

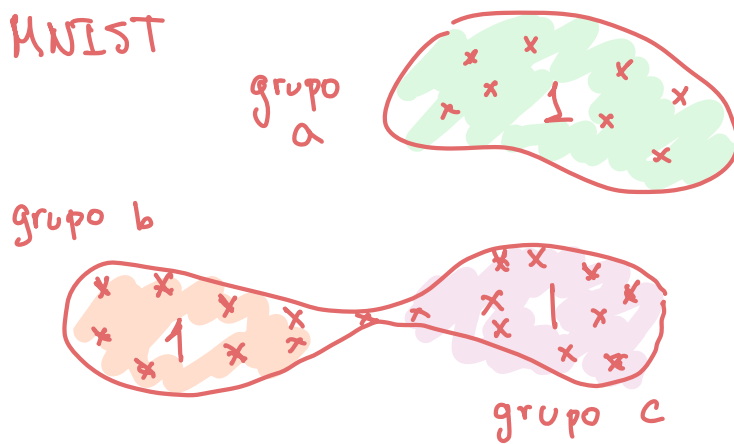
Ayer:

\* Datos  $(X, d_X)$

conjunto  $\downarrow$   $X$

distancia  $\swarrow$   $d_X$

\* MNIST



Topología	Ciencia de datos
Componentes arco-conexas	Clusters (agrupamientos)

Hoy:

- \* Componentes arco-conexas
- \* 0-homología
- \* Clustering (agrupamiento)

— 11 —

Def: Un grafo (no dirigido) es un par

$$G = (V, E) \text{ donde}$$

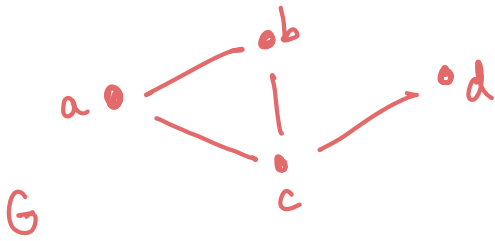
$V \leftarrow$  conjunto de vertices

$$E \subseteq \{ \{v, v'\} \mid v, v' \in V, v \neq v' \}$$

$\uparrow$  conjunto de aristas

## Ejemplos

1)



$$V = \{a, b, c, d\}$$

$$E = \{\{a, b\}, \{a, c\}, \{b, c\}, \{c, d\}\}$$

2) Sea  $(X, d_X)$  un conjunto de datos y  $\varepsilon > 0$ .

$$\text{Sea } E_\varepsilon(X) = \left\{ \{x, x'\} \subseteq X \mid 0 < d_X(x, x') < \varepsilon \right\}$$

y define

$$G_\varepsilon(X, d_X) = (X, E_\varepsilon(X))$$

$\uparrow$   
grafo de  
 $\varepsilon$ -vecindad

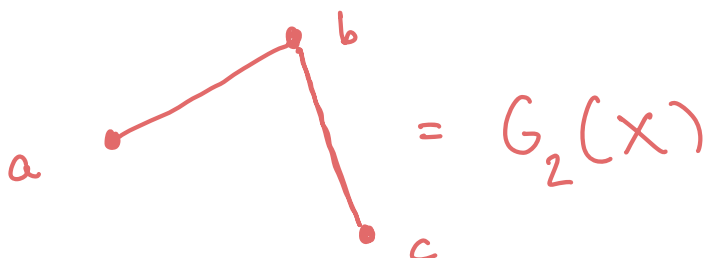
$\uparrow$   
vertices

$\uparrow$   
aristas

e.g.  $X = \{a, b, c\}$  ,

$d_X$	a	b	c
a	0	1	2
b	1	0	1.5
c	2	1.5	0

,  $\varepsilon = 2$



3) Sea  $(X, d_x)$  un espacio métrico finito,  
y  $1 \leq k < \#(X)$ . Para  $x \in X$  sea

$$E_k(x) = \left\{ \{x, x'\} \subseteq X \mid \begin{array}{l} x' \neq x \text{ es uno de los} \\ k \text{ vecinos mas cercanos} \\ \text{a } x \end{array} \right\}$$

