

# Classifying Clustering Schemes

Facundo Mémoli,  
joint with Gunnar Carlsson

IMA, 2014

# What is data clustering?

A (standard/flat) **clustering scheme** assigns to any finite set  $X$  a partition  $P_X$  of that set.

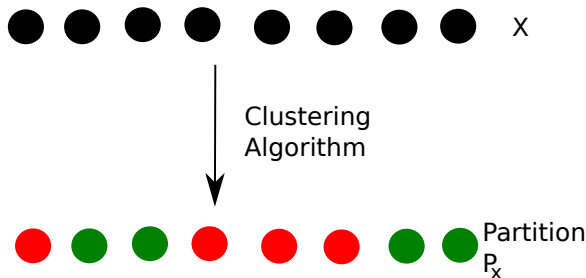


Figure : Elements with the same color are in the same block

# Clustering: setup

- Let  $\mathcal{M} = \{(X, d_X), \text{ finite metric spaces}\}$ .

# Clustering: setup

- Let  $\mathcal{M} = \{(X, d_X), \text{ finite metric spaces}\}$ .
- Given finite set  $X$ ,  $\mathbf{P}(X)$  denotes set of all **partitions** of  $X$ .

# Clustering: setup

- Let  $\mathcal{M} = \{(X, d_X), \text{ finite metric spaces}\}$ .
- Given finite set  $X$ ,  $\mathbf{P}(X)$  denotes set of all **partitions** of  $X$ .
- Let  $\mathcal{P} = \{(X, P_X), X \text{ finite and } P_X \in \mathbf{P}(X)\}$ . Partitioned spaces.

# Clustering: setup

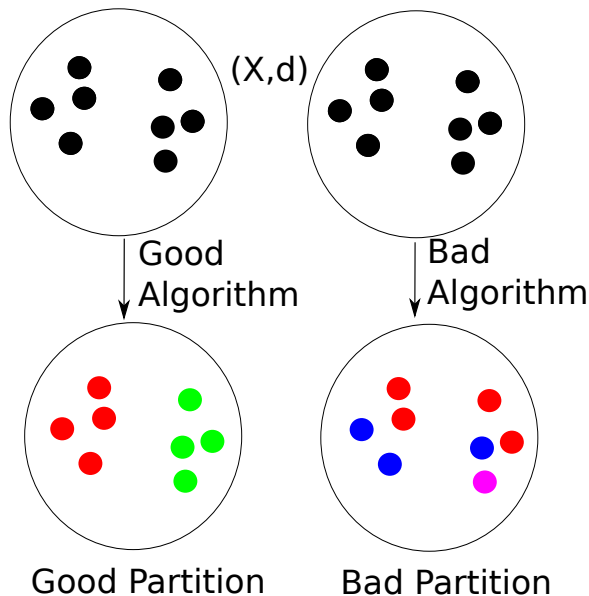
- Let  $\mathcal{M} = \{(X, d_X), \text{ finite metric spaces}\}$ .
- Given finite set  $X$ ,  $\mathbf{P}(X)$  denotes set of all **partitions** of  $X$ .
- Let  $\mathcal{P} = \{(X, P_X), X \text{ finite and } P_X \in \mathbf{P}(X)\}$ . Partitioned spaces.

A **standard clustering scheme** is a map:

$$\mathfrak{C} : \mathcal{M} \longrightarrow \mathcal{P}$$

$$(X, d_X) \mapsto (X, P_X)$$

# Goal of data clustering



# A theorem of J. Kleinberg

## Theorem (2002)

*J. Kleinberg states: there exists no standard clustering algorithm  $\mathcal{C}$  that can simultaneously satisfy the following three properties:*



# A theorem of J. Kleinberg

## Theorem (2002)

*J. Kleinberg states: there exists no standard clustering algorithm  $\mathfrak{C}$  that can simultaneously satisfy the following three properties:*

- *Scale Invariance: For all  $(X, d_X) \in \mathcal{M}$ ,  
 $\mathfrak{C}(X, d_X) = \mathfrak{C}(X, \lambda d_X), \forall \lambda > 0$ .*

# A theorem of J. Kleinberg

## Theorem (2002)

*J. Kleinberg states: there exists no standard clustering algorithm  $\mathfrak{C}$  that can simultaneously satisfy the following three properties:*

- *Scale Invariance: For all  $(X, d_X) \in \mathcal{M}$ ,  
 $\mathfrak{C}(X, d_X) = \mathfrak{C}(X, \lambda d_X)$ ,  $\forall \lambda > 0$ .*
- *Surjectivity:  $\forall X$  and  $P_X \in \mathbf{P}(X) \exists d_X$  s.t.  
 $\mathfrak{C}(X, d_X) = (X, P_X)$ .*

# A theorem of J. Kleinberg

## Theorem (2002)

*J. Kleinberg states: there exists no standard clustering algorithm  $\mathfrak{C}$  that can simultaneously satisfy the following three properties:*

- *Scale Invariance: For all  $(X, d_X) \in \mathcal{M}$ ,  
 $\mathfrak{C}(X, d_X) = \mathfrak{C}(X, \lambda d_X)$ ,  $\forall \lambda > 0$ .*
- *Surjectivity:  $\forall X$  and  $P_X \in \mathbf{P}(X) \exists d_X$  s.t.  
 $\mathfrak{C}(X, d_X) = (X, P_X)$ .*
- *Consistency:*

# A theorem of J. Kleinberg

## Theorem (2002)

*J. Kleinberg states: there exists no standard clustering algorithm  $\mathfrak{C}$  that can simultaneously satisfy the following three properties:*

- *Scale Invariance: For all  $(X, d_X) \in \mathcal{M}$ ,  
 $\mathfrak{C}(X, d_X) = \mathfrak{C}(X, \lambda d_X)$ ,  $\forall \lambda > 0$ .*
- *Surjectivity:  $\forall X$  and  $P_X \in \mathbf{P}(X) \exists d_X$  s.t.  
 $\mathfrak{C}(X, d_X) = (X, P_X)$ .*
- *Consistency: Partition your dataset into blocks.*

# A theorem of J. Kleinberg

## Theorem (2002)

*J. Kleinberg states: there exists no standard clustering algorithm  $\mathfrak{C}$  that can simultaneously satisfy the following three properties:*

- *Scale Invariance: For all  $(X, d_X) \in \mathcal{M}$ ,  
 $\mathfrak{C}(X, d_X) = \mathfrak{C}(X, \lambda d_X)$ ,  $\forall \lambda > 0$ .*
- *Surjectivity:  $\forall X$  and  $P_X \in \mathbf{P}(X) \exists d_X$  s.t.  
 $\mathfrak{C}(X, d_X) = (X, P_X)$ .*
- *Consistency: Partition your dataset into blocks.*
  - *Shrink intrablock distances*

# A theorem of J. Kleinberg

## Theorem (2002)

*J. Kleinberg states: there exists no standard clustering algorithm  $\mathfrak{C}$  that can simultaneously satisfy the following three properties:*

- *Scale Invariance: For all  $(X, d_X) \in \mathcal{M}$ ,  
 $\mathfrak{C}(X, d_X) = \mathfrak{C}(X, \lambda d_X)$ ,  $\forall \lambda > 0$ .*
- *Surjectivity:  $\forall X$  and  $P_X \in \mathbf{P}(X) \exists d_X$  s.t.  
 $\mathfrak{C}(X, d_X) = (X, P_X)$ .*
- *Consistency: Partition your dataset into blocks.*
  - *Shrink intrablock distances*
  - *Expand interblock distances*

# A theorem of J. Kleinberg

## Theorem (2002)

*J. Kleinberg states: there exists no standard clustering algorithm  $\mathfrak{C}$  that can simultaneously satisfy the following three properties:*

- *Scale Invariance: For all  $(X, d_X) \in \mathcal{M}$ ,  
 $\mathfrak{C}(X, d_X) = \mathfrak{C}(X, \lambda d_X)$ ,  $\forall \lambda > 0$ .*
- *Surjectivity:  $\forall X$  and  $P_X \in \mathbf{P}(X) \exists d_X$  s.t.  
 $\mathfrak{C}(X, d_X) = (X, P_X)$ .*
- *Consistency: Partition your dataset into blocks.*
  - *Shrink intrablock distances*
  - *Expand interblock distances*
  - *invariance under this transformation*

# A theorem of J. Kleinberg

## Theorem (2002)

*J. Kleinberg states: there exists no standard clustering algorithm  $\mathfrak{C}$  that can simultaneously satisfy the following three properties:*

- *Scale Invariance: For all  $(X, d_X) \in \mathcal{M}$ ,  
 $\mathfrak{C}(X, d_X) = \mathfrak{C}(X, \lambda d_X)$ ,  $\forall \lambda > 0$ .*
- *Surjectivity:  $\forall X$  and  $P_X \in \mathbf{P}(X) \exists d_X$  s.t.  
 $\mathfrak{C}(X, d_X) = (X, P_X)$ .*
- *Consistency: Partition your dataset into blocks.*
  - *Shrink intrablock distances*
  - *Expand interblock distances*
  - *invariance under this transformation*

Kleinberg's theorem was the inspiration for this work.



# This talk

→ study reformulation/variation of Kleinberg's point of view

# This talk

→ study reformulation/variation of Kleinberg's point of view

Roadmap for studying data clustering:

# This talk

→ study reformulation/variation of Kleinberg's point of view

Roadmap for studying data clustering:

- 1 Recast input/output spaces as categories

# This talk

→ study reformulation/variation of Kleinberg's point of view

Roadmap for studying data clustering:

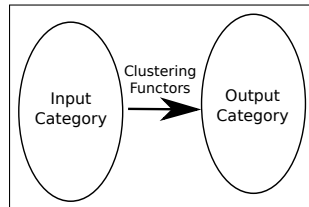
- 1 Recast input/output spaces as categories
- 2 Categories impose a structure on data

# This talk

→ study reformulation/variation of Kleinberg's point of view

Roadmap for studying data clustering:

- 1 Recast input/output spaces as categories
- 2 Categories impose a structure on data
- 3 Study “functors” between categories (respecting this structure)

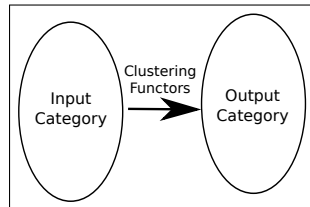


# This talk

→ study reformulation/variation of Kleinberg's point of view

Roadmap for studying data clustering:

- 1 Recast input/output spaces as categories
- 2 Categories impose a structure on data
- 3 Study “functors” between categories (respecting this structure)



We provide a range of input categories (small to large) and show analogues of Kleinberg's theorem in each.

# Roadmap for classification

- Motivation
- **Categories**
- Framework
  - Input Categories
  - Output Category
- Functors
- Results
  - $\mathcal{M}^{iso}$  (small)
  - $\mathcal{M}^{gen}$  (large)
  - $\mathcal{M}^{inj}$  (medium)
- Conclusion

# Category theory

Intuitively: category  $\simeq$  directed multigraph.



# Category theory

Intuitively: category  $\simeq$  directed multigraph.

Definition (A **category**  $\mathcal{C}$  consists of:)

# Category theory

Intuitively: category  $\simeq$  directed multigraph.

