

# Mathematical Basics

---

Machine Learning 1 — Lecture 2

12<sup>th</sup> March 2024

Robert Peharz

Institute of Theoretical Computer Science

Graz University of Technology

“Mathematics” comes from Greek:

μάθημα: knowledge, study, **learning**

μαθηματικός: on the matter of that which is **learned**

“Mathematics” comes from Greek:

μάθημα: knowledge, study, **learning**

μαθηματικός: on the matter of that which is **learned**

Wait, is machine **learning** just. . . mathematics?

“Mathematics” comes from Greek:

μάθημα: knowledge, study, **learning**

μαθηματικός: on the matter of that which is **learned**

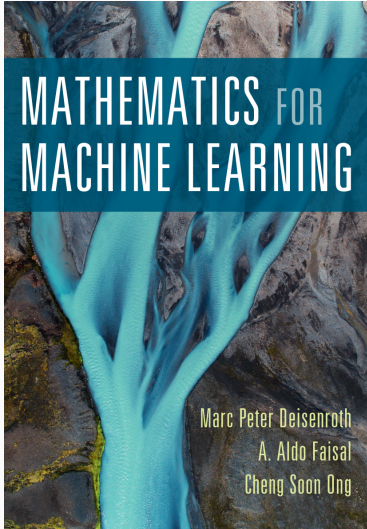
Wait, is machine **learning** just. . . mathematics?

Frankly, **yes!** Machine learning is mostly applied math.

# I do hope you will have a boring lecture. . .

- This lecture is intended as a **refresher** at a fairly quick pace, not as a complete introduction
- Most of today's lecture should be familiar to you
- If you severely struggle following today's lecture, the remainder of the course will be very challenging!

## Reading Material



I can recommend this excellent book, either for further reading or if you need to catch up.

<https://mml-book.github.io/>

The core mathematical disciplines in Machine Learning are:

- **Linear algebra**
- **Calculus**
- **Probability**

# Euclidean Vector Space

---



- The Euclidean space is the set of all  **$D$ -tuples**  $\mathbf{x}$  of real numbers:

$$\mathbb{R}^D = \mathbb{R} \times \mathbb{R} \times \cdots \times \mathbb{R} = \left\{ \begin{pmatrix} x_1 \\ \vdots \\ x_D \end{pmatrix} : x_1, \dots, x_D \in \mathbb{R} \right\}$$

- By default, we will interpret all vectors  $\mathbf{x} \in \mathbb{R}^D$  as **column-vectors**
- **Row-vectors** are denoted as  $\mathbf{x}^T$  ( $\mathbf{x}$  transposed)
- Thus, we also write  $\mathbf{x} = (x_1, \dots, x_D)^T$
- $D$  is the **dimensionality** of  $\mathbb{R}^D$

- We can **add** and **subtract** vectors  $\mathbf{x} = (x_1, x_2, \dots, x_D)^T$ ,  $\mathbf{y} = (y_1, y_2, \dots, y_D)^T$ :

$$\mathbf{x} + \mathbf{y} = (x_1 + y_1, x_2 + y_2, \dots, x_D + y_D)^T$$

$$\mathbf{x} - \mathbf{y} = (x_1 - y_1, x_2 - y_2, \dots, x_D - y_D)^T$$

- We can **scale** vectors, i.e. multiply with a scalar  $a$ :

$$a\mathbf{x} = (ax_1, ax_2, \dots, ax_D)^T$$

A **norm**  $\|\cdot\|$  defines a notion of length to a vector  $\mathbf{x} \in \mathbb{R}^D$ .

- **Manhattan norm**:  $\|\mathbf{x}\|_1 := \sum_i |x_i|$
- **Euclidean norm**:  $\|\mathbf{x}\|_2 := \sqrt{\sum_i x_i^2}$
- **Max norm**:  $\|\mathbf{x}\|_\infty = \max_i |x_i|$

These are special cases of the **p-norm**,  $p \geq 1$

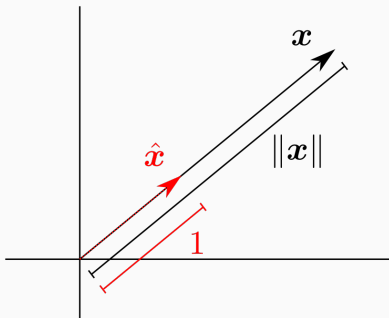
$$\|\mathbf{x}\|_p = \left( \sum_{i=1}^D |x_i|^p \right)^{1/p}$$

denoted  $\ell_p$ -norm,  $L_p$ -norm, etc.

