

## Cluster Analysis

Range of techniques aimed at sorting data into clusters, with objects in a cluster more similar to each other than to those in other clusters.

"Cluster" is not a well-defined term. There are many (elementary) approaches to clustering.

We'll focus only on "hierarchical clustering". This is a range of techniques which requires a notion of distance  $d(x, y)$  between objects  $x, y$  to be clustered. The results of these techniques can be represented as dendrograms / phylogenetic trees.

# Example Distances between objects

	h	m	r	c	w
h	0	11	10	14	22
m	11	0	3	13	21
r	10	3	0	12	20
c	14	13	12	0	16
w	22	21	20	16	0

$$V = \{h, m, r, c, w\}$$

Table defines  
a metric on  
 $C$ .

$G(V, t)$  = graph with vertex set  
 $V$ , and one undirected  
edge  $\{x, y\}$  for all  
 $x, y \in V$  with  
 $d(x, y) \leq t$ .

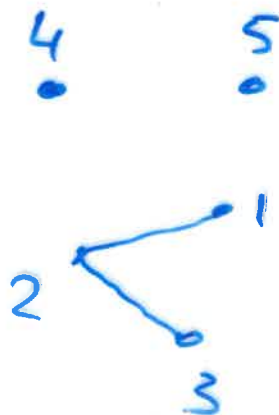
choose  
 $t \geq 0$

for  $t' > t$  we have an inclusion  
of graphs

$$G(V, t) \hookrightarrow G(V, t')$$

for instance,

$G(V, 10)$



[see computer]

$$\pi_0(G(V, 10)) = \{x, y, z\}$$

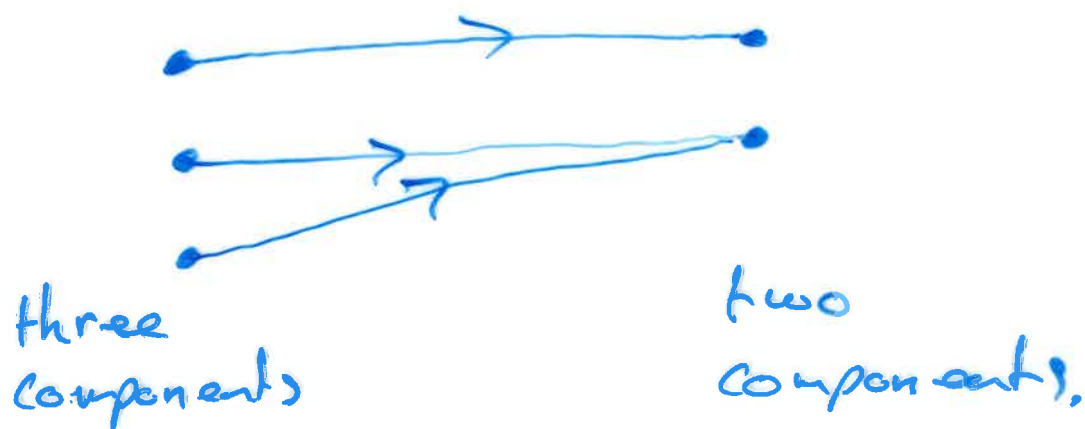
Say,

has three connected components.  
The construction  $\pi_0$  induces a set-theoretic function

$$\pi_0(G(V, t)) \longrightarrow \pi_0(G(V, t'))$$

for any  $t' > t$ . For instance

$$\pi_0(G(V, 10)) \longrightarrow \pi_0(G(V, 14))$$



The functions

$$\pi_0(G(V, t)) \longrightarrow \pi_0(G(V, t'))$$

for  $0 \leq t < t' < \infty$  can be represented as a dendrogram.

See computer

The leaves of the dendrogram represent the objects to be clustered. Two objects  $x$  and  $y$  are considered to be in different clusters if the path from leaf  $x$  to leaf  $y$  is "long".

A barcode is obtained from a dendrogram by removing those edges that represent the merging of clusters.

see computer

see pdf slides

$\pi_0(X)$  = set of connected components of  $X$

is a homotopy invariant:  $X \simeq Y \Rightarrow \pi_0(X) = \pi_0(Y)$ .

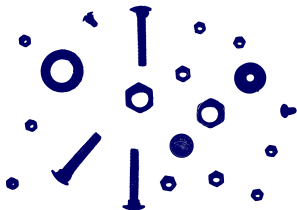
**Toy Application** How does one compute the number of objects in a digital image  $f: \mathbb{R}^2 \rightarrow \mathbb{R}^3$  ?

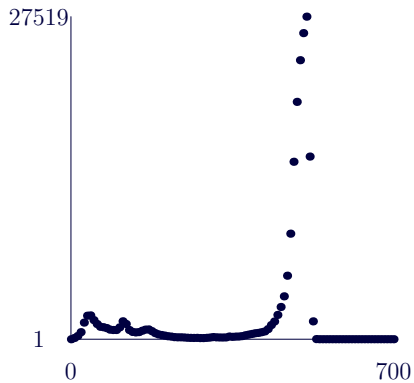


**Toy Application** How does one compute the number of objects in a digital image  $f: \mathbb{R}^2 \rightarrow \mathbb{R}^3$  ?