**Definition** (A **category**  $\mathcal{C}$  consists of:)

- A collection of **objects**  $\boxed{\text{ob}(\mathcal{C})}$  (*think “vertices”*)

# Category theory

Intuitively: category  $\simeq$  directed multigraph.

**Definition** (A **category**  $\mathcal{C}$  consists of:)

- A collection of **objects**  $\boxed{\text{ob}(\mathcal{C})}$  (*think “vertices”*)
- For each pair of objects  $X, Y$  a set  $\boxed{\text{Mor}_{\mathcal{C}}(X, Y)}$ ,  
the **morphisms** from  $X$  to  $Y$  (*think directed edges/arrows*)

Intuitively: category  $\simeq$  directed multigraph.

**Definition** (A **category**  $\mathcal{C}$  consists of:)

- A collection of **objects**  $\boxed{\text{ob}(\mathcal{C})}$  (*think “vertices”*)
- For each pair of objects  $X, Y$  a set  $\boxed{\text{Mor}_{\mathcal{C}}(X, Y)}$ , the **morphisms** from  $X$  to  $Y$  (*think directed edges/arrows*)
- **Composition operations:**
  - $: \text{Mor}_{\mathcal{C}}(X, Y) \times \text{Mor}_{\mathcal{C}}(Y, Z) \rightarrow \text{Mor}_{\mathcal{C}}(X, Z)$
  - such that ◦ is *associative* (*think function composition*)

Intuitively: category  $\simeq$  directed multigraph.

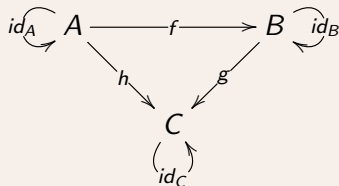
**Definition** (A **category**  $\mathcal{C}$  consists of:)

- A collection of **objects**  $\text{ob}(\mathcal{C})$  (*think “vertices”*)
- For each pair of objects  $X, Y$  a set  $\text{Mor}_{\mathcal{C}}(X, Y)$ , the **morphisms** from  $X$  to  $Y$  (*think directed edges/arrows*)
- **Composition operations:**
  - $\circ : \text{Mor}_{\mathcal{C}}(X, Y) \times \text{Mor}_{\mathcal{C}}(Y, Z) \rightarrow \text{Mor}_{\mathcal{C}}(X, Z)$   
such that  $\circ$  is *associative* (*think function composition*)
- **Identity morphisms:** special morphisms  $X \rightarrow X$ ,  $\text{id}_X \in \text{Mor}_{\mathcal{C}}(X, X)$  for each  $X$  (*think self loops/ identity maps*) such that if  $f \in \text{Mor}_{\mathcal{C}}(X, Y)$  then  $\text{id}_Y \circ f = f \circ \text{id}_X = f$

# Examples of categories

## Example (category $\underline{3}$ )

The category  $\underline{3}$  has exactly three objects  $A$ ,  $B$  and  $C$  and six morphisms: the identities for  $A$ ,  $B, C$ , and three more morphisms,  $\text{Mor}_{\underline{3}}(A, B) = f$ ,  $\text{Mor}_{\underline{3}}(B, C) = g$  and  $\text{Mor}_{\underline{3}}(A, C) = h$ :



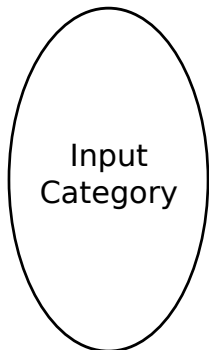
In order to satisfy composition:  $h = g \circ f$ .

# Roadmap for classification

- Motivation
- Categories
- Framework
  - **Input Categories**
  - Output Categories
- Functors
- Results
  - $\mathcal{M}^{iso}$  (small)
  - $\mathcal{M}^{gen}$  (large)
  - $\mathcal{M}^{inj}$  (medium)
- Conclusion

# Roadmap for classification

- Motivation
- Categories
- Framework
  - **Input Categories**
  - Output Categories
- Functors
- Results
  - $\mathcal{M}^{iso}$  (small)
  - $\mathcal{M}^{gen}$  (large)
  - $\mathcal{M}^{inj}$  (medium)
- Conclusion





# Input categories for our model

Consider  $(X, d_X), (Y, d_Y) \in \mathcal{M}$ , where  $\mathcal{M}$  is the collection of all finite metric spaces.

# Input categories for our model

Consider  $(X, d_X), (Y, d_Y) \in \mathcal{M}$ , where  $\mathcal{M}$  is the collection of all finite metric spaces.

Define  $N(X, Y)$  the set of **distance non-increasing maps**:

$$N(X, Y) = \{f : X \rightarrow Y \mid \forall x, x' \in X, d_X(x, x') \geq d_Y(f(x), f(x'))\}$$

# Input categories for our model

Consider  $(X, d_X), (Y, d_Y) \in \mathcal{M}$ , where  $\mathcal{M}$  is the collection of all finite metric spaces.

Define  $N(X, Y)$  the set of **distance non-increasing maps**:

$$N(X, Y) = \{f : X \rightarrow Y \mid \forall x, x' \in X, d_X(x, x') \geq d_Y(f(x), f(x'))\}$$

**Note:** composition works,  $\circ : N(X, Y) \times N(Y, Z) \rightarrow N(X, Z)$ ,  
and  $\text{id}_X \in N(X, X)$ .

# Input categories for our model

Consider  $(X, d_X), (Y, d_Y) \in \mathcal{M}$ , where  $\mathcal{M}$  is the collection of all finite metric spaces.

Define  $N(X, Y)$  the set of **distance non-increasing maps**:

$$N(X, Y) = \{f : X \rightarrow Y \mid \forall x, x' \in X, d_X(x, x') \geq d_Y(f(x), f(x'))\}$$

**Note:** composition works,  $\circ : N(X, Y) \times N(Y, Z) \rightarrow N(X, Z)$ , and  $\text{id}_X \in N(X, X)$ .

We now define three categories:

# Input categories for our model

Consider  $(X, d_X), (Y, d_Y) \in \mathcal{M}$ , where  $\mathcal{M}$  is the collection of all finite metric spaces.

Define  $N(X, Y)$  the set of **distance non-increasing maps**:

$$N(X, Y) = \{f : X \rightarrow Y \mid \forall x, x' \in X, d_X(x, x') \geq d_Y(f(x), f(x'))\}$$

**Note:** composition works,  $\circ : N(X, Y) \times N(Y, Z) \rightarrow N(X, Z)$ , and  $\text{id}_X \in N(X, X)$ .

We now define three categories:

- $\mathcal{M}^{gen}$ :  $\text{ob}(\mathcal{M}^{gen}) = \mathcal{M}$  and  $\text{Mor}_{\mathcal{M}^{gen}}(X, Y) = N(X, Y)$ .

# Input categories for our model

Consider  $(X, d_X), (Y, d_Y) \in \mathcal{M}$ , where  $\mathcal{M}$  is the collection of all finite metric spaces.

Define  $N(X, Y)$  the set of **distance non-increasing maps**:

$$N(X, Y) = \{f : X \rightarrow Y \mid \forall x, x' \in X, d_X(x, x') \geq d_Y(f(x), f(x'))\}$$

**Note:** composition works,  $\circ : N(X, Y) \times N(Y, Z) \rightarrow N(X, Z)$ , and  $\text{id}_X \in N(X, X)$ .

We now define three categories:

- $\mathcal{M}^{gen}$ :  $\text{ob}(\mathcal{M}^{gen}) = \mathcal{M}$  and  $\text{Mor}_{\mathcal{M}^{gen}}(X, Y) = N(X, Y)$ .
- $\mathcal{M}^{inj}$ : the subcategory whose morphisms are also **injective**.

# Input categories for our model

Consider  $(X, d_X), (Y, d_Y) \in \mathcal{M}$ , where  $\mathcal{M}$  is the collection of all finite metric spaces.

Define  $N(X, Y)$  the set of **distance non-increasing maps**:

$$N(X, Y) = \{f : X \rightarrow Y \mid \forall x, x' \in X, d_X(x, x') \geq d_Y(f(x), f(x'))\}$$

**Note:** composition works,  $\circ : N(X, Y) \times N(Y, Z) \rightarrow N(X, Z)$ , and  $\text{id}_X \in N(X, X)$ .

We now define three categories:

- $\mathcal{M}^{gen}$ :  $\text{ob}(\mathcal{M}^{gen}) = \mathcal{M}$  and  $\text{Mor}_{\mathcal{M}^{gen}}(X, Y) = N(X, Y)$ .
- $\mathcal{M}^{inj}$ : the subcategory whose morphisms are also **injective**.
- $\mathcal{M}^{iso}$ : the subcategory whose morphisms are **isometries**:  
 $f : X \rightarrow Y$  is an *isometry* if  $f$  is bijective and  
 $d_Y(f(x), f(x')) = d_X(x, x')$  for all  $x$  and  $x'$ .

# Input categories for our model

Consider  $(X, d_X), (Y, d_Y) \in \mathcal{M}$ , where  $\mathcal{M}$  is the collection of all finite metric spaces.

Define  $N(X, Y)$  the set of **distance non-increasing maps**:

$$N(X, Y) = \{f : X \rightarrow Y \mid \forall x, x' \in X, d_X(x, x') \geq d_Y(f(x), f(x'))\}$$

**Note:** composition works,  $\circ : N(X, Y) \times N(Y, Z) \rightarrow N(X, Z)$ , and  $\text{id}_X \in N(X, X)$ .

We now define three categories:

- $\mathcal{M}^{gen}$ :  $\text{ob}(\mathcal{M}^{gen}) = \mathcal{M}$  and  $\text{Mor}_{\mathcal{M}^{gen}}(X, Y) = N(X, Y)$ .
- $\mathcal{M}^{inj}$ : the subcategory whose morphisms are also **injective**.
- $\mathcal{M}^{iso}$ : the subcategory whose morphisms are **isometries**:  
 $f : X \rightarrow Y$  is an *isometry* if  $f$  is bijective and  
 $d_Y(f(x), f(x')) = d_X(x, x')$  for all  $x$  and  $x'$ .

$$\mathcal{M}^{gen} \supseteq \mathcal{M}^{inj} \supseteq \mathcal{M}^{iso}$$



# Input categories for our model

Consider  $(X, d_X), (Y, d_Y) \in \mathcal{M}$ , where  $\mathcal{M}$  is the collection of all finite metric spaces.