A **unit vector** has length (norm) 1. We can **normalize** a vector  $\mathbf{x}$  by dividing by its norm:

$$\hat{\mathbf{x}} = \frac{1}{\|\mathbf{x}\|} \mathbf{x} = \frac{\mathbf{x}}{\|\mathbf{x}\|}$$

$\hat{\mathbf{x}}$  is the unit vector “showing in the same direction” as  $\mathbf{x}$ .



# Unit Spheres

**Unitspheres** (sets containing all unit vectors) for various  $p$ -norms:

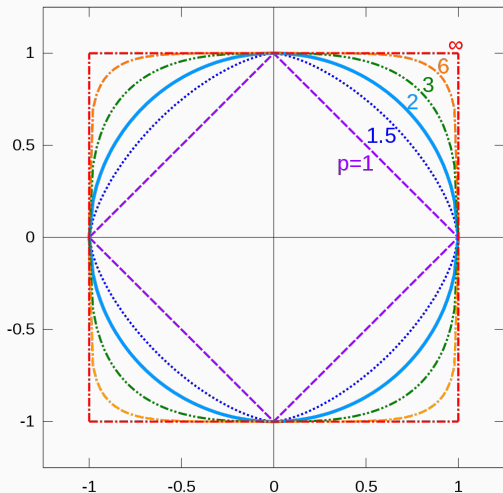


Image: wikipedia

As soon as we have defined a norm  $\|\cdot\|$ , we can define a **distance**  $d(\cdot, \cdot)$  between two vectors:

$$d(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|$$

Of course the distance depends on the used norm. The Euclidean norm  $\|\cdot\|_2$  measure distance differently than the max norm  $\|\cdot\|_\infty$ .

- For any two vectors  $\mathbf{x} = (x_1, \dots, x_D)^T$ ,  $\mathbf{y} = (y_1, \dots, y_D)^T$ , the **dot product** is defined as

$$\mathbf{x}^T \mathbf{y} := \sum_{i=1}^D x_i y_i$$

- E.g., for  $\mathbf{x} = (3, 1, 2.2)^T$  and  $\mathbf{y} = (4, 0, -1)^T$

$$\mathbf{x}^T \mathbf{y} = 3 \times 4 + 1 \times 0 + 2.2 \times (-1) = 9.8$$

- Note that the dot product and the Euclidean norm ( $\ell_2$ -norm) are related via

$$\|\mathbf{x}\|_2 = \sqrt{\mathbf{x}^T \mathbf{x}}$$

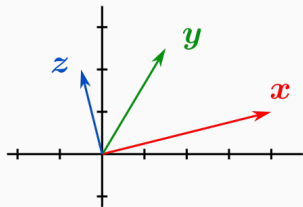
The dot products measures the **angle** between two vectors. In particular, the cosine of the angle  $\omega$  between two vectors  $\mathbf{x}$  and  $\mathbf{y}$  is given via:

$$-1 \leq \cos \omega =: \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\|_2 \|\mathbf{y}\|_2} \leq 1$$

If  $\mathbf{x}^T \mathbf{y} = 0 = \cos(\pi/2)$ , then  $\mathbf{x}$  and  $\mathbf{y}$  are **orthogonal**, written  $\mathbf{x} \perp \mathbf{y}$ .

