

Density-based Cluster Algorithms in Low-dimensional and High-dimensional Applications

Benno Stein

Bauhaus University Weimar
Web-based Information Systems

Michael Busch

IBM Silicon Valley Laboratory
(WebSphere II Project)

Introduction

Cluster
Algorithms

Density-based
Algorithms

Analysis

Summary

Introduction

Given a set of objects (documents) $D = \{d_1, \dots, d_n\}$.

Clustering is the *unsupervised* classification of d_i into groups.

Result is a partitioning \mathcal{C} of D .

Introduction

Cluster
Algorithms

Density-based
Algorithms

Analysis

Summary

Introduction

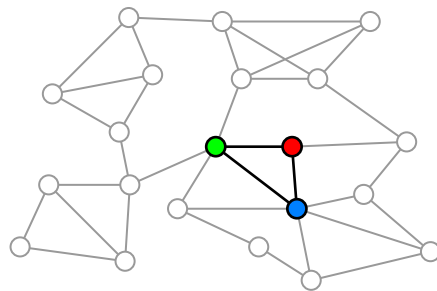
Given a set of objects (documents) $D = \{d_1, \dots, d_n\}$.

Clustering is the *unsupervised* classification of d_i into groups.

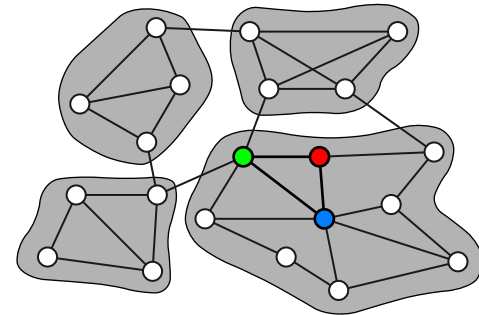
Result is a partitioning \mathcal{C} of D .

Objective: Maximize intra-group similarity.

Minimize inter-group similarity.



Similarity graph



Clustered graph

Introduction

Cluster Algorithms

Density-based Algorithms

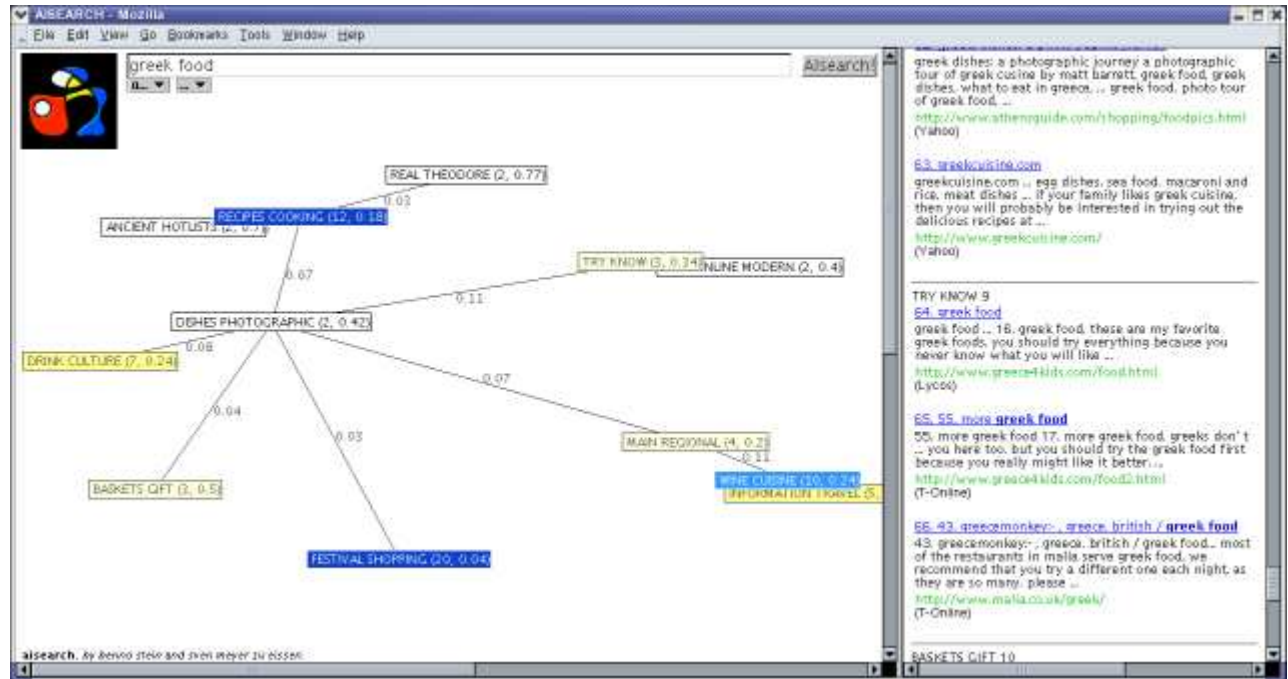
Analysis

Summary

Introduction

Cluster algorithms form the backbone of document categorization.

