# Analytic Geometry 1

Jo Ciucă

Australian National University

comp36706670@anu.edu.au
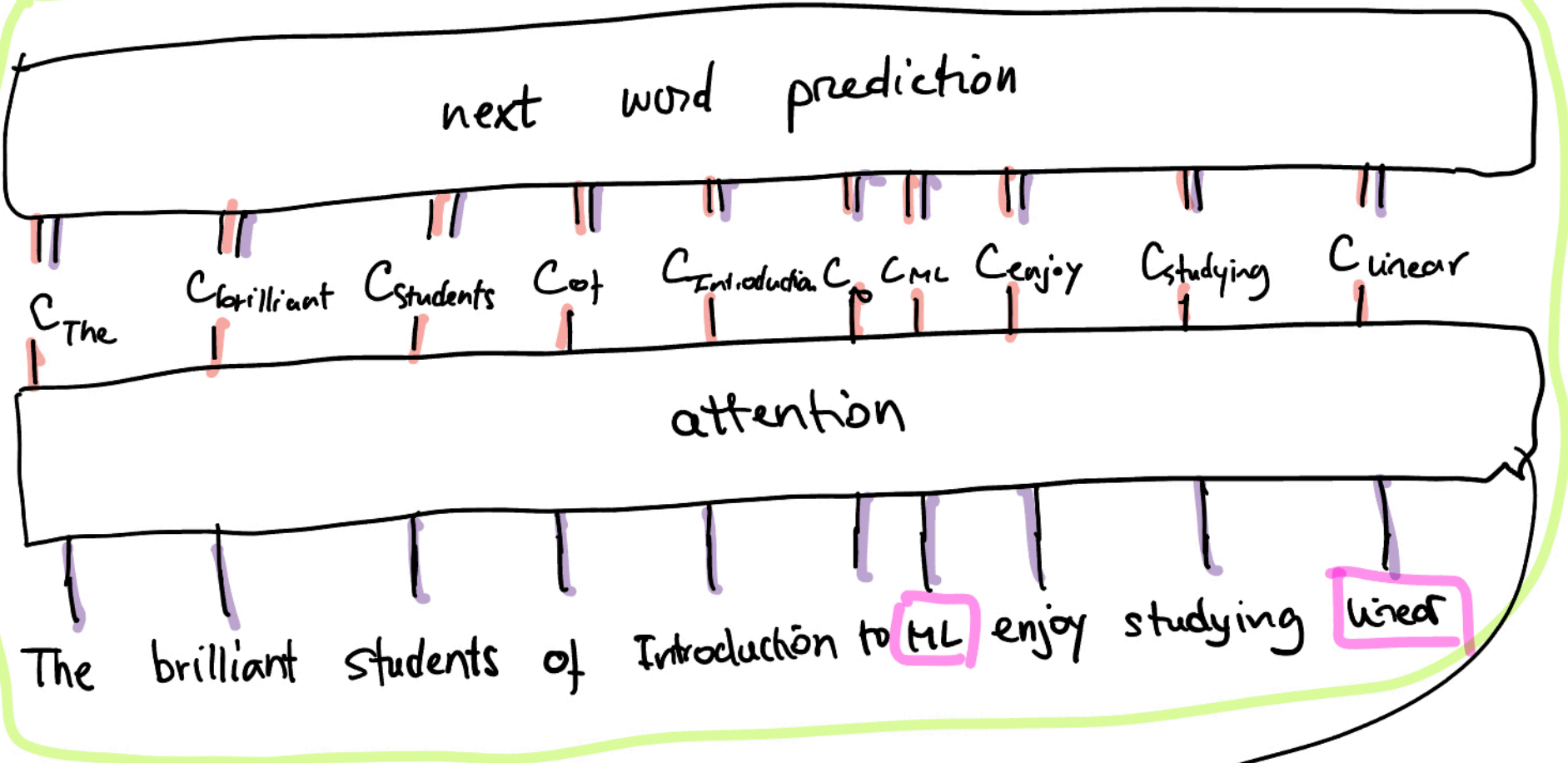
# Attention is all you need