Define  $N(X, Y)$  the set of **distance non-increasing maps**:

$$N(X, Y) = \{f : X \rightarrow Y \mid \forall x, x' \in X, d_X(x, x') \geq d_Y(f(x), f(x'))\}$$

**Note:** composition works,  $\circ : N(X, Y) \times N(Y, Z) \rightarrow N(X, Z)$ , and  $\text{id}_X \in N(X, X)$ .

We now define three categories:

- $\mathcal{M}^{gen}$ :  $\text{ob}(\mathcal{M}^{gen}) = \mathcal{M}$  and  $\text{Mor}_{\mathcal{M}^{gen}}(X, Y) = N(X, Y)$ .
- $\mathcal{M}^{inj}$ : the subcategory whose morphisms are also **injective**.
- $\mathcal{M}^{iso}$ : the subcategory whose morphisms are **isometries**:  
 $f : X \rightarrow Y$  is an *isometry* if  $f$  is bijective and  
 $d_Y(f(x), f(x')) = d_X(x, x')$  for all  $x$  and  $x'$ .

$$\mathcal{M}^{gen} \supseteq \mathcal{M}^{inj} \supseteq \mathcal{M}^{iso}$$

$\mathcal{M}$  will denote any of these categories

Nested categories: same object set, but filter morphisms

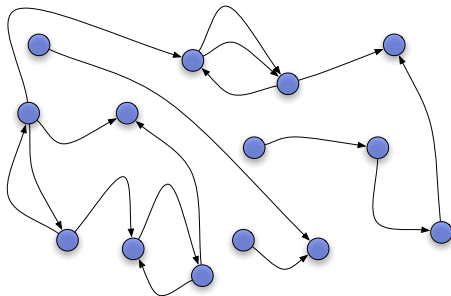


Figure :  $\mathcal{M}^{gen} \supseteq \mathcal{M}^{inj} \supseteq \mathcal{M}^{iso}$ . Same object set, nested morphism sets!

# Nested categories: same object set, but filter morphisms

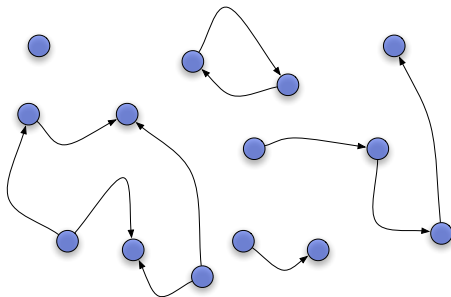


Figure :  $\mathcal{M}^{gen} \supseteq \mathcal{M}^{inj} \supseteq \mathcal{M}^{iso}$ . Same object set, nested morphism sets!

# Nested categories: same object set, but filter morphisms

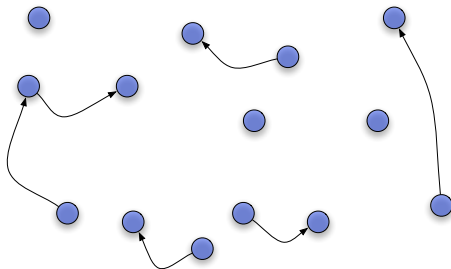
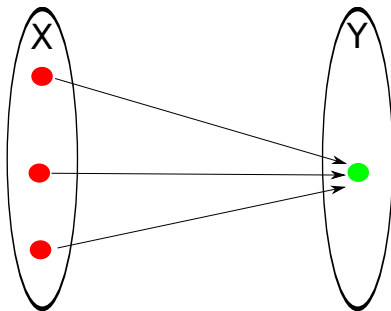


Figure :  $\mathcal{M}^{gen} \supseteq \mathcal{M}^{inj} \supseteq \mathcal{M}^{iso}$ . Same object set, nested morphism sets!

$\text{Mor}_{\mathcal{M}^{gen}}(X, Y)$  is never empty

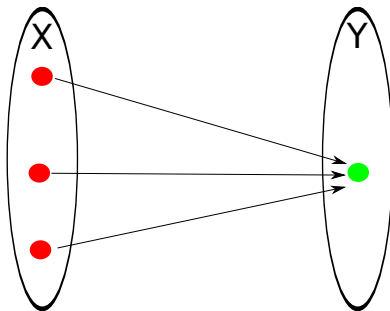
$\mathcal{M}^{gen}$  is special in that  $\forall X, Y \text{ Mor}_{\mathcal{M}^{gen}}(X, Y) \neq \emptyset$



**Figure :** Constant maps are always distance non-increasing. Hence,  $\mathcal{M}^{gen}$  is fully connected! Thus, can always have diagram  $X \rightarrow Y \rightarrow X$  in  $\mathcal{M}^{gen}$ .

# $\text{Mor}_{\mathcal{M}^{gen}}(X, Y)$ is never empty

$\mathcal{M}^{gen}$  is special in that  $\forall X, Y \text{ Mor}_{\mathcal{M}^{gen}}(X, Y) \neq \emptyset$

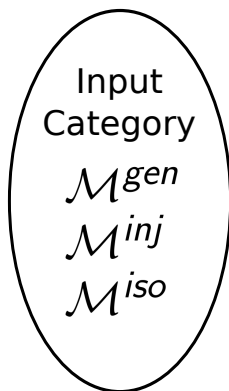


**Figure :** Constant maps are always distance non-increasing. Hence,  $\mathcal{M}^{gen}$  is fully connected! Thus, can always have diagram  $X \rightarrow Y \rightarrow X$  in  $\mathcal{M}^{gen}$ .

This property fails in  $\mathcal{M}^{inj}$  and  $\mathcal{M}^{iso}$ , since morphism  $\phi$  could not send two different elements in  $X$ , to one element in  $Y$ .

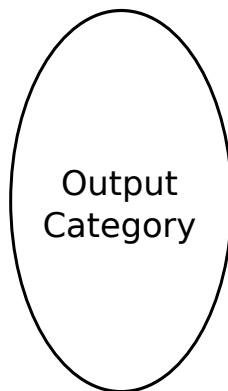
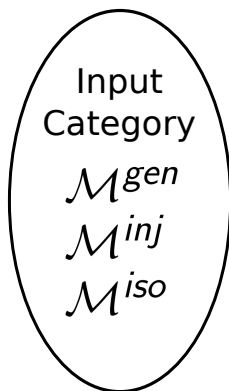
# Roadmap for classification

- Motivation
- Categories
- Framework
  - Input Categories
  - **Output Category**
- Functors
- Results
  - $\mathcal{M}^{iso}$  (small)
  - $\mathcal{M}^{gen}$  (large)
  - $\mathcal{M}^{inj}$  (medium)
- Conclusion



# Roadmap for classification

- Motivation
- Categories
- Framework
  - Input Categories
  - **Output Category**
- Functors
- Results
  - $\mathcal{M}^{iso}$  (small)
  - $\mathcal{M}^{gen}$  (large)
  - $\mathcal{M}^{inj}$  (medium)
- Conclusion





For a finite set  $X$ ,  $\mathbf{P}(X)$  denotes set of all partitions of  $X$ .

For a finite set  $X$ ,  $\mathbf{P}(X)$  denotes set of all partitions of  $X$ .

**Definition ( $\mathcal{P}$ , a category of outputs of standard clustering schemes)**

The category  $\mathcal{P}$  has

For a finite set  $X$ ,  $\mathbf{P}(X)$  denotes set of all partitions of  $X$ .

**Definition ( $\mathcal{P}$ , a category of outputs of standard clustering schemes)**

The category  $\mathcal{P}$  has

- $\text{ob}(\mathcal{P})$  equal to all possible pairs  $(X, P_X)$  where  $X$  is a finite set and  $P_X \in \mathbf{P}(X)$ .

For a finite set  $X$ ,  $\mathbf{P}(X)$  denotes set of all partitions of  $X$ .

## Definition ( $\mathcal{P}$ , a category of outputs of standard clustering schemes)

The category  $\mathcal{P}$  has

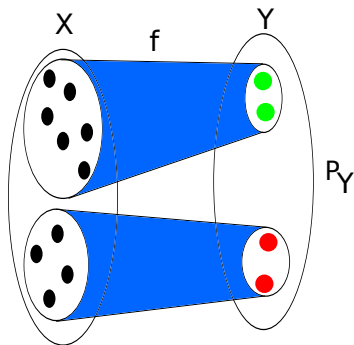
- $\text{ob}(\mathcal{P})$  equal to all possible pairs  $(X, P_X)$  where  $X$  is a finite set and  $P_X \in \mathbf{P}(X)$ .
- $\text{Mor}_{\mathcal{P}}((X, P_X), (Y, P_Y))$  is the set of all maps  $f : X \rightarrow Y$  where  $P_X$  is a refinement of  $f^*(P_Y)$ .

What is  $f^*(P_Y)$ ?

# What is $f^*(P_Y)$ ?

**Definition ( $f^*(P_Y)$ : pullback partition)**

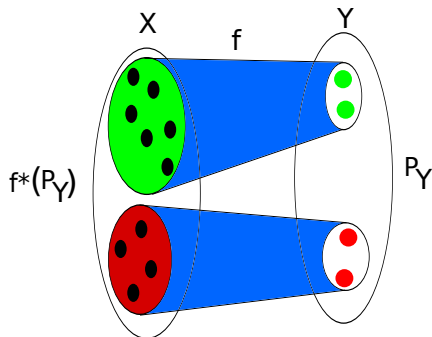
Let  $Y$  be a finite set,  $P_Y \in \mathbf{P}(Y)$ , and  $f : X \rightarrow Y$  be a set map. We define  $f^*(P_Y) = \{f^{-1}(B) : B \in P_Y\} \in \mathbf{P}(X)$ .



# What is $f^*(P_Y)$ ?

**Definition ( $f^*(P_Y)$ : pullback partition)**

Let  $Y$  be a finite set,  $P_Y \in \mathbf{P}(Y)$ , and  $f : X \rightarrow Y$  be a set map. We define  $f^*(P_Y) = \{f^{-1}(B) : B \in P_Y\} \in \mathbf{P}(X)$ .



# Morphisms in $\mathcal{P}$

$\text{Mor}_{\mathcal{P}}((X, P_X), (Y, P_Y))$  is the set of all maps  $f : X \rightarrow Y$  where  $P_X$  is a refinement of  $f^*(P_Y)$

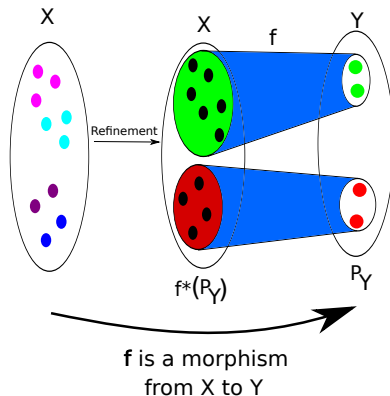


Figure :  $P_X$  refines  $f^*(P_Y)$



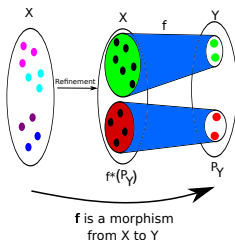
# Repeat: output category

For a finite set  $X$  we denote by  $\mathbf{P}(X)$  the set of all partitions of  $X$ .

## Definition ( $\mathcal{P}$ , a category of outputs of standard clustering schemes)

The category  $\mathcal{P}$  has

- $\text{ob}(\mathcal{P})$  equal to all possible pairs  $(X, P_X)$  where  $X$  is a finite set and  $P_X \in \mathbf{P}(X)$ .
- $\text{Mor}_{\mathcal{P}}((X, P_X), (Y, P_Y))$  is the set of all maps  $f : X \rightarrow Y$  where  $P_X$  is a refinement of  $f^*(P_Y)$ .



## Example (forcing)

- Assume that  $f \in \text{Mor}_{\mathcal{P}}((\{a, b\}, \{\{a, b\}\}), (Y, P_Y))$ . Then,  
 $f(a)$  and  $f(b)$  must be in same block of  $P_Y$ .

