|\cdot\| : V \rightarrow \mathbb{R}$, $\mathbf{x} \rightarrow \|\mathbf{x}\|$ which assigns to each vector \mathbf{x} a length $\|\mathbf{x}\|$ such that for all $\lambda \in \mathbb{R}$ and $\mathbf{x}, \mathbf{y} \in V$ the following properties hold:
 - ▶ Absolutely homogeneous: $\|\lambda\mathbf{x}\| = |\lambda|\|\mathbf{x}\|$
 - ▶ Triangle inequality: $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$
 - ▶ Positive definite: $\|\mathbf{x}\| \geq 0$ and $\|\mathbf{x}\| = 0 \implies \mathbf{x} = \mathbf{0}$

EXAMPLES OF NORMS

Manhattan norm : $\|\mathbf{x}\|_1 = \sum_{i=1}^{i=n} |x_i|$

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

Spectral norm: For a matrix $A \in \mathbb{R}^{n \times m}$,

$$\|A\|_2 = \max_{\mathbf{x} \neq 0} \frac{\|A\mathbf{x}\|_2}{\|\mathbf{x}\|_2} = \sigma_{\max}(A)$$

where $\sigma_{\max}(A)$ is the largest singular value of A .

TABLE OF CONTENTS

- 1 Vector Spaces & Norms
- 2 Inner Products**
- 3 Orthogonality & Gram-Schmidt

INNER PRODUCTS

- Dot product in \mathbb{R}^n is given by $\mathbf{x}^T \mathbf{y} = \sum_{i=1}^{i=n} x_i y_i$
- A bilinear mapping Ω is a mapping with two arguments and is linear in both arguments: Let V be a vector space such that $\mathbf{x}, \mathbf{y}, \mathbf{z} \in V$, and let $\lambda, \psi \in \mathbb{R}$. Then we have $\Omega(\lambda \mathbf{x} + \psi \mathbf{y}, \mathbf{z}) = \lambda \Omega(\mathbf{x}, \mathbf{z}) + \psi \Omega(\mathbf{y}, \mathbf{z})$, and $\Omega(\mathbf{x}, \lambda \mathbf{y} + \psi \mathbf{z}) = \lambda \Omega(\mathbf{x}, \mathbf{y}) + \psi \Omega(\mathbf{x}, \mathbf{z})$.
- Let V be a vector space and $\Omega : V \times V \rightarrow \mathbb{R}$ be a bilinear mapping that takes two vectors as arguments and returns a real number. Then Ω is called symmetric if $\Omega(\mathbf{x}, \mathbf{y}) = \Omega(\mathbf{y}, \mathbf{x})$. Also Ω is called positive-definite if $\forall \mathbf{x} \in V \setminus \{0\}, \Omega(\mathbf{x}, \mathbf{x}) > 0$ and $\Omega(\mathbf{0}, \mathbf{0}) = 0$.

INNER PRODUCTS

- A positive-definite, symmetric bilinear mapping $\Omega : V \times V \rightarrow \mathbb{R}$ is called an inner product. To denote an inner product on V we generally write $\langle \mathbf{x}, \mathbf{y} \rangle$.
- The pair $(V, \langle \cdot, \cdot \rangle)$ is called an inner product space.
- Next we introduce the concept of symmetric, positive-definite matrices and show we can express an inner product using such matrices.
- We recall that in a vector space V any vector \mathbf{x} can be written as linear combination of the basis vectors. We use this to express an inner product in terms of a matrix.

SYMMETRIC, POSITIVE-DEFINITE MATRICES

- A matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is *symmetric* if $\mathbf{A} = \mathbf{A}^T$.
- An $n \times n$ symmetric real matrix \mathbf{A} is said to be *positive-definite* if $\mathbf{x}^T \mathbf{A} \mathbf{x} > 0$ for all non-zero $\mathbf{x} \in \mathbb{R}^n$.
- For a real-valued, finite-dimensional vector space V and an ordered basis B of V , it holds that $\langle \cdot, \cdot \rangle : V \times V \rightarrow \mathbb{R}$ is an inner product if and only if there exists a symmetric, positive definite matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ with $\langle \mathbf{x}, \mathbf{y} \rangle = \hat{\mathbf{x}}^T \mathbf{A} \hat{\mathbf{y}}$.

LENGTHS AND DISTANCES

- Inner products and norms are closely related in the sense that any inner product induces a norm: $\|\mathbf{x}\| = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}$
- Not every norm is induced by an inner product, for example the Manhattan norm.
- For an inner product vector space $(V, \langle \cdot, \cdot \rangle)$, the induced norm $\|\cdot\|$ satisfies the Cauchy-Schwarz inequality:
 $|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \|\mathbf{y}\|.$

METRIC SPACE

- Consider an inner product space $(V, \langle \cdot, \cdot \rangle)$. Define $d(\mathbf{x}, \mathbf{y})$ the distance between two vectors \mathbf{x} and \mathbf{y} to be
$$d(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\| = \sqrt{\langle \mathbf{x} - \mathbf{y}, \mathbf{x} - \mathbf{y} \rangle}.$$
- If we use the dot product as the inner product, then the distance is called the Euclidean distance.
- The mapping $d : V \times V \rightarrow \mathbb{R}$ is called a metric.