If additionally,  $\|\mathbf{x}\|_2 = \|\mathbf{y}\|_2 = 1$ , then  $\mathbf{x}$  and  $\mathbf{y}$  are **orthonormal**.





$$\mathbf{x} = (4, 1)^T \quad \mathbf{y} = (1.5, 2.5)^T \quad \mathbf{z} = (-0.5, 2)^T$$

$$\frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|} = \frac{4 \times 1.5 + 1 \times 2.5}{4.123 \times 2.915} = 0.7071 = \cos\left(\frac{\pi}{4}\right) \triangleq 45^\circ$$

$$\frac{\mathbf{x}^T \mathbf{z}}{\|\mathbf{x}\| \|\mathbf{z}\|} = \frac{4 \times (-0.5) + 1 \times 2}{4.123 \times 2.062} = 0 = \cos\left(\frac{\pi}{2}\right) \triangleq 90^\circ$$

$$\frac{\mathbf{y}^T \mathbf{z}}{\|\mathbf{y}\| \|\mathbf{z}\|} = \frac{1.5 \times (-0.5) + 2.5 \times 2}{2.915 \times 2.062} = 0.7071 = \cos\left(\frac{\pi}{4}\right) \triangleq 45^\circ$$

# Vector Projection

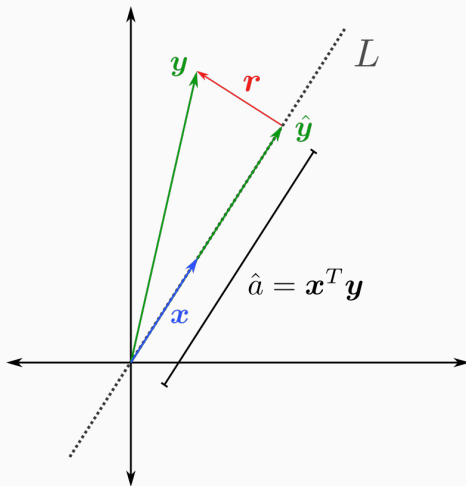
Assume a normalized vector  $\mathbf{x}$ . Then  $f(a) = a\mathbf{x}$  parametrizes a **line**  $L$  going through the origin ( $-\infty < a < \infty$ ).

Assume another vector  $\mathbf{y}$ , not necessarily normalized. The **projection**  $\hat{\mathbf{y}}$  of  $\mathbf{y}$  onto  $L$  is given as

$$\hat{\mathbf{y}} = (\mathbf{y}^T \mathbf{x}) \mathbf{x} = \hat{a} \mathbf{x}$$

where  $\hat{a} = \mathbf{y}^T \mathbf{x} = \mathbf{x}^T \mathbf{y}$  is given by the dot product.  $\hat{\mathbf{y}}$  is the **closest** vector (in Euclidean distance) to  $\mathbf{y}$  which lies on  $L$ .

$\mathbf{r} = \mathbf{y} - \hat{\mathbf{y}}$  is the **residual vector**, since  $\mathbf{y} = \hat{\mathbf{y}} + \mathbf{r}$ . Vectors  $\hat{\mathbf{y}}$  and  $\mathbf{r}$  are orthogonal.



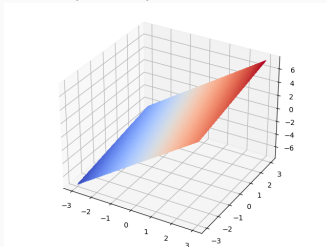
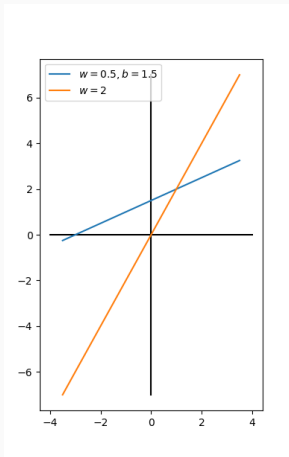
A vector  $\mathbf{w} \in \mathbb{R}^D$  defines a **linear function**  $f: \mathbb{R}^D \mapsto \mathbb{R}$  via the dot product:

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$$

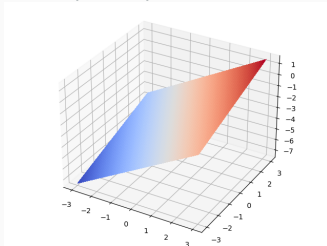
By adding a **bias**  $b$ , we get an **affine function**:

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$$

$$\mathbf{w} = (2, 0.5)^T, b = 0$$



$$\mathbf{w} = (1, 0.5)^T, b = -3$$



A vector  $\mathbf{w} \in \mathbb{R}^D$  and a bias  $b \in \mathbb{R}$  define an affine function  $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$ . The equation

$$f(\mathbf{x}) = 0$$

defines a **hyperplane** (a  $D - 1$  dimensional sub-space).