## Example (forcing)

- Assume that  $f \in \text{Mor}_{\mathcal{P}}((\{a, b\}, \{\{a, b\}\}), (Y, P_Y))$ . Then,

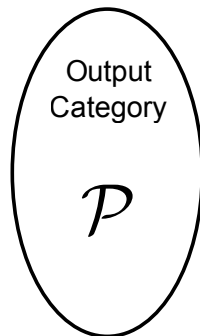
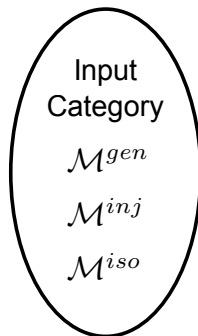
$f(a)$  and  $f(b)$  must be in same block of  $P_Y$ .

- Let  $f \in \text{Mor}_{\mathcal{P}}((X, P_X), (\{a, b\}, \{\{a\}, \{b\}\}))$ . Then,

$f(x) \neq f(x')$ , implies that  $x, x'$  must be in different blocks of  $P_X$ .

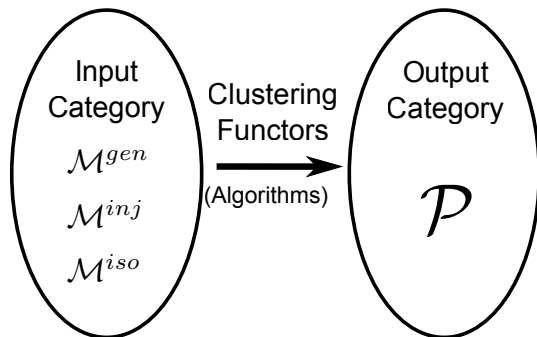
# Roadmap for classification

- Motivation
- Categories
- Framework
  - Input Categories
  - Output Categories
- **Functors**
- Results
  - $\mathcal{M}^{iso}$  (small)
  - $\mathcal{M}^{gen}$  (large)
  - $\mathcal{M}^{inj}$  (medium)
- Conclusion



# Roadmap for classification

- Motivation
- Categories
- Framework
  - Input Categories
  - Output Categories
- **Functors**
- Results
  - $\mathcal{M}^{iso}$  (small)
  - $\mathcal{M}^{gen}$  (large)
  - $\mathcal{M}^{inj}$  (medium)
- Conclusion



## Definition (**Functor**)

Let  $\mathcal{C}$  and  $\mathcal{D}$  be categories. A **functor** from  $\mathcal{C}$  to  $\mathcal{D}$  consists of:

## Definition (**Functor**)

Let  $\mathcal{C}$  and  $\mathcal{D}$  be categories. A **functor** from  $\mathcal{C}$  to  $\mathcal{D}$  consists of:

- **A mapping of the objects:**  $\Phi : \text{ob}(\mathcal{C}) \rightarrow \text{ob}(\mathcal{D})$ .

## Definition (**Functor**)

Let  $\mathcal{C}$  and  $\mathcal{D}$  be categories. A **functor** from  $\mathcal{C}$  to  $\mathcal{D}$  consists of:

- **A mapping of the objects:**  $\Phi : \text{ob}(\mathcal{C}) \rightarrow \text{ob}(\mathcal{D})$ .
- **A mapping of the morphisms:**  $\forall X, Y \in \text{ob}(\mathcal{C})$   
 $\Phi : \text{Mor}_{\mathcal{C}}(X, Y) \rightarrow \text{Mor}_{\mathcal{D}}(\Phi(X), \Phi(Y))$  so that



## Definition (**Functor**)

Let  $\mathcal{C}$  and  $\mathcal{D}$  be categories. A **functor** from  $\mathcal{C}$  to  $\mathcal{D}$  consists of:

- **A mapping of the objects:**  $\Phi : \text{ob}(\mathcal{C}) \rightarrow \text{ob}(\mathcal{D})$ .
- **A mapping of the morphisms:**  $\forall X, Y \in \text{ob}(\mathcal{C})$   
 $\Phi : \text{Mor}_{\mathcal{C}}(X, Y) \rightarrow \text{Mor}_{\mathcal{D}}(\Phi(X), \Phi(Y))$  so that
  - ① **identities are respected:**  $\Phi(\text{id}_X) = \text{id}_{\Phi(X)}$

## Definition (**Functor**)

Let  $\mathcal{C}$  and  $\mathcal{D}$  be categories. A **functor** from  $\mathcal{C}$  to  $\mathcal{D}$  consists of:

- **A mapping of the objects:**  $\Phi : \text{ob}(\mathcal{C}) \rightarrow \text{ob}(\mathcal{D})$ .
- **A mapping of the morphisms:**  $\forall X, Y \in \text{ob}(\mathcal{C})$   
 $\Phi : \text{Mor}_{\mathcal{C}}(X, Y) \rightarrow \text{Mor}_{\mathcal{D}}(\Phi(X), \Phi(Y))$  so that
  - 1 **identities are respected:**  $\Phi(\text{id}_X) = \text{id}_{\Phi(X)}$
  - 2 **composition is respected:**  $\Phi(g \circ f) = \Phi(g) \circ \Phi(f)$

# Clustering methods as functors

- For  $\mathcal{M}$  being each of our three input categories, we are going to require that for all morphisms (test functions)  $f \in \text{Mor}_{\mathcal{M}}(X, Y)$ :

$$\begin{array}{ccc} (X, d_X) & \xrightarrow{f} & (Y, d_Y) \\ \downarrow \mathfrak{c} & & \downarrow \mathfrak{c} \\ (X, P_X) & \xrightarrow{f} & (Y, P_Y) \end{array}$$

# Clustering methods as functors

- For  $\mathcal{M}$  being each of our three input categories, we are going to require that for all morphisms (test functions)  $f \in \text{Mor}_{\mathcal{M}}(X, Y)$ :

$$\begin{array}{ccc} (X, d_X) & \xrightarrow{f} & (Y, d_Y) \\ \downarrow \mathfrak{C} & & \downarrow \mathfrak{C} \\ (X, P_X) & \xrightarrow{f} & (Y, P_Y) \end{array}$$

- Take for example  $(X, d)$  and  $(X, d')$ , two metrics on  $X$  s.t.  $d \geq d'$ . If  $P$  and  $P'$  are the respective partitions produced by a clustering functor  $\mathfrak{C}$ , then if  $x \sim_P x'$ , it also has to happen that  $x \sim_{P'} x'$  as well. Namely, reducing distances has to make it easier to cluster two points!

# Clustering methods as functors

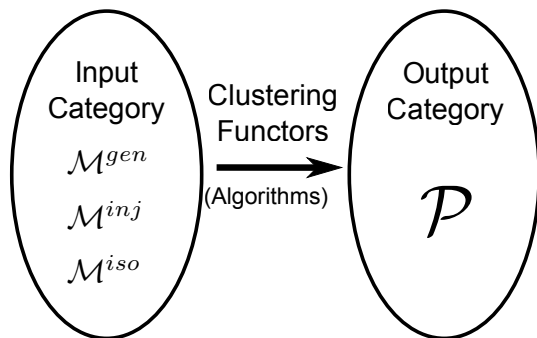
- For  $\mathcal{M}$  being each of our three input categories, we are going to require that for all morphisms (test functions)  $f \in \text{Mor}_{\mathcal{M}}(X, Y)$ :

$$\begin{array}{ccc} (X, d_X) & \xrightarrow{f} & (Y, d_Y) \\ \downarrow \mathfrak{C} & & \downarrow \mathfrak{C} \\ (X, P_X) & \xrightarrow{f} & (Y, P_Y) \end{array}$$

- Take for example  $(X, d)$  and  $(X, d')$ , two metrics on  $X$  s.t.  $d \geq d'$ . If  $P$  and  $P'$  are the respective partitions produced by a clustering functor  $\mathfrak{C}$ , then if  $x \sim_P x'$ , it also has to happen that  $x \sim_{P'} x'$  as well. **Namely, reducing distances has to make it easier to cluster two points!**
- Notice that it is more difficult to find functors  $\mathcal{M}^{gen} \rightarrow \mathcal{P}$  than to find functors  $\mathcal{M}^{inj} \rightarrow \mathcal{P}$ , than to find functors  $\mathcal{M}^{iso} \rightarrow \mathcal{P}$ .

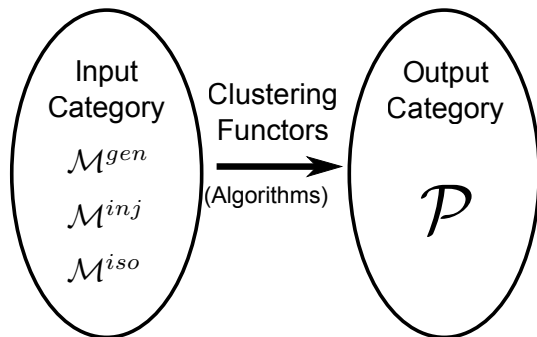
# Roadmap for classification

- Motivation
- Categories
- Framework
  - Input Categories
  - Output Categories
- Functors
- Results
  - $\mathcal{M}^{iso}$  (**small**)
  - $\mathcal{M}^{gen}$  (large)
  - $\mathcal{M}^{inj}$  (medium)
- Conclusion

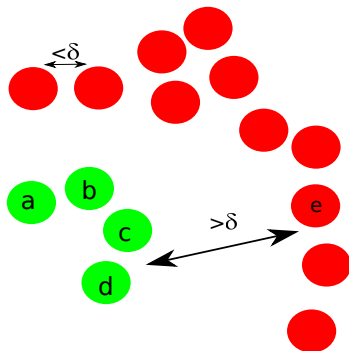


# Roadmap for classification

- Motivation
- Categories
- Framework
  - Input Categories
  - Output Categories
- Functors
- Results
  - $\mathcal{M}^{iso}$  (small)
  - $\mathcal{M}^{gen}$  (**large**)
  - $\mathcal{M}^{inj}$  (medium)
- Conclusion



# Single linkage clustering



## Definition

On  $(X, d_X) \in \mathcal{M}$ , for each  $\delta \geq 0$ , define the **equivalence relation**  $\sim_\delta$ , where  $x \sim_\delta x' \iff$  there is a sequence  $x_0, x_1, \dots, x_k \in X$  so that  $x_0 = x, x_k = x'$ , and  $d_X(x_i, x_{i+1}) \leq \delta$  for all  $i$ . Let  $\mathbf{P}_X(\delta)$  be the resulting partition.



# $\mathcal{M}^{gen}$ functorial clustering algorithms

## Definition (**Vietoris-Rips clustering functor** on $\mathcal{M}^{gen}$ )

For each  $\delta \geq 0$  define

$$\mathfrak{R}_\delta : \mathcal{M}^{gen} \rightarrow \mathcal{P}$$

by

$$\mathfrak{R}_\delta(X, d_X) = (X, P_X(\delta)).$$

A.K.A. the **single linkage clustering** on parameter  $\delta$ .

# $\mathcal{M}^{gen}$ functorial clustering algorithms

## Definition (**Vietoris-Rips clustering functor on $\mathcal{M}^{gen}$** )

For each  $\delta \geq 0$  define

$$\mathfrak{R}_\delta : \mathcal{M}^{gen} \rightarrow \mathcal{P}$$

by

$$\mathfrak{R}_\delta(X, d_X) = (X, P_X(\delta)).$$

A.K.A. the **single linkage clustering** on parameter  $\delta$ .