The **brilliant students** of Introduction to **ML** enjoy **studying linear algebra.**

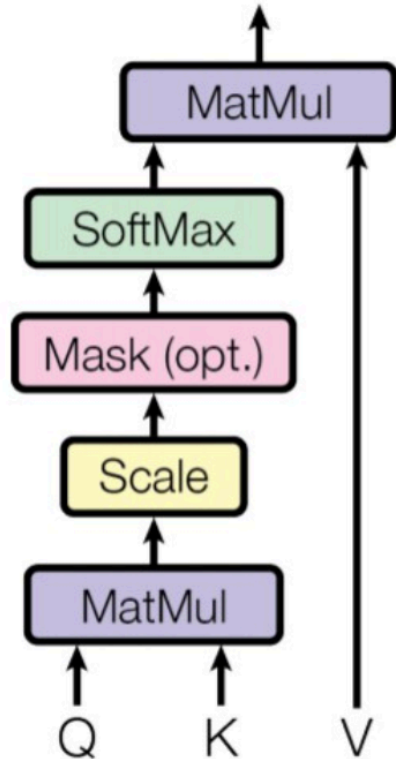# Causal Language Modelling Objective



algebra

Train

next word prediction

$C_{The}$  $C_{brilliant}$  $C_{students}$  $C_{of}$  $C_{Introduction}$  $C_{to}$  $C_{ML}$  $C_{enjoy}$  $C_{studying}$  $C_{linear}$

