

Análisis Topológico de Datos

* <https://github.com/joperca/TDA-Fall2025>

Conjuntos de Datos (X, d_X) :

a) Un conjunto X (finito o no)

ejemplos: Vectores, imágenes, sonidos, videos,
documentos (texto), moléculas, etc.

b) Una función (distancia)

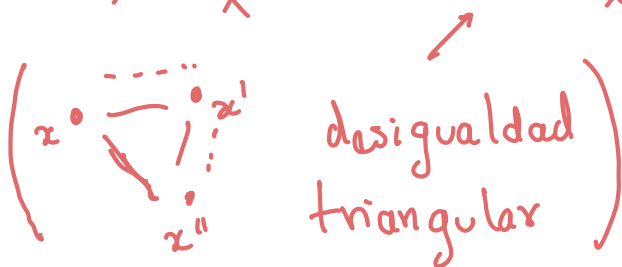
$$d_X : X \times X \rightarrow [0, \infty] \quad \text{tal que}$$

i) $d_X(x, x) = 0$ para todo x en X ($\forall x \in X$)

ii) $d_X(x, x') = d_X(x', x) \quad \forall x, x' \in X$ (simétrica)

Si además

iii) $d_X(x, x'') \leq d_X(x, x') + d_X(x', x'') \quad \forall x, x', x'' \in X$



diremos que d_X es una
pseudo métrica extendida

$(d_X(x, x') = 0, x \neq x')$ $(d_X(x, x') = \infty)$

iv) Si i) - iii) y $d_X(x, x') < \infty \quad \forall x, x' \in X$,

d_X es una pseudo métrica

v) Si i) - iv) y $d_X(x, x') = 0$ solo si $x = x'$,

entonces d_X es una métrica y (X, d_X)

un espacio métrico.

Ejemplos

Distancias entre vectores:

1) $X = \mathbb{R}^n$ y $d_2(x, y) = \|x - y\|_2$

distancia Euclídea
norma Euclídea

donde $\| \underset{x}{(x_1, \dots, x_n)} \|_2 = \left(x_1^2 + \dots + x_n^2 \right)^{1/2}$.

Por tanto $d_2(x, y) = \left((x_1 - y_1)^2 + \dots + (x_n - y_n)^2 \right)^{1/2}$.

Ejercicio:

a) Muestre que $(\mathbb{R}^n, d_2(\cdot, \cdot))$ es un espacio métrico.

b) Sea $p \in \mathbb{N} = \{1, 2, 3, \dots\}$ y para $x, y \in \mathbb{R}^n$ define

$$d_p(x, y) = \|x - y\|_p \leftarrow p\text{-norma}$$

$$= \left(|x_1 - y_1|^p + \dots + |x_n - y_n|^p \right)^{1/p}$$

Muestre que (\mathbb{R}^n, d_p) es un espacio métrico.

c) Para $(x_1, \dots, x_n) \in \mathbb{R}^n$ sea

$$\|(x_1, \dots, x_n)\|_\infty = \max \{ |x_1|, |x_2|, \dots, |x_n| \}$$

↑
norma del supremo

$$\text{y define } d_\infty(x, y) = \|x - y\|_\infty$$

$$= \max \{ |x_1 - y_1|, \dots, |x_n - y_n| \}$$

Muestre que (\mathbb{R}^n, d_∞) es un espacio métrico,

$$\text{y que } \lim_{p \rightarrow \infty} d_p(x, y) = d_\infty(x, y) \quad \forall x, y \in \mathbb{R}^n$$

$d_p(\cdot, \cdot)$:

- $p = 1$: Distancia Manhattan (del taxista)

- $p = 2$: Distancia Euclídea

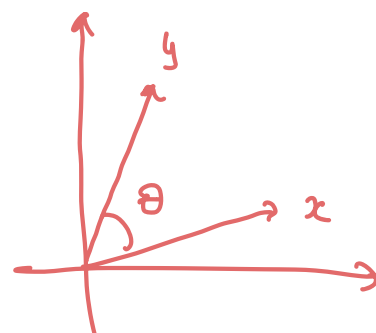
- $p = \infty$: Distancia del supremo.

2) Para $x, y \in \mathbb{R}^n$ sea $x \cdot y = x_1 y_1 + x_2 y_2 + \dots + x_n y_n$
 \uparrow producto punto (interno)

Si $x, y \neq \vec{0}$, define:

$$d_{\cos}(x, y) = \frac{x \cdot y}{\|x\|_2 \cdot \|y\|_2} = \cos \theta$$

\uparrow distancia del coseno.