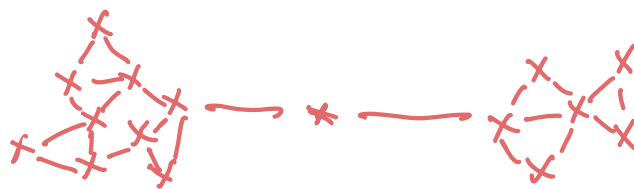
$$G_k(x, d_x) = \left( \underset{\substack{\uparrow \\ \text{grafos de} \\ k\text{-vecinos} \\ \text{mas cercanos}}}{X}, \underset{\substack{\uparrow \\ \text{aristas}}}{E_k(x)} \right)$$

vértices

a.g.

$k=2$

X



Def: Un camino en un grafo  $G = (V, E)$  es  
una sucesión finita

$$\Gamma = \left( \{v_0, v_1\}, \{v_1, v_2\}, \dots, \{v_{n-1}, v_n\} \right)$$

de aristas  $\{v_i, v_{i+1}\} \in E$ ,  $i = 0, \dots, n-1$ .

$v_0 = i(\Gamma) \leftarrow$  inicio,  $v_n = f(\Gamma) \leftarrow$  final.

Def: Dado un grafo  $G = (V, E)$ , diremos que dos vertices  $v, v' \in V$  están conectados por caminos ( $v \sim v'$ ) si  $v = v'$ , o  $v \neq v'$  y existe un camino  $\Gamma$  con  $v = i(\Gamma)$  y  $v' = f(\Gamma)$ .

Prop:  $\sim$  es una relación de equivalencia en  $V$ .

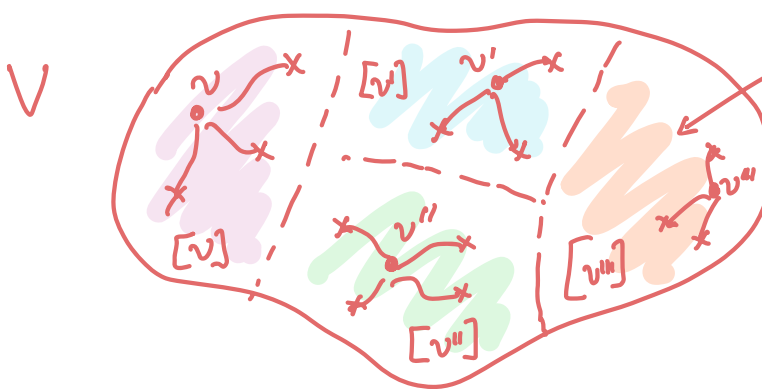
→ reflexiva:  $v \sim v$

→ simétrica: si  $v \sim v'$  entonces  $v' \sim v$

→ transitiva: si  $v \sim v'$  y  $v' \sim v''$  entonces  $v \sim v''$

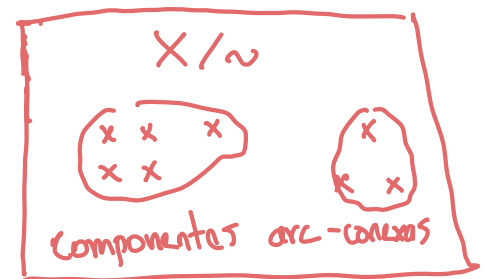
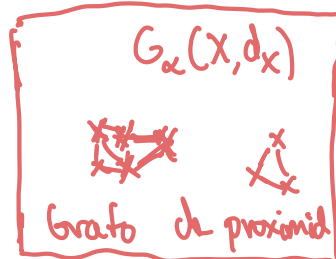
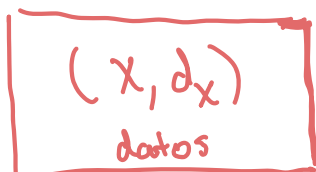
Note:  $\sim$  particiona  $V$  en clases (disjuntas) de equivalencia, llamadas componentes

arco-conexas.