attention

The  brilliant  students  of  Introduction  to  ML  enjoy  studying  linear

$C$ : context vector

encodes how much each word pays attention to others

Credit: @g5min

# Linear Algebra 😍



$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V$$

Key

Query          Value

```python
import torch
import torch.nn.functional as F

def scaled_dot_product_attention(query, key, value):
    # Calculate the dot product between query and key
    scores = torch.matmul(query, key.transpose(-2, -1))

    # Scale the scores by square root of 'dk' (the dimension of the keys)
    dk = key.size(-1)  # get the size of the key's last dimension
    scaled_scores = scores / torch.sqrt(torch.tensor(dk).float())

    # Apply the softmax function to the scaled scores to get the attention weights
    attention_weights = F.softmax(scaled_scores, dim=-1)

    # Multiply the weights by the value vectors to get the output
    output = torch.matmul(attention_weights, value)

    return output, attention_weights

# Test the function
device = "cuda" if torch.cuda.is_available() else "cpu"
query = torch.randn(3, 8, device=device)
key = torch.randn(3, 8, device=device)
value = torch.randn(3, 8, device=device)

output, attention_weights = scaled_dot_product_attention(query, key, value)
print("Output shape: ", output.shape)
print("Attention weights shape: ", attention_weights.shape)
```
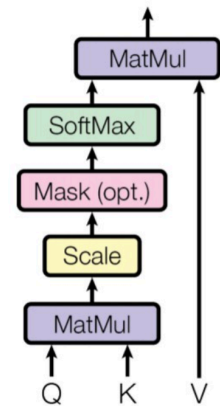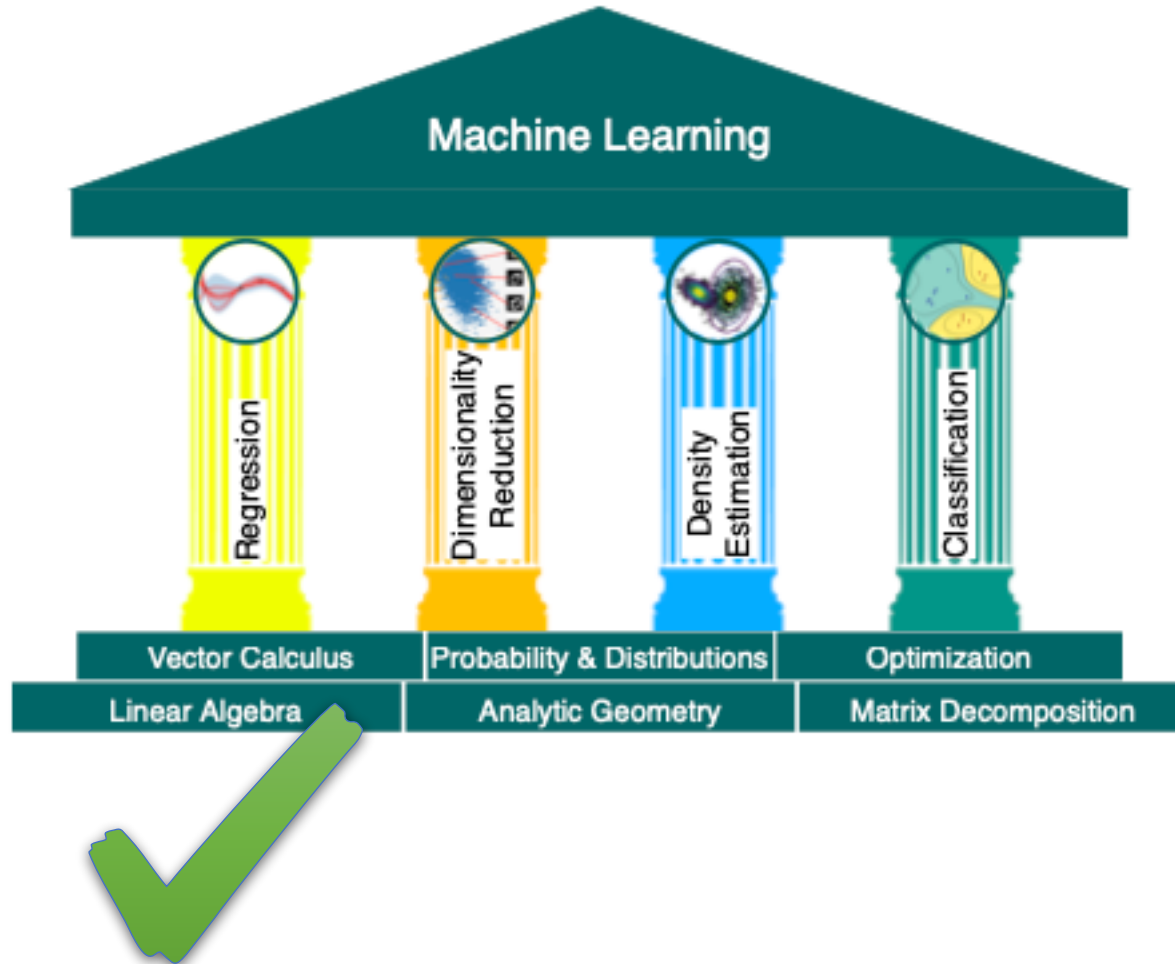
Thank you, next.

# Check your understanding

- Which of the following statements is correct?

(A) In a vector space, any vector can be represented as a linear combination of a certain set of vectors in this space.

(B) The dimension of a vector equals the dimension of the space it is in.

(C) $U$ is a vector subspace of $V$. Then vectors in $U$ have lower dimension than vectors in $V$.

(D) Set $\left\{ \begin{bmatrix} 2 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 2 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 7 \\ 7 \\ -1 \end{bmatrix} \right\}$ forms a basis for $\mathbb{R}^3$.

(E) $U = \left\{ (x, y) : x = y, \ x \in \mathbb{R}, \ y \in \mathbb{R} \right\}$ is a subspace of $\mathbb{R}^2$.

(F) The vector $\mathbf{0}$ is linearly dependent on any vector in the same vector space.

# Outline

- Bilinear Mappings
- Inner Product
- Lengths & distances
- Angles & Orthogonality

# 3.1 Norms

- A norm on a vector space $V$ is a function
$$\| \bullet \| : V \to \mathbb{R},$$

$$x \mapsto \|x\|,$$

  which assigns each vector $x$ its length $\|x\| \in \mathbb{R}$.