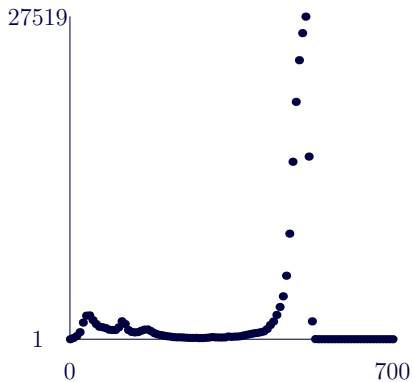
$$S_t = \{(x, y) \in \mathbb{R}^2 : \|f(x, y)\| \leq t\}$$





Plot of  $|\pi_0(S_t)|$  as a function of  $t$



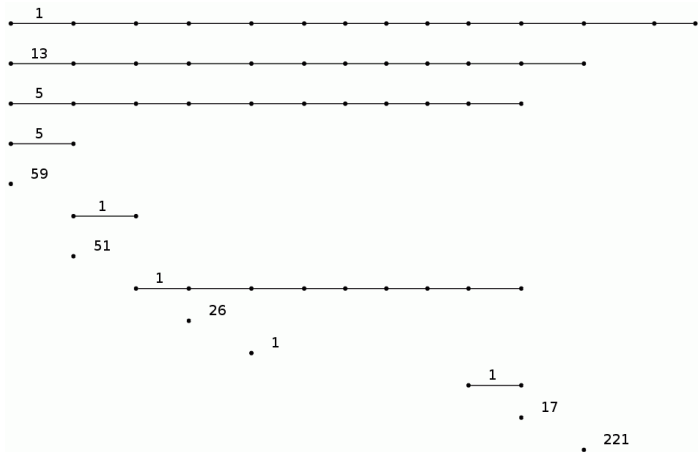


Plot of  $|\pi_0(S_t)|$  as a function of  $t$

$$t_1 < t_2 < \dots < t_T \text{ implies } S_{t_1} \subset S_{t_2} \subset \dots \subset S_{t_T}$$

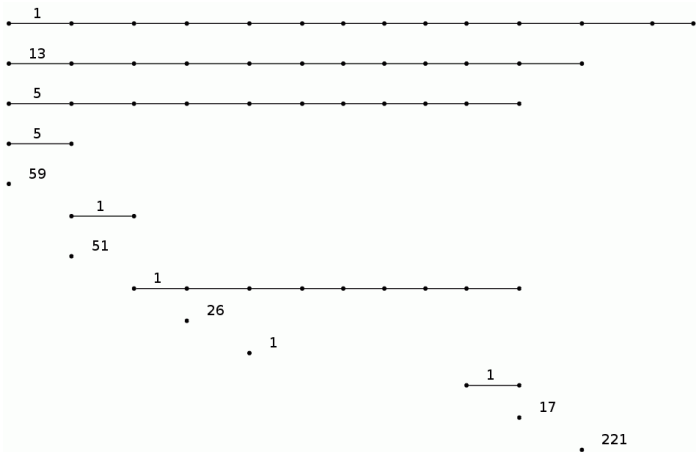
$$\beta_0^{t,t'} = |\text{image}(\pi_0(S_t) \rightarrow \pi_0(S_{t'}))| \text{ for } t \leq t'$$

$$\beta_0^{t,t'} = |\text{image}(\pi_0(S_t) \rightarrow \pi_0(S_{t'}))| \text{ for } t \leq t'$$



$r$  lines from column  $t$  to column  $t'$  if  $\beta_0^{t,t'} = r$

$$\beta_0^{t,t'} = |\text{image}(\pi_0(S_t) \rightarrow \pi_0(S_{t'}))| \text{ for } t \leq t'$$



$r$  lines from column  $t$  to column  $t'$  if  $\beta_0^{t,t'} = r$

There are 20 objects in the photo.

To see how many objects have holes in them, consider

$$S_t^{comp} = \mathbb{R}^2 \setminus S_t,$$

To see how many objects have holes in them, consider

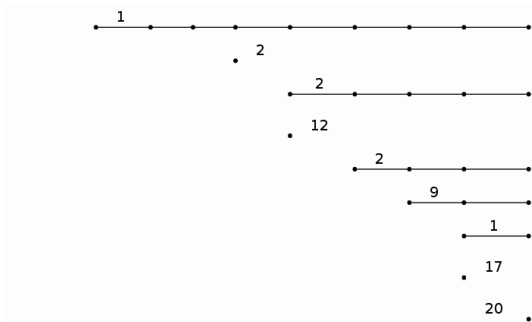
$$S_t^{comp} = \mathbb{R}^2 \setminus S_t,$$

$$\dots \supset S_3^{comp} \supset S_2^{comp} \supset S_1^{comp}.$$

To see how many objects have holes in them, consider

$$S_t^{comp} = \mathbb{R}^2 \setminus S_t,$$

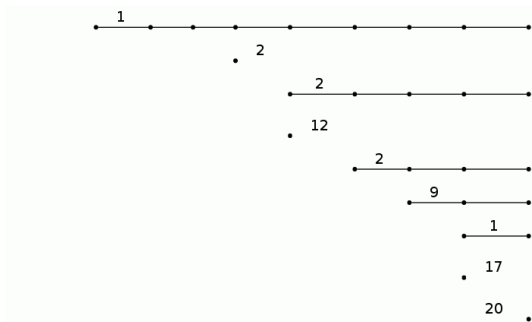
$$\dots \supset S_3^{comp} \supset S_2^{comp} \supset S_1^{comp}.$$



To see how many objects have holes in them, consider

$$S_t^{comp} = \mathbb{R}^2 \setminus S_t,$$

$$\dots \supset S_3^{comp} \supset S_2^{comp} \supset S_1^{comp}.$$



The photo has 14 objects with holes.