Example AIssearch [www.aishsearch.de] :



Introduction

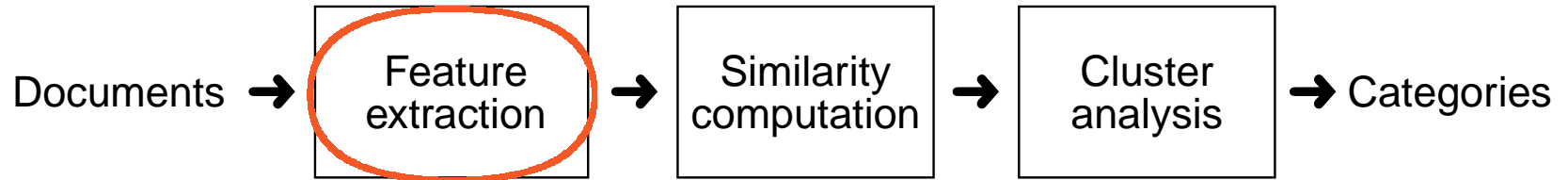
Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Introduction



Indexing (includes parsing, stopword elimination, stemming):

Vector representation
with weighting scheme:

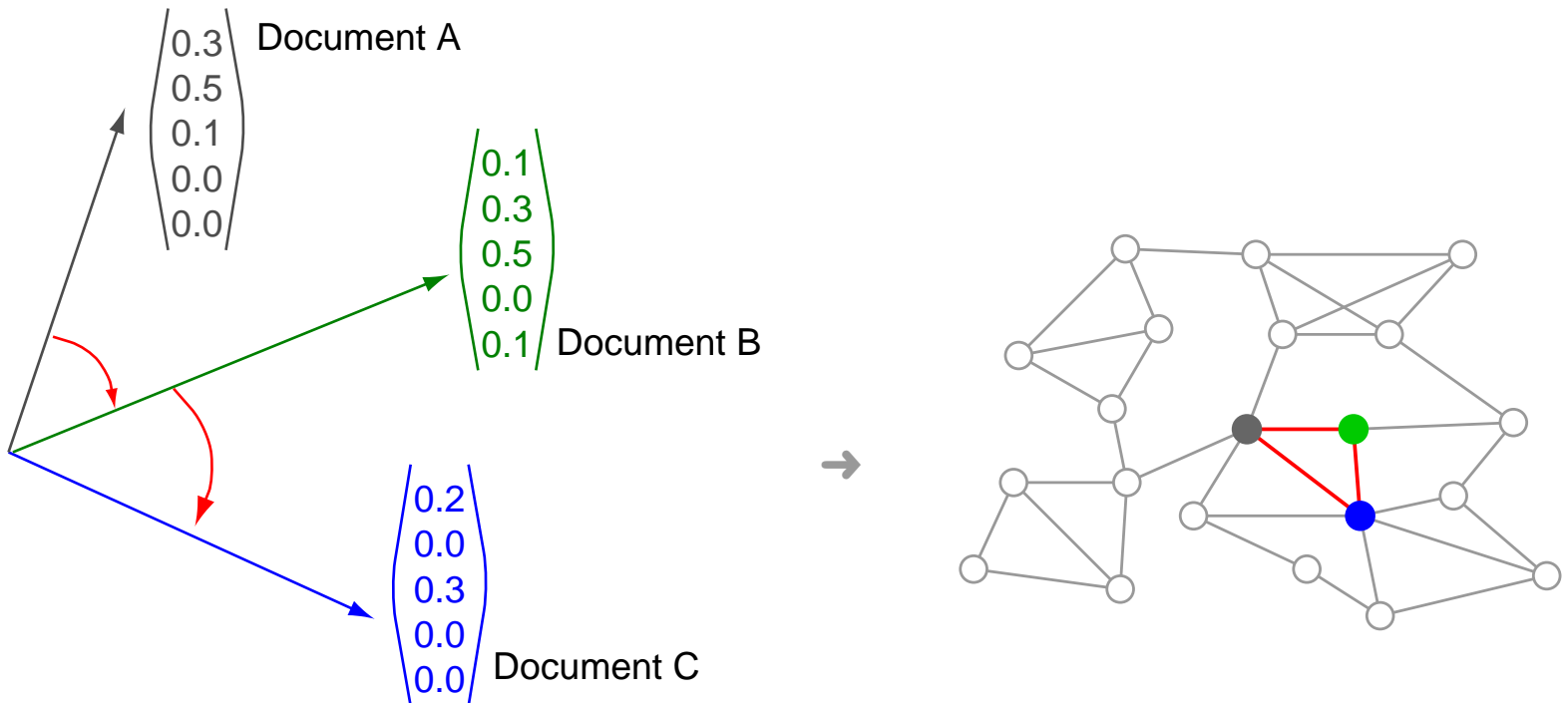
$$\underbrace{\#abs}_{tf} \cdot \log \frac{|X|}{\underbrace{\#docs}_{idf}}$$