# Examples

- The Manhattan norm on $\mathbb{R}^n$ is defined for $x \in \mathbb{R}^n$ as
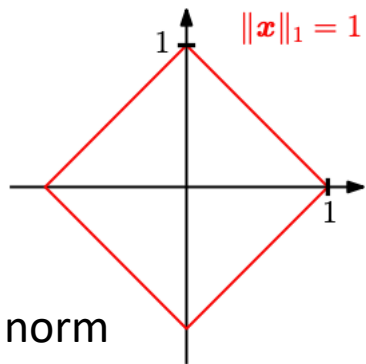
$$\|x\|_1 := \sum_{i=1}^{n} |x_i|,$$

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

where $|\bullet|$ is the absolute value. It is also called $\ell_1$ norm.
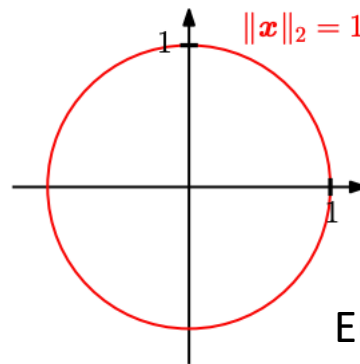
- The Euclidean norm of $x \in \mathbb{R}^n$ is defined as

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

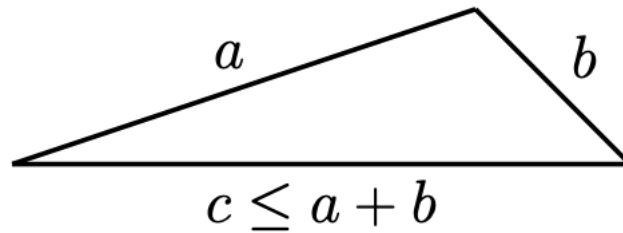It is the Euclidean distance of $x$ from the origin; also called $\ell_2$ norm



Manhattan norm

Euclidean norm

# 3.1 Norms

For all $\lambda \in \mathbb{R}$, and $x, y \in V$ the following holds:

- Absolutely homogeneous: $\|\lambda x\| = |\lambda| \|x\|$

- Triangle inequality: $\|x + y\| \leq \|x\| + \|y\|$

- Positive definite: $\|x\| \geq 0$ and $\|x\| = 0 \iff x = 0$



$$c \leq a + b$$

# 3.2. Inner products

Dot Product

- Scalar product/dot product in $\mathbb{R}^n$ is given by

$$\underset{1 \times n}{x^{\mathrm{T}}} \underset{n \times 1}{y} = \sum_{i=1}^{n} x_i y_i$$

# Bilinear mapping

- A bilinear mapping $\Omega$ is a mapping with two arguments, and it is linear in each argument. Consider a vector space $V$, for all $x, y, z \in V, \lambda, \varphi \in \mathbb{R}$,

$$\Omega(\lambda x + \varphi y, \ z) = \lambda \Omega(x, z) + \varphi \Omega(y, z)$$

$\Omega$ is linear in the first argument

$$\Omega(x, \lambda y + \varphi z) = \lambda \Omega(x, y) + \varphi \Omega(x, z).$$

$\Omega$ is linear in the second argument

13

# Inner product

- Let $V$ be a vector space and $\Omega: V \times V \to \mathbb{R}$ be a bilinear mapping.

- $\Omega$ is called symmetric if $\Omega(x, y) = \Omega(y, x)$

- $\Omega$ is called positive definite if

$$\forall x \in V \smallsetminus \{0\} : \Omega(x, x) > 0, \ \Omega(0, 0) = 0$$

- A positive definite, symmetric bilinear mapping $\Omega: V \times V \to \mathbb{R}$ is called an inner product on $V$. We write $\langle x, y \rangle$ instead of $\Omega(x, y)$.

- The pair $(V, \langle \bullet, \bullet \rangle)$ is called is called an inner product vector space. If we use the dot product, we call $(V, \langle \bullet, \bullet \rangle)$ a Euclidean vector space.