PROPERTIES OF A METRIC SPACE

A metric d has the following properties:

- d is positive-definite which means $d(\mathbf{x}, \mathbf{y}) \geq 0 \quad \forall \mathbf{x}, \mathbf{y} \in V$.
 $d(\mathbf{x}, \mathbf{y}) = 0 \implies \mathbf{x} = \mathbf{y}$.
- d is symmetric which means $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x}) \quad \forall \mathbf{x}, \mathbf{y} \in V$.
- d obeys the triangle inequality as follows:
 $d(\mathbf{x}, \mathbf{z}) \leq d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z}) \quad \forall \mathbf{x}, \mathbf{y}, \mathbf{z} \in V$

Inner products and metrics seem to be very similar in terms of their properties - however there is one important difference. When \mathbf{x} and \mathbf{y} are close to each other the inner product is large but the distance metric is small. On the other hand when \mathbf{x} and \mathbf{y} are far apart, then the inner product is small but the distance metric is large.

ANGLE BETWEEN TWO VECTORS

- In addition to being able to capture the lengths of vectors and the distance between vectors, inner products can also capture the **angle** ω between two vectors and can thus capture the geometry of a vector space.
- The key to using the inner product to characterize the angle between two vectors is the Cauchy-Schwarz inequality.
- Assume that \mathbf{x} and \mathbf{y} are not the $\mathbf{0}$ vector. Then the Cauchy-Schwarz inequality tells us that

$$-1 \leq \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|} \leq 1 \quad (1)$$

ANGLE BETWEEN TWO VECTORS

- Since the Cauchy-Schwarz ratio lies between -1 and 1 we can set it equal to the cosine of a unique angle $\omega \in [0, \pi]$ such that

$$\cos(\omega) = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{x}\| \|\mathbf{y}\|} \quad (2)$$

- The angle ω is the angle between two vectors.
- The notion of angle captures similarity of orientation between two vectors. When the dot product is close to zero, the vectors are more or less pointing in orthogonal directions and $\omega \approx \pi/2$.

TABLE OF CONTENTS

- 1 Vector Spaces & Norms
- 2 Inner Products
- 3 Orthogonality & Gram-Schmidt**

ORTHOGONALITY

- A key feature of the inner product is that we can use it to characterize vectors that are orthogonal.
- Two vectors \mathbf{x} and \mathbf{y} are orthogonal if and only if the inner product between them is 0. For an orthogonal pair of vectors \mathbf{x} , \mathbf{y} we can write $\mathbf{x} \perp \mathbf{y}$.
- If additionally $\|\mathbf{x}\| = 1 = \|\mathbf{y}\|$, then \mathbf{x} and \mathbf{y} are said to be *orthonormal*.
- By the above definition the $\mathbf{0}$ -vector is orthogonal to all vectors.
- Vectors which are orthogonal with respect to one inner product need not be orthogonal with respect to another inner product.

EXAMPLE

- Consider the vectors $\mathbf{x} = [1, 1]^T$ and $\mathbf{y} = [-1, 1]^T$
- With respect to the inner product defined as a dot product we see that $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^T \mathbf{y} = 1 * -1 + 1 * 1 = 0$.
- With respect to the inner product $\mathbf{x}^T \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \mathbf{y}$, the angle between the two vectors \mathbf{x} and \mathbf{y} becomes

$$\cos(\omega) = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{x}\| \|\mathbf{y}\|}$$

EXAMPLE



$$\begin{aligned}\cos(\omega) &= \frac{\mathbf{x}^T \mathbf{A} \mathbf{y}}{\sqrt{\mathbf{x}^T \mathbf{A} \mathbf{x}} \sqrt{\mathbf{y}^T \mathbf{A} \mathbf{y}}} \\ &= \frac{2x_1y_1 + x_2y_2}{\sqrt{(2x_1^2 + x_2^2)(2y_1^2 + y_2^2)}} \\ &= \frac{-1}{3}\end{aligned}$$

where $\mathbf{A} = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$.

- Thus, with respect to the second inner product the vectors \mathbf{x} and \mathbf{y} are no longer orthogonal.

ORTHONORMAL MATRIX

- A square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is an orthogonal matrix if and only if its columns are orthonormal:

$$\mathbf{A}^T \mathbf{A} = \mathbf{I} = \mathbf{A} \mathbf{A}^T$$

i.e., $\mathbf{A}^T = \mathbf{A}^{-1}$

- Transformations by an orthonormal matrix preserve lengths. This can be seen as follows, using the dot product as the definition of the inner product:

$$\|\mathbf{Ax}\|^2 = (\mathbf{Ax})^T \mathbf{Ax} = \mathbf{x}^T \mathbf{A}^T \mathbf{Ax} = \mathbf{x}^T \mathbf{I} \mathbf{x} = \mathbf{x}^T \mathbf{x} = \|\mathbf{x}\|^2.$$

ORTHONORMAL MATRIX

- Also the angle between two vectors \mathbf{x} and \mathbf{y} does not change after transformation by an orthonormal matrix. This can be seen as follows:

$$\begin{aligned}\cos(\omega) &= \frac{(\mathbf{Ax})^T \mathbf{Ay}}{\|\mathbf{Ax}\| \|\mathbf{Ay}\|} \\ &= \frac{\mathbf{x}^T \mathbf{A}^T \mathbf{Ay}}{\|\mathbf{x}\| \|\mathbf{y}\|} \\ &= \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|}\end{aligned}$$

EXAMPLE

- 2D-rotation matrix: $\begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}$ where θ is the angle of rotation.