One checks this *is* a functor (**existence**), and prove a **uniqueness** theorem for it.

# $\mathcal{M}^{gen}$ functorial clustering algorithms

## Definition (**Vietoris-Rips clustering functor** on $\mathcal{M}^{gen}$ )

For each  $\delta \geq 0$  define

$$\mathfrak{R}_\delta : \mathcal{M}^{gen} \rightarrow \mathcal{P}$$

by

$$\mathfrak{R}_\delta(X, d_X) = (X, P_X(\delta)).$$

A.K.A. the **single linkage clustering** on parameter  $\delta$ .

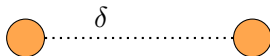
One checks this *is* a functor (**existence**), and prove a **uniqueness** theorem for it.

Functoriality is equivalent to:

$\phi \in \text{Mor}_{\mathcal{M}^{gen}}(X, Y)$ , then  $x \sim_\delta x'$  in  $X$ , implies  $\phi(x) \sim_\delta \phi(x')$  in  $Y$ .

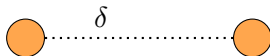
# Uniqueness

Let  $\Delta_2(\delta)$  be the metric space of two points, where  $d(x, y) = \delta$ :



# Uniqueness

Let  $\Delta_2(\delta)$  be the metric space of two points, where  $d(x, y) = \delta$ :

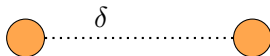


## Theorem

*Assume that  $\mathfrak{C} : \mathcal{M}^{gen} \rightarrow \mathcal{P}$  is a clustering functor for which there exists  $\delta_{\mathfrak{C}} > 0$  with the property that*

# Uniqueness

Let  $\Delta_2(\delta)$  be the metric space of two points, where  $d(x, y) = \delta$ :



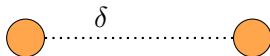
## Theorem

Assume that  $\mathfrak{C} : \mathcal{M}^{gen} \rightarrow \mathcal{P}$  is a clustering functor for which there exists  $\delta_{\mathfrak{C}} > 0$  with the property that

- 1  $\mathfrak{C}(\Delta_2(\delta))$  is in one piece for all  $\delta \in [0, \delta_{\mathfrak{C}}]$ , and

# Uniqueness

Let  $\Delta_2(\delta)$  be the metric space of two points, where  $d(x, y) = \delta$ :



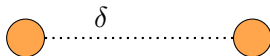
## Theorem

Assume that  $\mathfrak{C} : \mathcal{M}^{gen} \rightarrow \mathcal{P}$  is a clustering functor for which there exists  $\delta_{\mathfrak{C}} > 0$  with the property that

- 1  $\mathfrak{C}(\Delta_2(\delta))$  is in one piece for all  $\delta \in [0, \delta_{\mathfrak{C}}]$ , and
- 2  $\mathfrak{C}(\Delta_2(\delta))$  is in two pieces for all  $\delta > \delta_{\mathfrak{C}}$ .

# Uniqueness

Let  $\Delta_2(\delta)$  be the metric space of two points, where  $d(x, y) = \delta$ :



## Theorem

Assume that  $\mathfrak{C} : \mathcal{M}^{\text{gen}} \rightarrow \mathcal{P}$  is a clustering functor for which there exists  $\delta_{\mathfrak{C}} > 0$  with the property that

- 1  $\mathfrak{C}(\Delta_2(\delta))$  is in one piece for all  $\delta \in [0, \delta_{\mathfrak{C}}]$ , and
- 2  $\mathfrak{C}(\Delta_2(\delta))$  is in two pieces for all  $\delta > \delta_{\mathfrak{C}}$ .

Then,  $\mathfrak{C}$  is the Vietoris-Rips functor with parameter  $\delta_{\mathfrak{C}}$ . i.e.  
 $\mathfrak{C} = \mathfrak{R}_{\delta_{\mathfrak{C}}}$ .



- Great so SLC is unique in  $\mathcal{M}^{gen}$ !

- Great so SLC is unique in  $\mathcal{M}^{gen}$ !
- Recall Kleinberg's criteria:

- Great so SLC is unique in  $\mathcal{M}^{gen}$ !
- Recall Kleinberg's criteria:
  - *Scale Invariance*:  $\mathfrak{C}(X, d) = \mathfrak{C}(X, \lambda d) \forall \lambda > 0$
  - *Surjectivity*:  $\forall P_X \exists d_X$  such that:  $\mathfrak{C}(X, d_X) = P_X$
  - *Consistency*: Partition your dataset into blocks.
    - Shrink intrablock distances
    - Expand interblock distances
    - invariance under this transformation

- Great so SLC is unique in  $\mathcal{M}^{gen}$ !
- Recall Kleinberg's criteria:
  - *Scale Invariance*:  $\mathfrak{C}(X, d) = \mathfrak{C}(X, \lambda d) \forall \lambda > 0$
  - *Surjectivity*:  $\forall P_X \exists d_X$  such that:  $\mathfrak{C}(X, d_X) = P_X$
  - *Consistency*: Partition your dataset into blocks.
    - Shrink intrablock distances
    - Expand interblock distances
    - invariance under this transformation
- To what degree do these properties hold for the VR/SL functor?

## More on the VR functor

- The VR functor is surjective, this can be checked easily.

## More on the VR functor

- The VR functor is surjective, this can be checked easily.
- Functoriality is our substitute for consistency.

## More on the VR functor

- The VR functor is surjective, this can be checked easily.
- Functoriality is our substitute for consistency.
- The Vietoris-Rips functor is **not** scale invariant over  $\mathcal{P}$ !

## More on the VR functor

- The VR functor is surjective, this can be checked easily.
- Functoriality is our substitute for consistency.
- The Vietoris-Rips functor is **not** scale invariant over  $\mathcal{P}$ !

Theorem (Scale invariance in  $\mathcal{M}^{gen}$ )



## More on the VR functor

- The VR functor is surjective, this can be checked easily.
- Functoriality is our substitute for consistency.
- The Vietoris-Rips functor is **not** scale invariant over  $\mathcal{P}$ !

### Theorem (Scale invariance in $\mathcal{M}^{gen}$ )

*Let  $\sigma_\lambda$  be the mapping  $(X, d_X) \mapsto (X, \lambda \cdot d_X)$*

## More on the VR functor

- The VR functor is surjective, this can be checked easily.
- Functoriality is our substitute for consistency.
- The Vietoris-Rips functor is **not** scale invariant over  $\mathcal{P}$ !

### Theorem (Scale invariance in $\mathcal{M}^{gen}$ )

Let  $\sigma_\lambda$  be the mapping  $(X, d_X) \mapsto (X, \lambda \cdot d_X)$

Let  $\mathfrak{C} : \mathcal{M}^{gen} \rightarrow \mathcal{P}$  be a clustering functor s.t.  $\mathfrak{C} \circ \sigma_\lambda = \mathfrak{C}$  for all  $\lambda > 0$ . Then, either

## More on the VR functor

- The VR functor is surjective, this can be checked easily.
- Functoriality is our substitute for consistency.
- The Vietoris-Rips functor is **not** scale invariant over  $\mathcal{P}$ !

### Theorem (Scale invariance in $\mathcal{M}^{gen}$ )

Let  $\sigma_\lambda$  be the mapping  $(X, d_X) \mapsto (X, \lambda \cdot d_X)$

Let  $\mathfrak{C} : \mathcal{M}^{gen} \rightarrow \mathcal{P}$  be a clustering functor s.t.  $\mathfrak{C} \circ \sigma_\lambda = \mathfrak{C}$  for all  $\lambda > 0$ . Then, either

- $\mathfrak{C}$  assigns to each finite metric space  $X$  the partition of  $X$  into singletons, or

## More on the VR functor

- The VR functor is surjective, this can be checked easily.
- Functoriality is our substitute for consistency.
- The Vietoris-Rips functor is **not** scale invariant over  $\mathcal{P}$ !

### Theorem (Scale invariance in $\mathcal{M}^{gen}$ )

Let  $\sigma_\lambda$  be the mapping  $(X, d_X) \mapsto (X, \lambda \cdot d_X)$

Let  $\mathfrak{C} : \mathcal{M}^{gen} \rightarrow \mathcal{P}$  be a clustering functor s.t.  $\mathfrak{C} \circ \sigma_\lambda = \mathfrak{C}$  for all  $\lambda > 0$ . Then, either

- $\mathfrak{C}$  assigns to each finite metric space  $X$  the partition of  $X$  into singletons, or
- $\mathfrak{C}$  assigns to each finite metric the partition with only one block.

## More on the VR functor

- The VR functor is surjective, this can be checked easily.
- Functoriality is our substitute for consistency.
- The Vietoris-Rips functor is **not** scale invariant over  $\mathcal{P}$ !

### Theorem (Scale invariance in $\mathcal{M}^{gen}$ )

Let  $\sigma_\lambda$  be the mapping  $(X, d_X) \mapsto (X, \lambda \cdot d_X)$

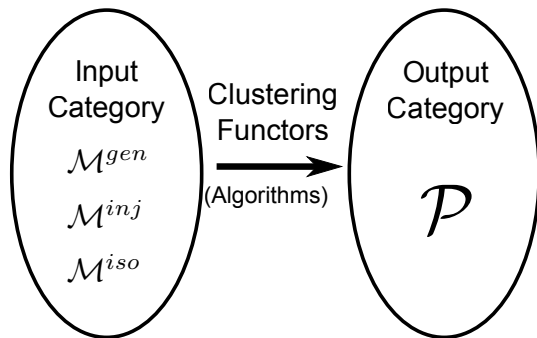
Let  $\mathfrak{C} : \mathcal{M}^{gen} \rightarrow \mathcal{P}$  be a clustering functor s.t.  $\mathfrak{C} \circ \sigma_\lambda = \mathfrak{C}$  for all  $\lambda > 0$ . Then, either

- $\mathfrak{C}$  assigns to each finite metric space  $X$  the partition of  $X$  into singletons, or
- $\mathfrak{C}$  assigns to each finite metric the partition with only one block.

A similar theorem holds in  $\mathcal{M}^{inj}$ .

# Roadmap for classification

- Motivation
- Categories
- Framework
  - Input Categories
  - Output Categories
- Functors
- Results
  - $\mathcal{M}^{iso}$  (small)
  - $\mathcal{M}^{gen}$  (large)
  - $\mathcal{M}^{inj}$  (**medium**)
- Conclusion



# A property for functorial algorithms

## Definition (**Excisive clustering functors**)

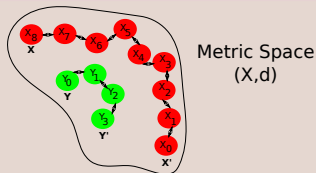
We say that a clustering functor  $\mathfrak{C}$  is **excisive** if for all  $(X, d_X) \in \text{ob}(\mathcal{M})$ , if we write  $\mathfrak{C}(X, d_X) = (X, \{X_\alpha\}_{\alpha \in A})$ , then  $\mathfrak{C}(X_\alpha, d_{X|_{X_\alpha \times X_\alpha}}) = (X_\alpha, \{X_\alpha\})$  for all  $\alpha \in A$ .