$V/\sim =$  conjunto de componentes arco-conexas

Clustering Topológico:



Algoritmos:  $G = (V, E)$  finito con  $n_E = \#(E)$ .

El algoritmo de Kruskal (para encontrar maximum spanning trees con union-find) toma

$\mathcal{O}(n_E \cdot \log(n_E))$  tiempo.

Versión algebraica: 0-homología

Sea  $(F, +, *)$  un campo (e.g.  $\mathbb{R}, \mathbb{Q}, \mathbb{C}$ ,  
&  $\mathbb{Z}_p = \mathbb{Z}/p\mathbb{Z}$  para  $p \in \mathbb{N}$  primo) y para  $G = (V, E)$

$$\text{sean: } C_0(G; F) = \text{Span}_F(V)$$

$\nearrow$   
0-cadenas  
(vertices)

$$= \left\{ \lambda_1 v_1 + \dots + \lambda_n v_n \mid \begin{array}{l} n \in \mathbb{N}, \lambda_i \in F \\ v_i \in V, 1 \leq i \leq n \end{array} \right\}$$

= Espacio vectorial sobre  $F$  con  
base (i.e., generado por)  $V$ .

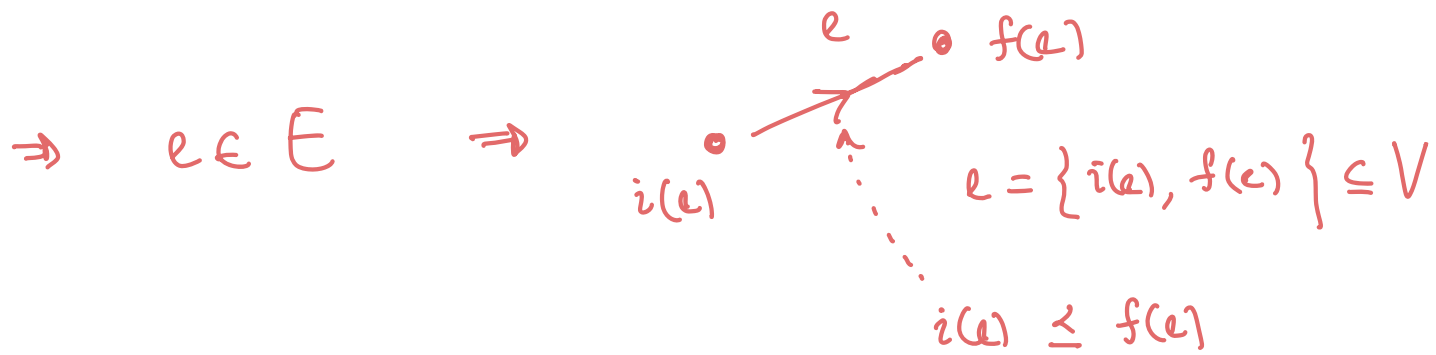
$$C_1(G; F) = \text{Span}_F(E)$$

$\nearrow$   
1-cadenas.  
(aristas)

$$= \left\{ \beta_1 e_1 + \dots + \beta_m e_m \mid \begin{array}{l} m \in \mathbb{N}, \beta_j \in F \\ e_j \in E, 1 \leq j \leq m \end{array} \right\}$$

= Espacio vectorial sobre  $F$  generado  
por  $E$ .