Moreover, two **half-spaces** are defined via

$$f(\mathbf{x}) < 0$$

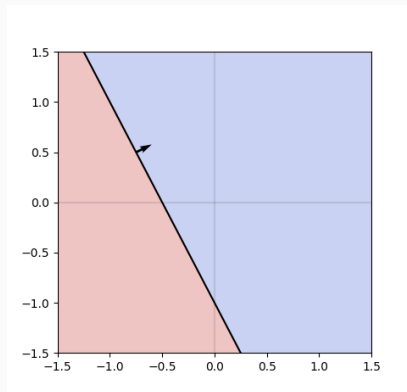
**negative half-space**

$$f(\mathbf{x}) \geq 0$$

**positive half-space**

The vector  $\mathbf{w}$  is **orthogonal** to the hyperplane and points into the positive half-space.

In 2-D, the hyperplane is a line, cutting  $\mathbb{R}^D$  in two parts. For  $\mathbf{w} = (2, 1)^T$  and  $b = 1$ :



blue: positive half-space  
red: negative half-space

In 3-D, a hyperplane is an actual 2-D plane. In higher dimensions it is a corresponding generalization, hence the name *hyperplane*.

# Matrices

---



A **matrix** is a rectangular array of real (or complex) numbers. With  $M, N \in \mathbb{N}$ , a matrix  $A \in \mathbb{R}^{M \times N}$  can be written as

$$A = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1N} \\ a_{21} & a_{22} & \dots & a_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ a_{M1} & a_{M2} & \dots & a_{MN} \end{pmatrix}$$

where  $M$  is the number of **rows** and  $N$  is the number of **columns**. The **entry**  $a_{ij}$  is a scalar located at the  $i^{\text{th}}$  row and  $j^{\text{th}}$  column.

Like vectors, matrices are added *element-wise*:

$$A + B = C = \begin{pmatrix} c_{11} & c_{12} & \dots & c_{1N} \\ c_{21} & c_{22} & \dots & c_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ c_{M1} & c_{M2} & \dots & c_{MN} \end{pmatrix}$$

where  $c_{ij} = a_{ij} + b_{ij}$ .

# Row and Column Vectors

A matrix can be understood as a collection of  $M$  **row-vectors**, or as a collection of  $N$  **column-vectors**:

$$A = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1N} \\ a_{21} & a_{22} & \dots & a_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ a_{M1} & a_{M2} & \dots & a_{MN} \end{pmatrix} = \begin{pmatrix} \mathbf{a}_{1:}^T \\ \mathbf{a}_{2:}^T \\ \vdots \\ \mathbf{a}_{M:}^T \end{pmatrix} = \begin{pmatrix} \mathbf{a}_{:1} & \mathbf{a}_{:2} & \dots & \mathbf{a}_{:N} \end{pmatrix}$$

where

$$\mathbf{a}_{i:}^T = (a_{i1}, a_{i2}, \dots, a_{iN})$$

**row-vectors**

$$\mathbf{a}_{:j} = (a_{1j}, a_{2j}, \dots, a_{Mj})^T$$

**column-vectors**

Note: we use the “colon-notation”  $\mathbf{a}_{i:}$ ,  $\mathbf{a}_{:j}$  to address rows and columns, similar as in python. Read “ $i$ -everything” or “everything- $j$ ”.