# A property for functorial algorithms

## Definition (**Excisive clustering functors**)

We say that a clustering functor  $\mathfrak{C}$  is **excisive** if for all  $(X, d_X) \in \text{ob}(\mathcal{M})$ , if we write  $\mathfrak{C}(X, d_X) = (X, \{X_\alpha\}_{\alpha \in A})$ , then  $\mathfrak{C}(X_\alpha, d_{X|_{X_\alpha \times X_\alpha}}) = (X_\alpha, \{X_\alpha\})$  for all  $\alpha \in A$ .

## The Vietoris-Rips functor is excisive.



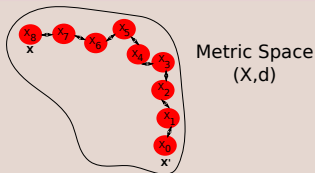


# A property for functorial algorithms

## Definition (**Excisive clustering functors**)

We say that a clustering functor  $\mathfrak{C}$  is **excisive** if for all  $(X, d_X) \in \text{ob}(\mathcal{M})$ , if we write  $\mathfrak{C}(X, d_X) = (X, \{X_\alpha\}_{\alpha \in A})$ , then  $\mathfrak{C}(X_\alpha, d_{X|_{X_\alpha \times X_\alpha}}) = (X_\alpha, \{X_\alpha\})$  for all  $\alpha \in A$ .

## The Vietoris-Rips functor is excisive.

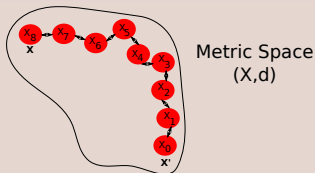


# A property for functorial algorithms

## Definition (**Excisive clustering functors**)

We say that a clustering functor  $\mathfrak{C}$  is **excisive** if for all  $(X, d_X) \in \text{ob}(\mathcal{M})$ , if we write  $\mathfrak{C}(X, d_X) = (X, \{X_\alpha\}_{\alpha \in A})$ , then  $\mathfrak{C}(X_\alpha, d_{X|_{X_\alpha \times X_\alpha}}) = (X_\alpha, \{X_\alpha\})$  for all  $\alpha \in A$ .

## The Vietoris-Rips functor is excisive.



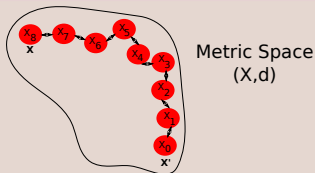
Do non-excisive methods exist?

# A property for functorial algorithms

## Definition (**Excisive clustering functors**)

We say that a clustering functor  $\mathfrak{C}$  is **excisive** if for all  $(X, d_X) \in \text{ob}(\mathcal{M})$ , if we write  $\mathfrak{C}(X, d_X) = (X, \{X_\alpha\}_{\alpha \in A})$ , then  $\mathfrak{C}(X_\alpha, d_{X|_{X_\alpha \times X_\alpha}}) = (X_\alpha, \{X_\alpha\})$  for all  $\alpha \in A$ .

## The Vietoris-Rips functor is excisive.



**Do non-excisive methods exist?** Definitely not in  $\mathcal{M}^{gen}$ !

## Excisiveness in $\mathcal{M}^{inj}$

In  $\mathcal{M}^{inj}$  there exist **non-excisive** clustering functors: for each finite metric space  $X$  (with at least two points):

## Excisiveness in $\mathcal{M}^{inj}$

In  $\mathcal{M}^{inj}$  there exist **non-excisive** clustering functors: for each finite metric space  $X$  (with at least two points): let

$$\eta(X) := (\text{sep}(X))^{-1},$$

where  $\text{sep}(X) = \min_{x \neq x'} d_X(x, x')$  (separation of a metric space).

# Excisiveness in $\mathcal{M}^{inj}$

In  $\mathcal{M}^{inj}$  there exist **non-excisive** clustering functors: for each finite metric space  $X$  (with at least two points): let

$$\eta(X) := (\text{sep}(X))^{-1},$$

where  $\text{sep}(X) = \min_{x \neq x'} d_X(x, x')$  (separation of a metric space).  
Let  $\hat{\mathfrak{K}}$  be functor assigning to each  $X$  partition into equivalence classes of  $\sim_{\eta(X)}$  (now parameter  $\delta$  depends on the space!).

# Excisiveness in $\mathcal{M}^{inj}$

In  $\mathcal{M}^{inj}$  there exist **non-excisive** clustering functors: for each finite metric space  $X$  (with at least two points): let

$$\eta(X) := (\text{sep}(X))^{-1},$$

where  $\text{sep}(X) = \min_{x \neq x'} d_X(x, x')$  (separation of a metric space).

Let  $\hat{\mathfrak{K}}$  be functor assigning to each  $X$  partition into equivalence classes of  $\sim_{\eta(X)}$  (now parameter  $\delta$  depends on the space!).

That this indeed defines a functor in  $\mathcal{M}^{inj}$  follows from:

$$\text{Mor}_{\mathcal{M}^{inj}}(X, Y) \neq \emptyset \implies \text{sep}(X) \geq \text{sep}(Y).$$

# Excisiveness in $\mathcal{M}^{inj}$

In  $\mathcal{M}^{inj}$  there exist **non-excisive** clustering functors: for each finite metric space  $X$  (with at least two points): let

$$\eta(X) := (\text{sep}(X))^{-1},$$

where  $\text{sep}(X) = \min_{x \neq x'} d_X(x, x')$  (separation of a metric space).

Let  $\hat{\mathfrak{R}}$  be functor assigning to each  $X$  partition into equivalence classes of  $\sim_{\eta(X)}$  (now parameter  $\delta$  depends on the space!).

That this indeed defines a functor in  $\mathcal{M}^{inj}$  follows from:

$$\text{Mor}_{\mathcal{M}^{inj}}(X, Y) \neq \emptyset \implies \text{sep}(X) \geq \text{sep}(Y).$$

Indeed, assume that  $\phi \in \text{Mor}_{\mathcal{M}^{inj}}(X, Y)$ , then for all  $x \neq x'$ :

$$d_X(x, x') \geq d_Y(\phi(x), \phi(x')) \boxed{\geq} \min_{y \neq y'} d_Y(y, y') = \text{sep}(Y).$$



# Excisiveness in $\mathcal{M}^{inj}$

In  $\mathcal{M}^{inj}$  there exist **non-excisive** clustering functors: for each finite metric space  $X$  (with at least two points): let

$$\eta(X) := (\text{sep}(X))^{-1},$$

where  $\text{sep}(X) = \min_{x \neq x'} d_X(x, x')$  (separation of a metric space).

Let  $\hat{\mathfrak{R}}$  be functor assigning to each  $X$  partition into equivalence classes of  $\sim_{\eta(X)}$  (now parameter  $\delta$  depends on the space!).

That this indeed defines a functor in  $\mathcal{M}^{inj}$  follows from:

$$\text{Mor}_{\mathcal{M}^{inj}}(X, Y) \neq \emptyset \implies \text{sep}(X) \geq \text{sep}(Y).$$

Indeed, assume that  $\phi \in \text{Mor}_{\mathcal{M}^{inj}}(X, Y)$ , then for all  $x \neq x'$ :

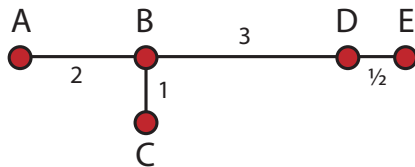
$$d_X(x, x') \geq d_Y(\phi(x), \phi(x')) \boxed{\geq} \min_{y \neq y'} d_Y(y, y') = \text{sep}(Y).$$

Then,  $\text{sep}(X) \geq \text{sep}(Y)$ , and since  $\mathfrak{R}_\delta$  is functorial and  $\eta(Y) \geq \eta(X)$ ,

$$x \sim_{\eta(X)} x' \implies \phi(x) \sim_{\eta(X)} \phi(x') \implies \phi(x) \sim_{\eta(Y)} \phi(x').$$

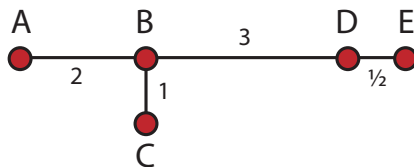
## Excisiveness in $\mathcal{M}^{inj}$

Consider  $X$  to be the metric space below (graph distance):



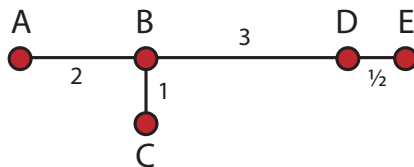
## Excisiveness in $\mathcal{M}^{inj}$

Consider  $X$  to be the metric space below (graph distance):



Note that  $\eta(X) = 2$

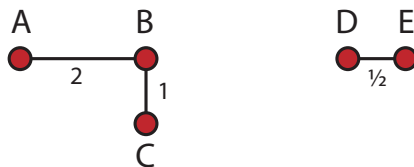
Consider  $X$  to be the metric space below (graph distance):



Note that  $\eta(X) = 2 \implies$

$$\hat{\mathfrak{R}}(X, d_X) = \mathfrak{R}_2(X, d_X) = (X, \{\{A, B, C\}, \{D, E\}\}).$$

Consider  $X$  to be the metric space below (graph distance):

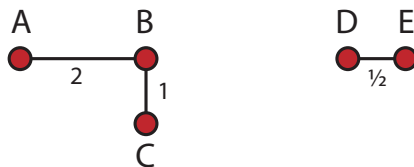


Note that  $\eta(X) = 2 \implies$

$$\hat{\mathfrak{R}}(X, d_X) = \mathfrak{R}_2(X, d_X) = (X, \{\{A, B, C\}, \{D, E\}\}).$$

Now, if  $Y = \{A, B, C\}$ , then  $\eta(Y) = 1 \implies \{A, B, C\}$  gets further partitioned! Then,  $\hat{\mathfrak{R}}$  is **not-excisive**.

Consider  $X$  to be the metric space below (graph distance):

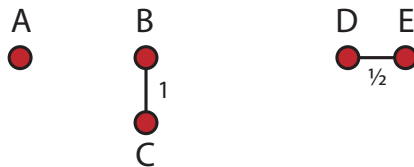


Note that  $\eta(X) = 2 \implies$

$$\hat{\mathfrak{R}}(X, d_X) = \mathfrak{R}_2(X, d_X) = (X, \{\{A, B, C\}, \{D, E\}\}).$$

Now, if  $Y = \{A, B, C\}$ , then  $\eta(Y) = 1 \implies \{A, B, C\}$  gets further partitioned! Then,  $\hat{\mathfrak{R}}$  is **not-excisive**.

Consider  $X$  to be the metric space below (graph distance):



Note that  $\eta(X) = 2 \implies$

$$\hat{\mathfrak{R}}(X, d_X) = \mathfrak{R}_2(X, d_X) = (X, \{\{A, B, C\}, \{D, E\}\}).$$

Now, if  $Y = \{A, B, C\}$ , then  $\eta(Y) = 1 \implies \{A, B, C\}$  gets further partitioned! Then,  $\hat{\mathfrak{R}}$  is **not-excisive**.