ORTHONORMAL BASIS

- We already looked at the concept of a basis of a vector space, and found that for the vector space \mathbb{R}^n we need n basis vectors.
- Our basis vectors needed only to be linearly independent - we can ensure linear independence by ensuring that our basis vectors point in different directions, so that a linear combination of $n - 1$ basis vectors cannot cancel out the n th basis vector.
- Now we will look at a special case of a basis where the vectors are all mutually orthogonal in the sense of the inner product, and each vector is of unit length. We call such a basis an orthonormal basis.

ORTHONORMAL BASIS

- Question: Can you immediately think of an orthonormal basis for \mathbb{R}^n ? Is an orthonormal basis for a vector space unique?
- Formal definition of an orthonormal basis: Consider an n -dimensional vector space V and n basis vectors $\{\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_n\}$. If it is true that $\forall i, j = 1, \dots, n, i \neq j \langle \mathbf{b}_i, \mathbf{b}_j \rangle = 0$ and $\langle \mathbf{b}_i, \mathbf{b}_i \rangle = 1$, then the basis is called an orthonormal basis.
- If the basis vectors are only mutually orthogonal but not of length unity, then we have an orthogonal basis.

GRAM-SCHMIDT PROCESS

- Given a set of basis vectors for a vector space, can we convert the given basis into an orthogonal basis?
- **Idea:** We construct new vectors step by step by *removing* from each vector its projections onto the previously constructed vectors.
- Starting with linearly independent vectors $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$:
 - ▶ Keep $\mathbf{u}_1 = \mathbf{v}_1$
 - ▶ Subtract from \mathbf{v}_2 its projection onto \mathbf{u}_1 to get

$$\mathbf{u}_2 = \mathbf{v}_2 - \frac{\langle \mathbf{v}_2, \mathbf{u}_1 \rangle}{\langle \mathbf{u}_1, \mathbf{u}_1 \rangle} \mathbf{u}_1$$

GRAM-SCHMIDT PROCESS

- Subtract from \mathbf{v}_k its projections onto $\mathbf{u}_1, \dots, \mathbf{u}_{k-1}$ to get

$$\mathbf{u}_k = \mathbf{v}_k - \sum_{j=1}^{k-1} \frac{\langle \mathbf{v}_k, \mathbf{u}_j \rangle}{\langle \mathbf{u}_j, \mathbf{u}_j \rangle} \mathbf{u}_j$$

- Each step ensures the new vector is orthogonal to all previous ones.
- Normalizing each \mathbf{u}_k gives an orthonormal basis:

$$\mathbf{e}_k = \frac{\mathbf{u}_k}{\|\mathbf{u}_k\|}$$

GRAM-SCHMIDT PROCESS: EXAMPLE

Consider the vectors in \mathbb{R}^2 : $\mathbf{v}_1 = [1, 2]^T$ $\mathbf{v}_2 = [3, 4]^T$

- Step 1: $\mathbf{u}_1 = \mathbf{v}_1 = [1, 2]^T$
- Step 2: Subtract the projection of \mathbf{v}_2 onto \mathbf{u}_1

$$\mathbf{u}_2 = \mathbf{v}_2 - \frac{\langle \mathbf{v}_2, \mathbf{u}_1 \rangle}{\langle \mathbf{u}_1, \mathbf{u}_1 \rangle} \mathbf{u}_1$$

$$= \begin{bmatrix} 3 \\ 4 \end{bmatrix} - \frac{(3)(1) + (4)(2)}{(1)^2 + (2)^2} \begin{bmatrix} 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 3 \\ 4 \end{bmatrix} - \frac{11}{5} \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

- Resulting orthogonal vectors: $\mathbf{u}_1 = [1, 2]^T$, $\mathbf{u}_2 = [\frac{4}{5}, -\frac{2}{5}]^T$

KEY TAKEAWAYS

- A vector space $(V, +, \cdot)$ is closed under addition and scalar multiplication; a basis provides unique coordinates for every vector in the space
- A norm $\| \cdot \|$ measures the length of a vector; an inner product $\langle \cdot, \cdot \rangle$ encodes both length and angle (and hence orthogonality)
- Gram-Schmidt orthogonalisation converts any linearly independent set into an orthonormal basis by iteratively subtracting projections
- These geometric foundations — norms, inner products, orthonormal bases — underlie every matrix decomposition and optimisation algorithm in the course

Thank you :)