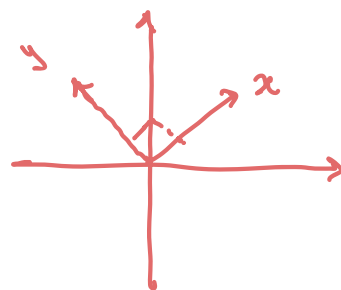
Note: $d_{\cos}(\cdot, \cdot)$ no es una métrica

i) $x = (1, 1)$, $y = (-1, 1)$

$$\|x\|_2 = \|y\|_2 = \sqrt{2}$$

$$x \cdot y = 1 \cdot (-1) + 1 \cdot 1 = 0$$

$$d_{\cos}((1, 1), (-1, 1)) = \frac{0}{\sqrt{2} \cdot \sqrt{2}} = 0$$



ii) Ejercicio: Muestre que $d_{\cos}(\cdot, \cdot)$ no satisface la desigualdad triangular.

Note: $d_{\cos}(\cdot, \cdot)$ es útil para comparar textos.

Diccionario :	palabra 1, palabra 2, ..., palabra n
Texto 1	$(t_{11}, t_{12}, \dots, t_{1n}) = t_1$
\vdots	\vdots
Texto i	$(t_{i1}, t_{i2}, \dots, t_{in}) = t_i$
\vdots	\vdots
Texto m	$(t_{m1}, t_{m2}, \dots, t_{mn}) = t_m$

$t_{ik} = \#$ de veces que la palabra k aparece en el Texto i

$$\text{similitud}(\text{Texto } i, \text{Texto } j) = d_{\cos}(t_i, t_j)$$

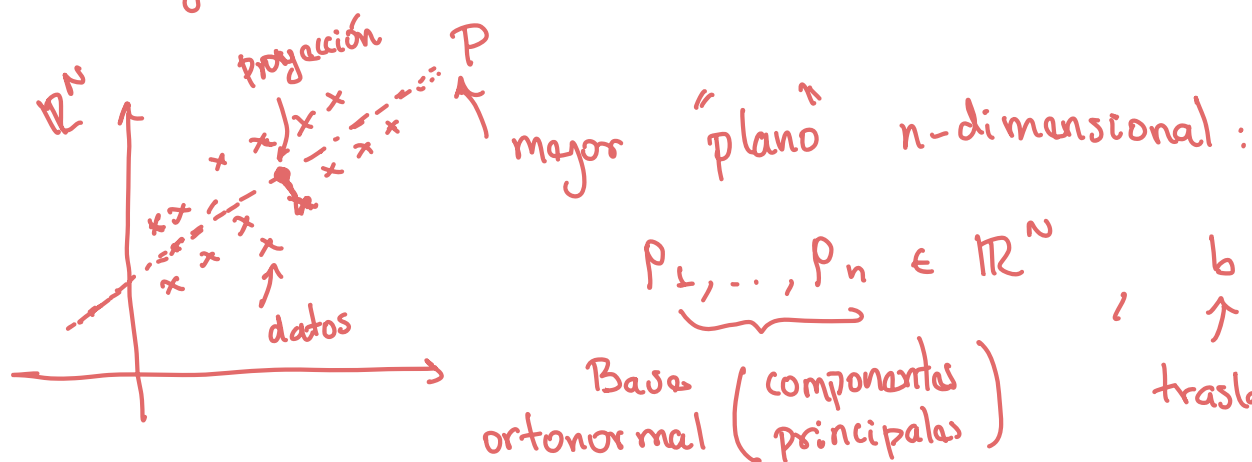
Visualización y Reducción de Dimensión

Problema: Dado (X, d_X) , $X = \{x_1, \dots, x_L\}$,
 encontrar $n \in \mathbb{N}$ (usualmente $1 \leq n \leq 5$)
 y $Y = \{y_1, \dots, y_L\} \subseteq \mathbb{R}^n$ tal que
 $d_X(x_i, x_j) \approx \|y_i - y_j\|_2$, $1 \leq i, j \leq L$

Análisis de Componentes Principales (PCA)

En el caso $X = \{x_1, \dots, x_L\} \subseteq \mathbb{R}^N$ $(N \gg 1)$
 $d_X(x_i, x_j) = \|x_i - x_j\|_2$
 \uparrow
mucho mas grande que

idea: "regresion lineal" (mínimos cuadrados)



$$P = \left\{ b + \lambda_1 P_1 + \dots + \lambda_n P_n \mid \lambda_i \in \mathbb{R}, 1 \leq i \leq n \right\}$$

$\text{proj}_P : \mathbb{R}^N \rightarrow P$ proyección ortogonal

$$y_\ell = \text{proj}_P(x_\ell), \quad 1 \leq \ell \leq L$$

\uparrow \uparrow
proyección Datos

$$(\lambda_{\ell 1}, \lambda_{\ell 2}, \dots, \lambda_{\ell n}) \in \mathbb{R}^n$$

$$(x_{\ell 1}, x_{\ell 2}, \dots, x_{\ell N}) \in \mathbb{R}^N$$