The **transpose**  $A^T$  of a matrix  $A$  is a matrix  $A^T$  whose rows are the columns of  $A$ , (or, whose columns are the rows of  $A$ ). For example:

$$A = \begin{pmatrix} 2 & 3 & 5 & 7 \\ 11 & 13 & 17 & 19 \\ 23 & 29 & 31 & 37 \end{pmatrix}$$

$$A^T = \begin{pmatrix} 2 & 11 & 23 \\ 3 & 13 & 29 \\ 5 & 17 & 31 \\ 7 & 19 & 37 \end{pmatrix}$$

A matrix  $A$  which has the same number of rows and columns is called a **square matrix**. For example:

$$A = \begin{pmatrix} 0 & 1 & 1 \\ 2 & 3 & 5 \\ 8 & 13 & 21 \end{pmatrix}$$

A square matrix  $A$  for which it holds that  $A = A^T$  is called **symmetric**. For example:

$$A = \begin{pmatrix} 3 & 1 & 7 \\ 1 & 2 & 9 \\ 7 & 9 & 1 \end{pmatrix}$$

A square matrix which contains only 0's, except on the main diagonal, is called a **diagonal** matrix. For example:

$$A = \begin{pmatrix} 3.1415 & 0 & 0 \\ 0 & -7 & 0 \\ 0 & 0 & 42 \end{pmatrix}$$

Of course, a diagonal matrix is symmetric.

**Matrix multiplication** is defined in terms of the **dot products** between row-vectors and column vectors. In particular let  $A \in \mathbb{R}^{M \times R}$  and  $B \in \mathbb{R}^{R \times N}$  (**note the same  $R$** ):

$$A = \begin{pmatrix} \mathbf{a}_{1:}^T \\ \vdots \\ \mathbf{a}_{M:}^T \end{pmatrix}, \quad B = \begin{pmatrix} \mathbf{b}_{:1} & \dots & \mathbf{b}_{:N} \end{pmatrix}$$

Then the **matrix product**  $C = AB$  is given as a matrix  $C \in \mathbb{R}^{M \times N}$

$$C = AB = \begin{pmatrix} c_{11} & c_{12} & \dots & c_{1N} \\ c_{21} & c_{22} & \dots & c_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ c_{M1} & c_{M2} & \dots & c_{MN} \end{pmatrix}$$

where  $c_{ij} = \mathbf{a}_{i:}^T \mathbf{b}_{:j} = \sum_{k=1}^R a_{ik} b_{kj}$



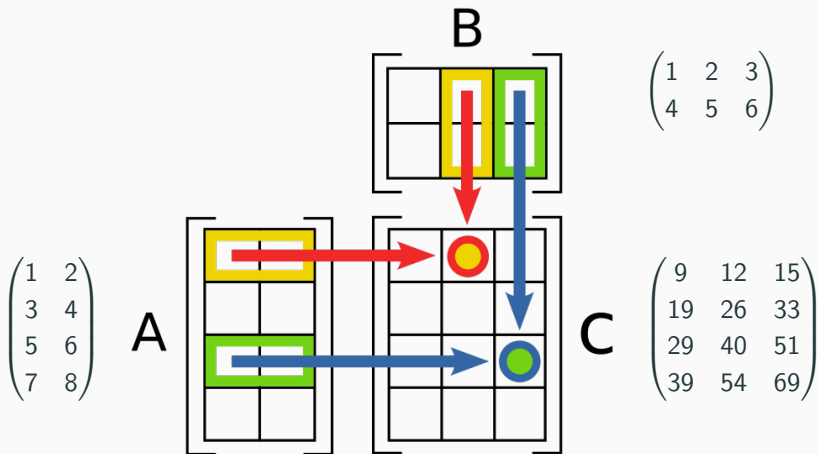


Image: wikipedia

# Properties of Matrix Product

For matrices  $A, B, C, D$  with compatible dimensions, the following properties hold:

- $(AB)C = A(BC)$  associativity
- $A(B + C) = AB + AC$  distributivity
- $(AB)^T = B^T A^T$
- $AA^T$  is always symmetric