Fije un buen orden  $\leq$  en  $V$  (todo subconjunto no vacío tiene un elemento mínimo = Teorema de Zermelo  $\Leftrightarrow$  Axioma de elección)



Defina la frontera de  $e$  como

$$\partial(e) = f(e) - i(e) \in C_0(G; \mathbb{F})$$

Como  $E$  es una base para  $C_1(G; \mathbb{F})$ , entonces

$\partial$  define una única transformación lineal

$$\partial : C_1(G; \mathbb{F}) \longrightarrow C_0(G; \mathbb{F})$$

Dada por :

$$\begin{aligned} \partial\left(\sum_{j=1}^m \beta_j e_j\right) &= \sum_{j=1}^m \beta_j \cdot \partial(e_j) = \sum_{j=1}^m \beta_j \cdot (f(e_j) - i(e_j)) \\ &= \sum_{j=1}^m \beta_j f(e_j) - \sum_{j=1}^m \beta_j i(e_j) \end{aligned}$$

Ejemplo:   $G$ ,  $\partial(\{a, c\} - \{b, c\}) = b - a$

Teorema: Sea  $G = (V, E)$  un grafo y  $v, v' \in V$ .

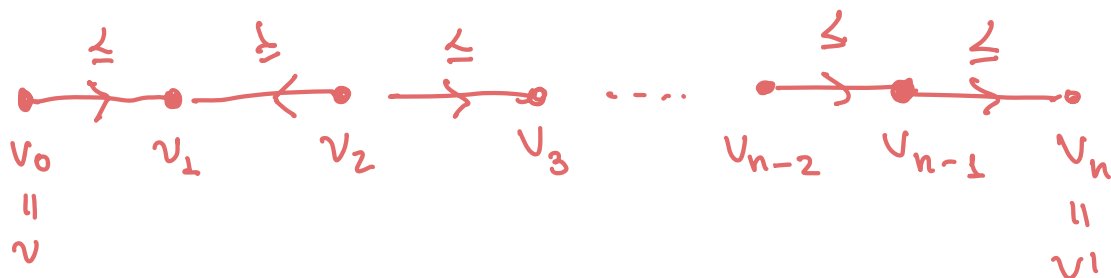
$v \sim v'$  (conectados por caminos) si y solo si

$$v' - v \in \text{img}(\partial).$$

Prueba: ( $\Rightarrow$ ) Si  $v \sim v'$ , entonces existe

un camino  $\Gamma = (\{v_0, v_1\}, \{v_1, v_2\}, \dots, \{v_{n-1}, v_n\})$

en  $G$  con  $v = v_0 = i(\Gamma)$  y  $v' = v_n = f(\Gamma)$ .



$$\text{Sea } s_j = \begin{cases} 1 & \text{si } v_{j-1} \leq v_j \\ -1 & \text{si } v_{j-1} \geq v_j \end{cases}$$

note que  $\partial(s_j \{v_{j-1}, v_j\}) = v_j - v_{j-1}$

$j = 1, \dots, n$

$$\text{Si } \gamma = \sum_{j=1}^n \lambda_j \cdot [v_{j-1}, v_j] \in C_1(G; \mathbb{F}),$$

entonces

$$\begin{aligned} \text{img}(\partial) \ni \partial(\gamma) &= \partial\left(\sum_{j=1}^n \lambda_j [v_{j-1}, v_j]\right) \\ &= \sum_{j=1}^n \partial(\lambda_j [v_{j-1}, v_j]) \\ &= \sum_{j=1}^n v_j - v_{j-1} = v_n - v_0 \\ &= v' - v \end{aligned}$$

$$\Rightarrow v' - v \in \text{img}(\partial).$$

$$(\Leftarrow) \text{ Si } v' - v = \partial\left(\sum_{j=1}^m \beta_j e_j\right), \quad \begin{array}{l} \beta_j \in \mathbb{F} \\ e_j \in E \end{array}$$

entonces uno puede utilizar algunos de los  $e_j$ 's para construir un camino  $\Gamma$  de  $v$  a  $v'$

(ver: Munkres, Elements of Algebraic Topology, Thm 7.1)



Corolario: El espacio vectorial cociente

$$H_0(G; \mathbb{F}) := C_0(G; \mathbb{F}) / \text{img}(\partial)$$

↑  
0-homología  
de  $G$  con coeficientes  
en  $\mathbb{F}$

← base = vertices de  $G$

← difieren por  
un camino

tiene dimension

$$\#(V/\sim) = \# \text{ de componentes arco-conexas de } G.$$