# Example

- Consider $V = \mathbb{R}^2$. If we define

$$\langle x, y \rangle := x_1 y_1 - \left( x_1 y_2 + x_2 y_1 \right) + 2 x_2 y_2$$

- then $\langle \bullet, \bullet \rangle$ is an inner product but different from the dot product.

This mapping is symmetric: it is easy to derive $\langle x, y \rangle = \langle y, x \rangle$
Is it positive definite?
$$\forall x \in V \setminus \{\mathbf{0}\}, \ \langle x, x \rangle = x_1^2 - \left( x_1 x_2 + x_2 x_1 \right) + 2 x_2^2 = \left( x_1 - x_2 \right)^2 + x_2^2 > 0$$

# 3.2.3 Symmetric, Positive Definite Matrices

- Consider an $n$-dimensional vector space $V$ with an inner product $\langle \bullet, \bullet \rangle$: $V \times V \to \mathbb{R}$, and a basis $B = \left( b_1, \cdots, b_n \right)$ of $V$.

$$\langle x, y \rangle = \left\langle \sum_{i=1}^{n} \varphi_i b_i, \sum_{j=1}^{n} \lambda_j b_j \right\rangle = \sum_{i=1}^{n} \sum_{j=1}^{n} \varphi_i \left\langle b_i, b_j \right\rangle \lambda_j = \hat{x}^T A \hat{y}$$

where $A_{ij} := \left\langle b_i, b_j \right\rangle$ and $\hat{x}$, $\hat{y}$ are the coordinates of $x, y$ with respect to the basis $B$.

- The inner product $\langle \bullet, \bullet \rangle$ is uniquely determined through $A$. The symmetry of the inner product also means that $A$ is symmetric.

- The positive definiteness of the inner product implies that

$$\forall x \in V \smallsetminus \left\{ 0 \right\} : \langle x, x \rangle = x^T A x > 0$$

# 3.2.3 Symmetric, Positive Definite Matrices

- A symmetric matrix $A \in \mathbb{R}^{n \times n}$ that satisfies $\forall x \in V \smallsetminus \{0\} : x^T A x > 0$ is called symmetric, positive definite, or just positive definite. If only $\geq$ holds, then $A$ is called symmetric, positive semidefinite.

- Example

$$A_1 = \begin{bmatrix} 3 & 1 \\ 1 & 4 \end{bmatrix}, \quad A_2 = \begin{bmatrix} 1 & 3 \\ 3 & 3 \end{bmatrix}$$

- $A_1$ is positive definite because it is symmetric and

$$x^T A_1 x = [x_1 \quad x_2] \begin{bmatrix} 3 & 1 \\ 1 & 4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$= 3x_1^2 + 2x_1 x_2 + 4x_2^2 = (x_1 + x_2)^2 + 2x_1^2 + 3x_2^2 > 0$$

for all $x \in V \smallsetminus \{0\}$.

- $A_2$ is symmetric but not positive definite

$$x^T A_2 x = x_1^2 + 6x_1 x_2 + 3x_2^2 = (x_1 + 3x_2)^2 - 6x_2^2 \text{ can be less than 0}$$

# 3.2.3 Symmetric, Positive Definite Matrices

- For a real-valued, finite-dimensional vector space $V$ and a basis $B$ of $V$, it holds that $\langle \bullet, \bullet \rangle : V \times V \to \mathbb{R}$ is an inner product if and only if there exists a symmetric, positive definite matrix $A \in \mathbb{R}^{n \times n}$ with

$$\langle x, y \rangle = \hat{x}^T A \hat{y}$$

- If $A \in \mathbb{R}^{n \times n}$ is symmetric and positive definite,

  the diagonal elements $a_{ii}$ of $A$ are positive because $a_{ii} = e_i^T A e_i = \langle e_i, e_i \rangle > 0$, where $e_i$ is the $i$th vector of the standard basis in $\mathbb{R}^n$.