# a factory of monsters





# a factory of monsters



Let  $\iota : \mathcal{M}^{inj} \rightarrow (\mathbb{R}, \geq)$  by any non-constant contra-variant functor:

$$\text{Mor}_{\mathcal{M}^{inj}}(X, Y) \neq \emptyset \implies \iota(X) \leq \iota(Y).$$

For each of them know how to build a non-excusive monster:  $\mathfrak{R}_{\iota(\cdot)}$ .

# a factory of monsters



Let  $\iota : \mathcal{M}^{inj} \rightarrow (\mathbb{R}, \geq)$  by any non-constant contra-variant functor:

$$\text{Mor}_{\mathcal{M}^{inj}}(X, Y) \neq \emptyset \implies \iota(X) \leq \iota(Y).$$

For each of them know how to build a non-excusive monster:  $\mathfrak{R}_{\iota(\cdot)}$ .

We saw that no such  $\iota$  exist on  $\mathcal{M}^{gen}$ : uniqueness of VR.

# a factory of monsters



Let  $\iota : \mathcal{M}^{inj} \rightarrow (\mathbb{R}, \geq)$  by any non-constant contra-variant functor:

$$\text{Mor}_{\mathcal{M}^{inj}}(X, Y) \neq \emptyset \implies \iota(X) \leq \iota(Y).$$

For each of them know how to build a non-excisive monster:  $\mathfrak{R}_{\iota(\cdot)}$ .

We saw that no such  $\iota$  exist on  $\mathcal{M}^{gen}$ : uniqueness of VR.

Indirectly:  $X \rightarrow Y \rightarrow X$  is always a possible diagram on  $\mathcal{M}^{gen}$ .

Apply  $\iota$  and get  $\iota(X) = \iota(Y)$ .

And now something completely different..



# Generalize single linkage clustering

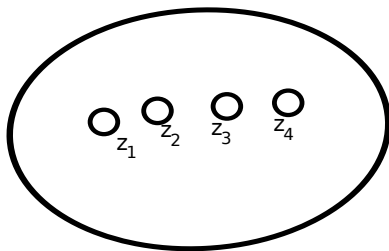
Notice the following: given  $\delta \geq 0$ ,  $(X, d_X) \in \mathcal{M}$  and  $x, x' \in X$ ,

$$d_X(x, x') \leq \delta \iff \exists \phi \in \text{Mor}_{\mathcal{M}}(\Delta_2(\delta), X) \text{ s.t. } \{x, x'\} \in \text{im}(\phi).$$

# Generalize single linkage clustering

Notice the following: given  $\delta \geq 0$ ,  $(X, d_X) \in \mathcal{M}$  and  $x, x' \in X$ ,

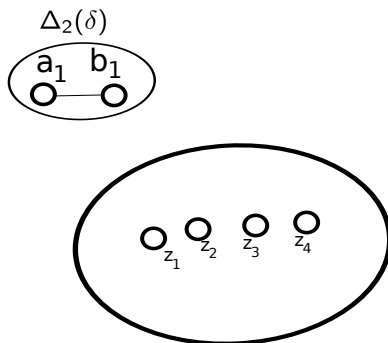
$$d_X(x, x') \leq \delta \iff \exists \phi \in \text{Mor}_{\mathcal{M}}(\Delta_2(\delta), X) \text{ s.t. } \{x, x'\} \in \text{im}(\phi).$$



# Generalize single linkage clustering

Notice the following: given  $\delta \geq 0$ ,  $(X, d_X) \in \mathcal{M}$  and  $x, x' \in X$ ,

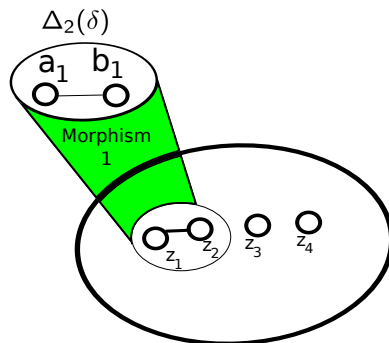
$$d_X(x, x') \leq \delta \iff \exists \phi \in \text{Mor}_{\mathcal{M}}(\Delta_2(\delta), X) \text{ s.t. } \{x, x'\} \in \text{im}(\phi).$$



# Generalize single linkage clustering

Notice the following: given  $\delta \geq 0$ ,  $(X, d_X) \in \mathcal{M}$  and  $x, x' \in X$ ,

$$d_X(x, x') \leq \delta \iff \exists \phi \in \text{Mor}_{\mathcal{M}}(\Delta_2(\delta), X) \text{ s.t. } \{x, x'\} \in \text{im}(\phi).$$

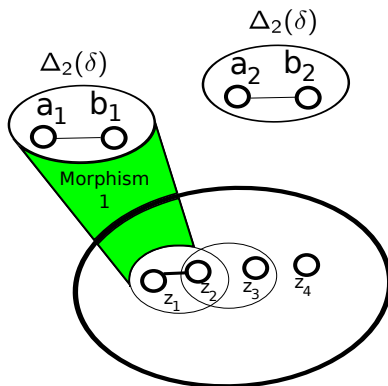




# Generalize single linkage clustering

Notice the following: given  $\delta \geq 0$ ,  $(X, d_X) \in \mathcal{M}$  and  $x, x' \in X$ ,

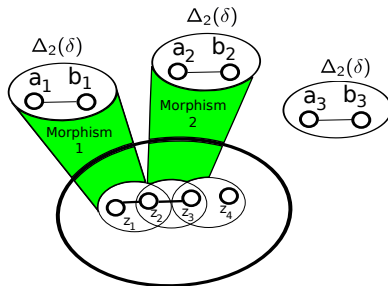
$$d_X(x, x') \leq \delta \iff \exists \phi \in \text{Mor}_{\mathcal{M}}(\Delta_2(\delta), X) \text{ s.t. } \{x, x'\} \in \text{im}(\phi).$$



# Generalize single linkage clustering

Notice the following: given  $\delta \geq 0$ ,  $(X, d_X) \in \mathcal{M}$  and  $x, x' \in X$ ,

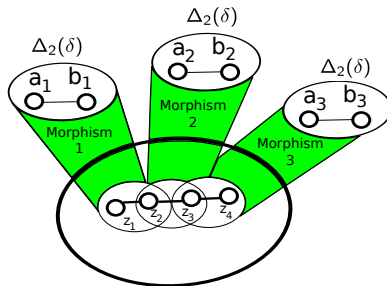
$$d_X(x, x') \leq \delta \iff \exists \phi \in \text{Mor}_{\mathcal{M}}(\Delta_2(\delta), X) \text{ s.t. } \{x, x'\} \in \text{im}(\phi).$$



# Generalize single linkage clustering

Notice the following: given  $\delta \geq 0$ ,  $(X, d_X) \in \mathcal{M}$  and  $x, x' \in X$ ,

$$d_X(x, x') \leq \delta \iff \exists \phi \in \text{Mor}_{\mathcal{M}}(\Delta_2(\delta), X) \text{ s.t. } \{x, x'\} \in \text{im}(\phi).$$



What if we use something else instead of  $\Delta_2(\delta)$ ?

# Generalize single linkage clustering

## Definition (representable functors)

We say a functor  $\mathfrak{C}$  is **representable** if there is a collection  $\Omega$  of finite metric spaces such that: points  $x$  and  $x'$  in  $X$  are in the same block of  $\mathfrak{C}(X, d_X)$  if and only if there exist:

# Generalize single linkage clustering

## Definition (representable functors)

We say a functor  $\mathfrak{C}$  is **representable** if there is a collection  $\Omega$  of finite metric spaces such that: points  $x$  and  $x'$  in  $X$  are in the same block of  $\mathfrak{C}(X, d_X)$  if and only if there exist:

- a sequence of points  $z_0, \dots, z_k \in X$  with  $z_0 = x$  and  $z_k = x'$ ,

## Definition (representable functors)

We say a functor  $\mathfrak{C}$  is **representable** if there is a collection  $\Omega$  of finite metric spaces such that: points  $x$  and  $x'$  in  $X$  are in the same block of  $\mathfrak{C}(X, d_X)$  if and only if there exist:

- a sequence of points  $z_0, \dots, z_k \in X$  with  $z_0 = x$  and  $z_k = x'$ ,
- a sequence of metric spaces  $\omega_1, \dots, \omega_k \in \Omega$ ,

# Generalize single linkage clustering

## Definition (representable functors)

We say a functor  $\mathfrak{C}$  is **representable** if there is a collection  $\Omega$  of finite metric spaces such that: points  $x$  and  $x'$  in  $X$  are in the same block of  $\mathfrak{C}(X, d_X)$  if and only if there exist:

- a sequence of points  $z_0, \dots, z_k \in X$  with  $z_0 = x$  and  $z_k = x'$ ,
- a sequence of metric spaces  $\omega_1, \dots, \omega_k \in \Omega$ ,
- for each  $i = 1, \dots, k$ , pairs of points  $(\alpha_i, \beta_i) \in \omega_i$  and morphisms  $f_i \in \text{Mor}_{\mathcal{M}}(\omega_i, X)$  s.t.  $f_i(\alpha_i) = z_{i-1}$  and  $f_i(\beta_i) = z_i$ .

# Generalize single linkage clustering

## Definition (representable functors)

We say a functor  $\mathfrak{C}$  is **representable** if there is a collection  $\Omega$  of finite metric spaces such that: points  $x$  and  $x'$  in  $X$  are in the same block of  $\mathfrak{C}(X, d_X)$  if and only if there exist:

- a sequence of points  $z_0, \dots, z_k \in X$  with  $z_0 = x$  and  $z_k = x'$ ,
- a sequence of metric spaces  $\omega_1, \dots, \omega_k \in \Omega$ ,
- for each  $i = 1, \dots, k$ , pairs of points  $(\alpha_i, \beta_i) \in \omega_i$  and morphisms  $f_i \in \text{Mor}_{\mathcal{M}}(\omega_i, X)$  s.t.  $f_i(\alpha_i) = z_{i-1}$  and  $f_i(\beta_i) = z_i$ .

If  $\Omega$  is finite, then we say that  $\mathfrak{C}$  is **finitely representable**. We write  $\mathfrak{C}^\Omega$  for a functor that is represented by  $\Omega$ .



# Generalize single linkage clustering

## Definition (representable functors)

We say a functor  $\mathfrak{C}$  is **representable** if there is a collection  $\Omega$  of finite metric spaces such that: points  $x$  and  $x'$  in  $X$  are in the same block of  $\mathfrak{C}(X, d_X)$  if and only if there exist:

- a sequence of points  $z_0, \dots, z_k \in X$  with  $z_0 = x$  and  $z_k = x'$ ,
- a sequence of metric spaces  $\omega_1, \dots, \omega_k \in \Omega$ ,
- for each  $i = 1, \dots, k$ , pairs of points  $(\alpha_i, \beta_i) \in \omega_i$  and morphisms  $f_i \in \text{Mor}_{\mathcal{M}}(\omega_i, X)$  s.t.  $f_i(\alpha_i) = z_{i-1}$  and  $f_i(\beta_i) = z_i$ .

If  $\Omega$  is finite, then we say that  $\mathfrak{C}$  is **finitely representable**. We write  $\mathfrak{C}^\Omega$  for a functor that is represented by  $\Omega$ .