Attention: matrix multiplication is **not commutative**, i.e. in general:

$$AB \neq BA$$

By interpreting vectors as “thin” matrices we can also multiply

- $A \in \mathbb{R}^{M \times N}$  and vector  $\mathbf{x} \in \mathbb{R}^N$ , yielding  $\mathbf{y} \in \mathbb{R}^M$ :

$$\mathbf{y} = A\mathbf{x}$$

- $\mathbf{x} \in \mathbb{R}^M$  and  $A \in \mathbb{R}^{M \times N}$ , yielding  $\mathbf{y} \in \mathbb{R}^N$ :

$$\mathbf{y}^T = \mathbf{x}^T A$$

# Matrices as Linear Functions

Recall that the **dot product** with vector  $\mathbf{w} \in \mathbb{R}^D$  defines a linear function  $f(\mathbf{x}): \mathbb{R}^D \mapsto \mathbb{R}$ ,  $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ .

Similarly, the **matrix-vector product** with matrix  $A \in \mathbb{R}^{M \times N}$  defines a **linear function**  $f: \mathbb{R}^N \mapsto \mathbb{R}^M$

$$f(\mathbf{x}) = A\mathbf{x}$$

We get an **affine function** by adding an  $M$ -dimensional bias vector  $\mathbf{b}$ :

$$f(\mathbf{x}) = A\mathbf{x} + \mathbf{b}$$

The diagonal matrix  $I$  which contains only 1's in the diagonal is called the **identity matrix**. There is an identity matrix for each  $M \in \mathbb{N}$ . For  $M = 3$ :

$$I = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

The identity matrix represents the **identity function**, since

$$I\mathbf{x} = \mathbf{x}$$

Let  $A$  be a square matrix. If there exists a matrix  $A^{-1}$  such that

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

then  $A^{-1}$  is called the **inverse (matrix)** of  $A$ .

If the inverse matrix exists, then it also holds that

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

The computational complexity of computing the inverse of a matrix  $A \in \mathbb{R}^{M \times M}$  is  $\mathcal{O}(M^3)$ .

# Linear Basis, Linear Subspace

---

Let  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K$  be  $K$  vectors in  $\mathbb{R}^D$ . A **linear combination** of these vectors yields another vector written as

$$\sum_{k=1}^K z_k \mathbf{x}_k$$

with coefficients  $z_k \in \mathbb{R}$ .

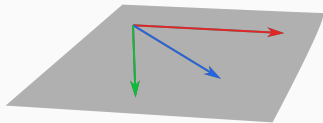
The vectors  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K$  are called **linearly independent**, when **none of them can be expressed as a linear combination of the other vectors**.

That is, for each  $1 \leq i \leq K$  and any coefficients  $z_k$

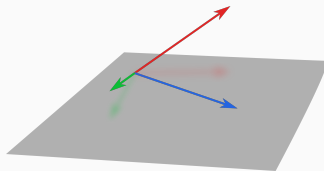
$$\mathbf{x}_i \neq \sum_{k=1, k \neq i}^K z_k \mathbf{x}_k$$



Linearly dependent vectors



Linearly independent vectors



Let  $\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_D$  be  $D$  **linearly independent vectors** in  $\mathbb{R}^D$ . Any such collection  $\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_D$  is called a **basis** of  $\mathbb{R}^D$ .

For any basis, it is possible to express any vector  $\mathbf{x} \in \mathbb{R}^D$  as a linear combination

$$\mathbf{x} = \sum_{d=1}^D \mathbf{b}_d z_d,$$

with **unique** coefficients  $z_d$ .

The coefficient vector  $\mathbf{z} = (z_1, \dots, z_D)^T$  represents  $\mathbf{x}$  in the basis  $\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_D$ , i.e. a **change of basis**.

How to find  $\mathbf{z}$ ?

Let  $\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_D$  be a basis of  $\mathbb{R}^D$ . We can collect the basis vectors as columns of a **basis matrix**  $B$ :

$$B = \begin{pmatrix} \mathbf{b}_1 & \mathbf{b}_2 & \dots & \mathbf{b}_D \end{pmatrix}$$

Then we can express any  $\mathbf{x} \in \mathbb{R}^D$  as

$$\mathbf{x} = \sum_{d=1}^D \mathbf{b}_d z_d = B\mathbf{z}.$$

Thus  $\mathbf{z}$  is given as:

$$B^{-1}\mathbf{x} = \mathbf{z}$$

Note:  $B^{-1}$  exists if and only if  $\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_D$  are a basis.