# 3.3 Lengths and Distances

- Any inner product induces a norm

$$\|x\| := \sqrt{\langle x, x \rangle}$$

- Cauchy-Schwarz Inequality
- For an inner product vector space $(V, \langle \bullet, \bullet \rangle)$ the induced norm $\|\bullet\|$ satisfies the Cauchy-Schwarz inequality

$$\left| \langle x, y \rangle \right| \leq \|x\| \|y\|$$

# Example - Lengths of Vectors Using Inner Products

- We can now use an inner product to compute vector lengths, using $\|x\| := \sqrt{\langle x, x \rangle}$. Consider $x = [1,1]^{\mathrm{T}} \in \mathbb{R}^2$. If we use the dot product as the inner product, we obtain

$$\|x\| = \sqrt{x^T x} = \sqrt{1^2 + 1^2} = \sqrt{2}$$

$x^T \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} y$ is dot product

as the length of $x$. Let us now choose a different inner product:

$$\langle x, y \rangle := x^T \begin{bmatrix} 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 \end{bmatrix} y = x_1 y_1 - \frac{1}{2}(x_1 y_2 + x_2 y_1) + x_2 y_2$$

With this inner product, we obtain

$$\langle x, x \rangle = x_1^2 - x_1 x_2 + x_2^2 = 1 - 1 + 1 = 1 \implies \|x\| = \sqrt{1} = 1$$

$x$ is "shorter" with this inner product than with the dot product.

# 3.3 Lengths and Distances

- Consider an inner product space $\left(V, \langle \bullet, \bullet \rangle\right)$, then

$$d\left(x, y\right) := \|x - y\| = \sqrt{\langle x - y, x - y \rangle}$$

  is called the distance between $x$ and $y$ for $x, y \in V$.

- If we use the dot product as the inner product, then the distance is called Euclidean distance.

# 3.3 Lengths and Distances

- The mapping
$$d : V \times V \to \mathbb{R}$$
$$(\boldsymbol{x}, \boldsymbol{y}) \mapsto d(\boldsymbol{x}, \boldsymbol{y})$$

  is called a <span style="color:red">metric</span>.

- A metric $d$ satisfies the following:

- $d$ is positive definite, i.e., $d(\boldsymbol{x}, \boldsymbol{y}) \geq 0$ for all $\boldsymbol{x}, \boldsymbol{y} \in V$ and $d(\boldsymbol{x}, \boldsymbol{y}) = 0 \Leftrightarrow \boldsymbol{x} = \boldsymbol{y}$

- $d$ is symmetric, i.e., $d(\boldsymbol{x}, \boldsymbol{y}) = d(\boldsymbol{y}, \boldsymbol{x})$ for all $\boldsymbol{x}, \boldsymbol{y} \in V$

- Triangle inequality: $d(\boldsymbol{x}, \boldsymbol{z}) \leq d(\boldsymbol{x}, \boldsymbol{y}) + d(\boldsymbol{y}, \boldsymbol{z})$ for all $\boldsymbol{x}, \boldsymbol{y}, \boldsymbol{z} \in V$

- Very similar $\boldsymbol{x}$ and $\boldsymbol{y}$ will result in a <span style="color:red">large value for the inner product</span> and a <span style="color:red">small value for the metric</span>.
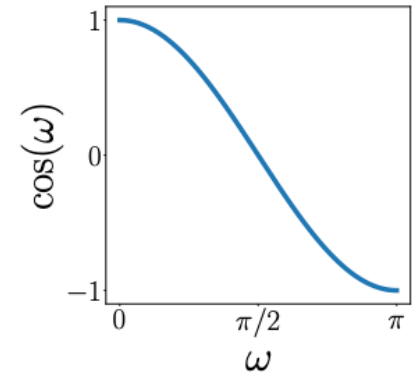
# 3.4 Angles and Orthogonality $\quad \left|\langle x, y \rangle\right| \leq \|x\|\|y\|$

- According to Cauchy-Schwarz inequality, assume $x \neq 0,\ y \neq 0$. Then,

$$-1 \leq \frac{\langle x, y \rangle}{\|x\|\|y\|} \leq 1$$

Therefore, there exists a unique $\omega \in \left[0, \pi\right]$, with

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



The number $\omega$ is the <span style="color:red">angle</span> between the vectors $x$ and $y$.