Notice if  $\Omega = \{\Delta_2(\delta)\}$  then we see that  $\mathfrak{C}^\Omega$  is the VR functor, thus it is representable.

# Generalize single linkage clustering

## Definition (representable functors)

We say a functor  $\mathfrak{C}$  is **representable** if there is a collection  $\Omega$  of finite metric spaces such that: points  $x$  and  $x'$  in  $X$  are in the same block of  $\mathfrak{C}(X, d_X)$  if and only if there exist:

- a sequence of points  $z_0, \dots, z_k \in X$  with  $z_0 = x$  and  $z_k = x'$ ,
- a sequence of metric spaces  $\omega_1, \dots, \omega_k \in \Omega$ ,
- for each  $i = 1, \dots, k$ , pairs of points  $(\alpha_i, \beta_i) \in \omega_i$  and morphisms  $f_i \in \text{Mor}_{\mathcal{M}}(\omega_i, X)$  s.t.  $f_i(\alpha_i) = z_{i-1}$  and  $f_i(\beta_i) = z_i$ .

If  $\Omega$  is finite, then we say that  $\mathfrak{C}$  is **finitely representable**. We write  $\mathfrak{C}^\Omega$  for a functor that is represented by  $\Omega$ .

Notice if  $\Omega = \{\Delta_2(\delta)\}$  then we see that  $\mathfrak{C}^\Omega$  is the VR functor, thus it is representable. Representability means: **parametric description of clustering methods**.

# Clustering in $\mathcal{M}^{inj}$ (and $\mathcal{M}^{gen}$ )

## Theorem (Excisive = Representable)

*A clustering functor is excisive if and only if it is representable.*

# Clustering in $\mathcal{M}^{inj}$ (and $\mathcal{M}^{gen}$ )

## Theorem (Excisive = Representable)

*A clustering functor is excisive if and only if it is representable.*

## Theorem (Factorization)

*If  $\mathcal{M}$  is either  $\mathcal{M}^{inj}$  or  $\mathcal{M}^{gen}$  then if  $\mathfrak{C}$  is any  $\mathcal{M}$ -functorial finitely represented functor, represented by  $\Omega$ ,*

*Then  $\mathfrak{C} = \mathfrak{R}_1 \circ \boxed{\mathfrak{T}^\Omega}$*

# Clustering in $\mathcal{M}^{inj}$ (and $\mathcal{M}^{gen}$ )

## Theorem (Excisive = Representable)

*A clustering functor is excisive if and only if it is representable.*

## Theorem (Factorization)

*If  $\mathcal{M}$  is either  $\mathcal{M}^{inj}$  or  $\mathcal{M}^{gen}$  then if  $\mathfrak{C}$  is any  $\mathcal{M}$ -functorial finitely represented functor, represented by  $\Omega$ ,*

*Then  $\mathfrak{C} = \mathfrak{R}_1 \circ \boxed{\mathfrak{T}^\Omega}$   $\Leftarrow$  a certain transformation that changes the metric:*

$$\mathfrak{T}^\Omega : \mathcal{M} \rightarrow \mathcal{M}$$

.

# Excisive clustering functors in $\mathcal{M}^{inj}$

Our factorization theorem suggests how to change metrics to account for **density**– ameliorate SL's **chaining effect**.

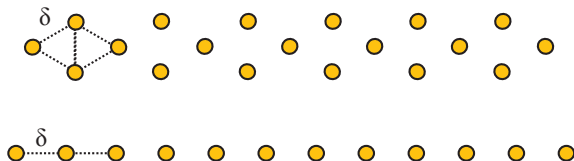


Figure :  $\mathfrak{C}^\Omega$ , with  $\Omega = \{\Delta_n(\delta)\}$  for increasing  $n$ , requires more and more density.

# Excisive clustering functors in $\mathcal{M}^{inj}$

Our factorization theorem suggests how to change metrics to account for **density**– ameliorate SL's **chaining effect**.

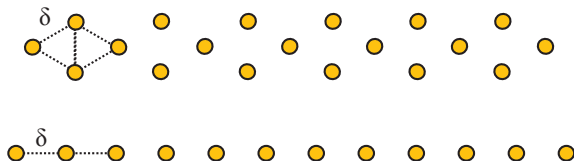
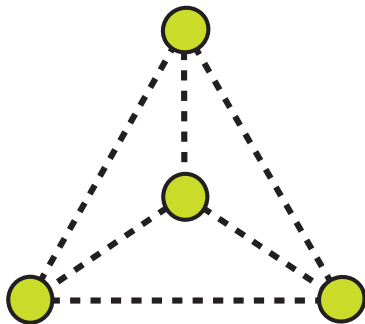


Figure :  $\mathfrak{C}^\Omega$ , with  $\Omega = \{\Delta_n(\delta)\}$  for increasing  $n$ , requires more and more density.

Related to **DBSCAN** (network clustering, and database mining)!

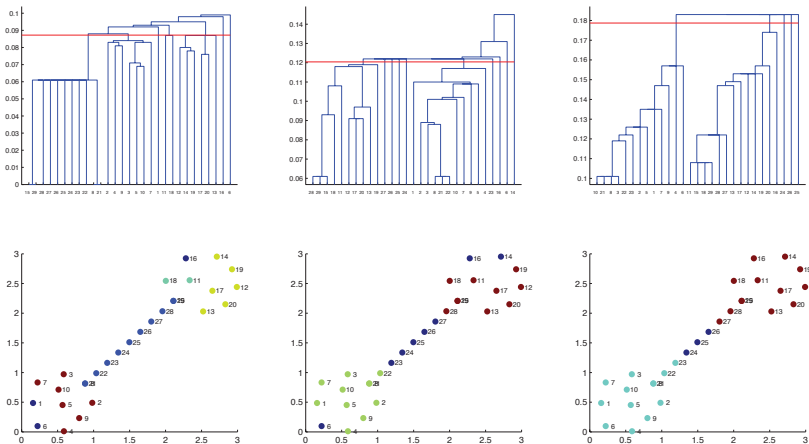
## More generality in $\mathcal{M}^{inj}$



**Figure :** A metric space that encodes a notion of “density” different from that provided by  $\Delta_3(\delta)$ .



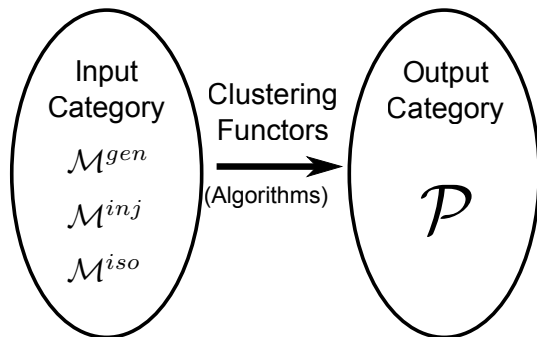
# Same thing for HC methods



**Figure :** Dendrograms arising from  $\mathfrak{H}^{\Delta m}$  applied to a randomly generated points in  $\mathbb{R}^2$ . Top: from left to right we show dendrograms corresponding to  $m = 2, 3$  and  $4$ . Bottom: partitioning induced on data by each dendrogram using parameters corresponding to red lines shown over dendrograms.

# Roadmap for classification

- Motivation
- Categories
- Framework
  - Input Categories
  - Output Categories
- Functors
- Results
  - $\mathcal{M}^{iso}$  (small)
  - $\mathcal{M}^{gen}$  (large)
  - $\mathcal{M}^{inj}$  (medium)
- **Discussion**



# Conclusion

- Motivation
- Categories
- Framework
  - Input Categories
  - Output Categories
- Functors
- Results
  - $\mathcal{M}^{iso}$  (small)
  - $\mathcal{M}^{gen}$  (large)
  - $\mathcal{M}^{inj}$  (medium)
- Conclusion

# Conclusion

- Kleinberg's nonexistence result.
- Motivation
- Categories
- Framework
  - Input Categories
  - Output Categories
- Functors
- Results
  - $\mathcal{M}^{iso}$  (small)
  - $\mathcal{M}^{gen}$  (large)
  - $\mathcal{M}^{inj}$  (medium)
- Conclusion

# Conclusion

- Motivation
  - Categories
  - Framework
    - Input Categories
    - Output Categories
  - Functors
  - Results
    - $\mathcal{M}^{iso}$  (small)
    - $\mathcal{M}^{gen}$  (large)
    - $\mathcal{M}^{inj}$  (medium)
  - Conclusion
- Kleinberg's nonexistence result.
  - category theory: language for **reasoning** about clustering.

# Conclusion

- Motivation
  - Categories
  - Framework
    - Input Categories
    - Output Categories
  - Functors
  - Results
    - $\mathcal{M}^{iso}$  (small)
    - $\mathcal{M}^{gen}$  (large)
    - $\mathcal{M}^{inj}$  (medium)
  - Conclusion
- Kleinberg's nonexistence result.
  - category theory: language for **reasoning** about clustering.
  - Future: better understanding of non-exciseive functors in  $\mathcal{M}^{inj}$ .

# Conclusion

- Motivation
  - Categories
  - Framework
    - Input Categories
    - Output Categories
  - Functors
  - Results
    - $\mathcal{M}^{iso}$  (small)
    - $\mathcal{M}^{gen}$  (large)
    - $\mathcal{M}^{inj}$  (medium)
  - Conclusion
- Kleinberg's nonexistence result.
  - category theory: language for **reasoning** about clustering.
  - Future: better understanding of non-exciseve functors in  $\mathcal{M}^{inj}$ .
  - note: also established results for **hierarchical clustering**

# Conclusion

- Motivation
  - Categories
  - Framework
    - Input Categories
    - Output Categories
  - Functors
  - Results
    - $\mathcal{M}^{iso}$  (small)
    - $\mathcal{M}^{gen}$  (large)
    - $\mathcal{M}^{inj}$  (medium)
  - Conclusion
- Kleinberg's nonexistence result.
  - category theory: language for **reasoning** about clustering.
  - Future: better understanding of non-exciseive functors in  $\mathcal{M}^{inj}$ .
  - note: also established results for **hierarchical clustering**
  - In context of HC: metric structures on  $\mathcal{M}$  lead to **stability** and **convergence** (Gromov-Hausdorff distances).



# Conclusion

- Motivation
  - Categories
  - Framework
    - Input Categories
    - Output Categories
  - Functors
  - Results
    - $\mathcal{M}^{iso}$  (small)
    - $\mathcal{M}^{gen}$  (large)
    - $\mathcal{M}^{inj}$  (medium)
  - Conclusion
- Kleinberg's nonexistence result.
  - category theory: language for **reasoning** about clustering.
  - Future: better understanding of non-exciseive functors in  $\mathcal{M}^{inj}$ .
  - note: also established results for **hierarchical clustering**
  - In context of HC: metric structures on  $\mathcal{M}$  lead to **stability** and **convergence** (Gromov-Hausdorff distances).
  - Paper in FoCM, 2013. Also JMLR, 2010.
  - Ongoing work on extending to **networks**. Joint with Carlsson, Ribeiro, Segarra. Some unexpected results.

# Acknowledgements

- DARPA grant HR0011-05-1-0007
- ONR grant N00014-09-1-0783
- Stanford University, The Ohio State University.
- **Ryan Lewis** for help preparing the slides!

Any questions?