Things become particularly nice if the basis vectors are **orthonormal**, i.e. if they are

- normalized (i.e.  $\|\mathbf{b}_i\|_2 = 1$ )
- orthogonal

This can be compactly described as

$$\mathbf{b}_i^T \mathbf{b}_j = \delta_{ij} = \underbrace{\begin{cases} 1 & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}}_{\text{Kronecker delta}}$$

In this case  $B = \begin{pmatrix} \mathbf{b}_1 & \mathbf{b}_2 & \dots & \mathbf{b}_D \end{pmatrix}$  is called an **orthonormal matrix**.

# Properties of Orthonormal Matrices

Let  $B$  be an orthonormal matrix.

- Since  $\mathbf{b}_i^T \mathbf{b}_j = \delta_{ij}$  it holds that

$$B^T B = I$$

- Thus  $B^{-1} = B^T$
- Thus also  $BB^T = I$

## Properties of Orthonormal Matrices cont'd

Multiplication with  $B$  corresponds to a **rotation** of a vector.

Hence, multiplication with  $B$  (or  $B^T$ ) leaves the dot product (angle) between two vectors unchanged:

$$(B\mathbf{x})^T(B\mathbf{y}) = \mathbf{x}^T(B^TB)\mathbf{y} = \mathbf{x}^T\mathbf{y}$$

$$(B^T\mathbf{x})^T(B^T\mathbf{y}) = \mathbf{x}^T(BB^T)\mathbf{y} = \mathbf{x}^T\mathbf{y}$$

## Properties of Orthonormal Matrices cont'd

Also the (quadratic) Euclidean norm of a vector is unchanged:

$$\|B\mathbf{x}\|_2^2 = (B\mathbf{x})^T(B\mathbf{x}) = \mathbf{x}^T\mathbf{x} = \|\mathbf{x}\|_2^2$$

$$\|B^T\mathbf{x}\|_2^2 = (B^T\mathbf{x})^T(B^T\mathbf{x}) = \mathbf{x}^T\mathbf{x} = \|\mathbf{x}\|_2^2$$

# Change of Basis for Orthonormal Matrices

Recall from before that for any  $\mathbf{x} \in \mathbb{R}^D$  as

$$\mathbf{x} = B\mathbf{z}$$

$$B^{-1}\mathbf{x} = \mathbf{z}$$

For orthonormal matrices  $B^T = B^{-1}$ , thus

$$\mathbf{x} = B\mathbf{z}$$

$$B^T\mathbf{x} = \mathbf{z}$$



# Linear Subspace

So far we considered a **complete** basis:

$$\mathbf{b}_1, \dots, \mathbf{b}_D$$

describing the same vector space  $\mathbb{R}^D$ .

If we use only  $K < D$  basis vectors

$$\mathbf{b}_1, \dots, \mathbf{b}_K, \cancel{\mathbf{b}_{K+1}}, \dots, \cancel{\mathbf{b}_D}$$

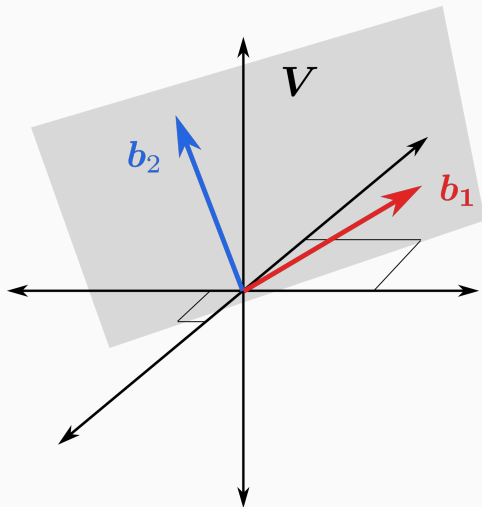
we describe a  $K$ -dimensional **linear subspace**.