- The angle between two vectors tells us how similar their orientations are.
- Using the dot product, the angle between $x$ and $y = 4x$ is $0$, so their orientation is the same.

$$cos\omega = \frac{\langle x, 4x \rangle}{\|x\|\|4x\|} = \frac{4\langle x, x \rangle}{\sqrt{x^T x}\sqrt{(4x)^T(4x)}} = \frac{4\langle x, x \rangle}{4\|x\|\|x\|} = \frac{\langle x, x \rangle}{\|x\|\|x\|}$$

# Example (Angle between Vectors)

- Let us compute the angle between $x = \begin{bmatrix} -1, 1 \end{bmatrix}^{\mathrm{T}} \in \mathbb{R}^2$ and $y = \begin{bmatrix} 2, 1 \end{bmatrix}^{\mathrm{T}} \in \mathbb{R}^2$. We use the dot product as the inner product. We get

$$cos\omega = \frac{\langle x, y \rangle}{\sqrt{\langle x, x \rangle \langle y, y \rangle}} = \frac{x^T y}{\sqrt{x^T x y^T y}} = \frac{-1}{\sqrt{10}}$$

- and the angle between the two vectors is $\arccos\left(\frac{-1}{\sqrt{10}}\right) \approx 1.89\mathrm{rad}$, which corresponds to about $108.4°$.

- We then use inner product to characterize orthogonality.

# 3.4 Angles and Orthogonality

- Two vectors $x$ and $y$ are <span style="color:red">orthogonal</span> if and only if $\langle x, y \rangle = 0$, and we write $x \perp y$. If additionally $\|x\| = \|y\| = 1$, i.e., the vectors are unit vectors, then $x$ and $y$ are <span style="color:red">orthonormal</span>.

- $\mathbf{0}$-vector is orthogonal to every vector in the vector space

- Example (Orthogonal Vectors)

- Consider $x = \begin{bmatrix} 1, 2 \end{bmatrix}^{\mathrm{T}}$ and $y = \begin{bmatrix} -4, 2 \end{bmatrix}^{\mathrm{T}}$

- Using dot product as inner product, we have
  - $\langle x, y \rangle = 0$, so $x \perp y$.

- if we choose the inner product
$$\langle x, y \rangle = x^T \begin{bmatrix} 1 & 1 \\ 1 & 3 \end{bmatrix} y$$

- the angle $\omega$ between $x$ and $y$ is given by

$$cos\omega = \frac{\langle x, y \rangle}{\|x\| \, |y|} = -\frac{2}{\sqrt{17 \times 12}} \quad \Longrightarrow \quad \omega \approx 1.43\text{rad} \approx 81.95°$$

25

# 3.4 Angles and Orthogonality

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

$$A\, A^T = I = A^T A$$

which implies that

$$A^{-1} = A^T$$

i.e., the inverse is obtained by simply transposing the matrix

# Properties - length

- The length of a vector $x$ is not changed when transforming it using an orthogonal matrix $A$. For dot product, we obtain

$$\|Ax\|^2 = (Ax)^T(Ax) = x^T A^T A x = x^T I x = x^T x = \|x\|^2$$

# Properties - angle

- The angle between any two vectors $x$ and $y$ as measured by their inner product, is also unchanged when transforming both of them using an orthogonal matrix $A$. We use the dot product as inner product

$$cos\omega = \frac{(Ax)^T(Ay)}{\|Ax\|\|Ay\|} = \frac{x^T A^T A y}{\sqrt{x^T A^T A x\, y^T A^T A y}} = \frac{x^T y}{\|x\|\|y\|}$$

- Orthogonal matrices $A$ with $A^{-1} = A^T$ preserve both angles and distances.

- Orthogonal matrices define transformations that are rotations.