Let  $B = (\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_K)$  be a  $D \times K$ -matrix with orthonormal vectors as columns, where  $K < D$ .

The **linear subspace** spanned by  $B$  is defined as all possible linear combinations of  $B$ 's columns:

$$\mathbf{V} = \left\{ \mathbf{v} = B\mathbf{z} \mid \mathbf{z} = (z_1, \dots, z_K)^T \in \mathbb{R}^K \right\}$$

Note that  $\mathbf{V}$  still contains tuples of length  $D$ , but it is a  $K$ -dimensional space, i.e. in one-to-one correspondence with  $\mathbb{R}^K$ .



2-dimensional subspace  $V$  of  $\mathbb{R}^3$  spanned by 2 vectors  $b_1$  and  $b_2$ .

Let  $B = (\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_K)$  be a matrix with  $K$  orthonormal vectors as columns, spanning a  $K$ -dimensional subspace  $\mathbf{V}$  of  $\mathbb{R}^D$ .

We can **project** an arbitrary vector  $\mathbf{x} \in \mathbb{R}^D$  onto  $\mathbf{V}$  by

- computing the **projection coefficients**  $B^T \mathbf{x} =: \mathbf{z} \in \mathbb{R}^K$
- computing the **projection/reconstruction**  $B\mathbf{z} =: \hat{\mathbf{x}} \in \mathbf{V}$
- thus,  $\hat{\mathbf{x}} = \underbrace{BB^T}_{\text{projection matrix}} \mathbf{x}$
- $\hat{\mathbf{x}}$  is the **closest point** in  $\mathbf{V}$  (in Euclidean distance) to  $\mathbf{x}$
- the **residual**  $\mathbf{r} = \mathbf{x} - \hat{\mathbf{x}}$  is always orthogonal to  $\hat{\mathbf{x}}$

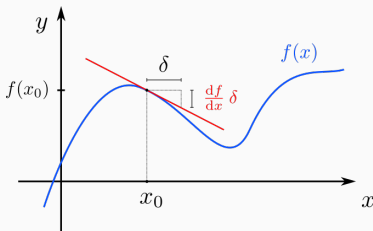
# Differential Calculus

---

Let  $f: \mathbb{R} \mapsto \mathbb{R}$  be a univariate function. If it exists, the **derivative** at a point  $x$  is defined as

$$\frac{df}{dx}(x) = f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

- rate of change
- slope of the tangent line at some point  $x_0$ , i.e. best local linear approximation:  $f(x) \approx f_{linear}(x) = \underbrace{f(x_0)}_b + \underbrace{f'(x_0)}_a \underbrace{(x - x_0)}_\delta$



# Standard Rules for Derivatives

Let  $f$  and  $g$  be differentiable univariate functions and  $a, b$  arbitrary constants.

- **Derivative is linear:**

$$(af + bg)'(x) = af'(x) + bg'(x)$$
$$\left( \frac{d(af + bg)}{dx} = a \frac{df}{dx} + b \frac{dg}{dx} \right)$$

- **Product rule:**

$$(fg)'(x) = f'g(x) + fg'(x)$$
$$\left( \frac{d fg}{dx} = \frac{df}{dx} g + \frac{dg}{dx} f \right)$$

- **Chain rule:**

$$f(g(x))' = f'(g(x))g'(x)$$
$$\left( \frac{df \circ g}{dx} = \frac{df}{dg} \frac{dg}{dx} \right)$$

Note:  $f \circ g$  denotes function composition:  $(f \circ g)(x) = f(g(x))$

## Derivatives of some well-known functions

Name	$f(x)$	$f'(x)$
Constant	$a$	$0$
Affine	$ax + b$	$a$
Polynomial	$x^k$	$k x^{k-1}$
Exponential	$\exp(x), e^x$	$\exp(x), e^x$
	$b^x$	$b^x \log b$
Logarithm	$\log x$	$\frac{1}{x}$
Sine	$\sin(x)$	$\cos(x)$
Cosine	$\cos(x)$	$-\sin(x)$

Note:  $\log$  denotes the **natural logarithm** in this course.