chrysler	0.12
deal	0.2
leav	0.1
amc	0.01
cat	0.0
sal	0.01
dog	0.0
⋮	⋮

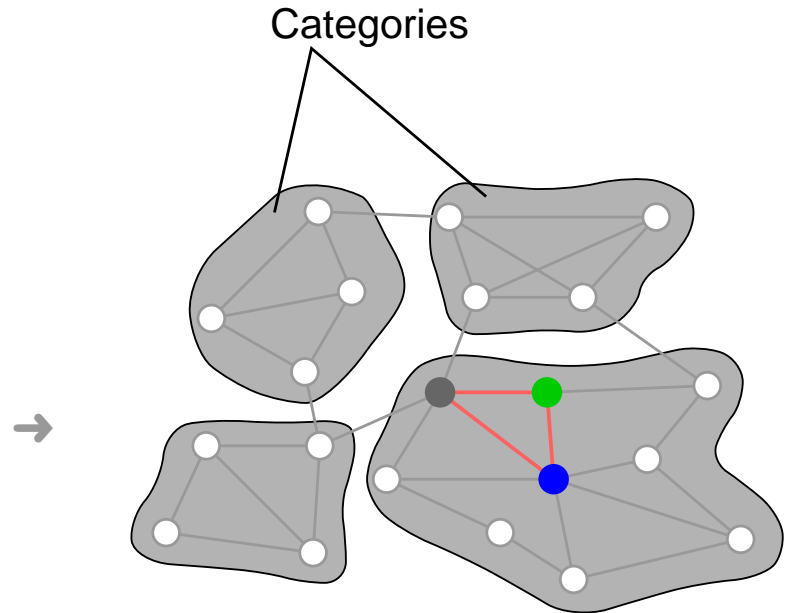
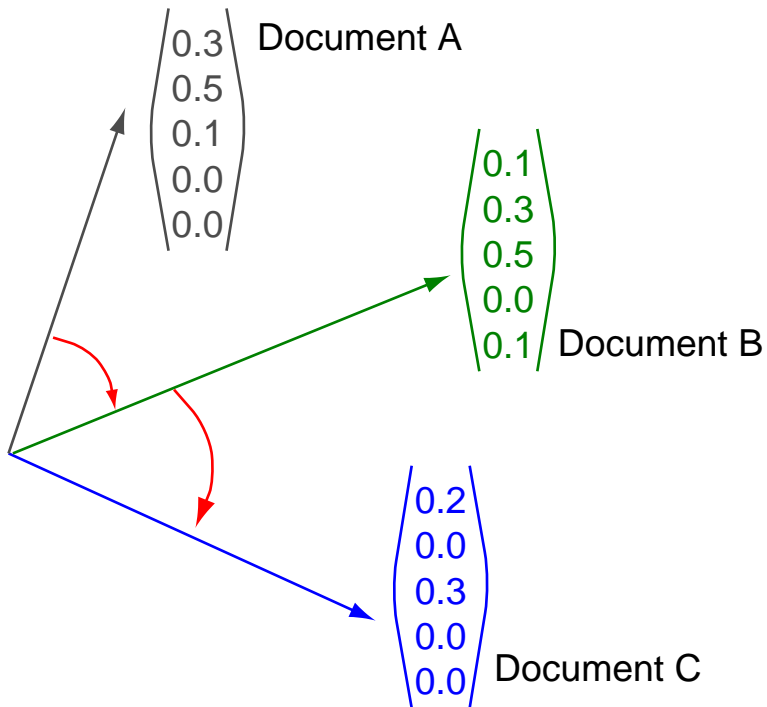
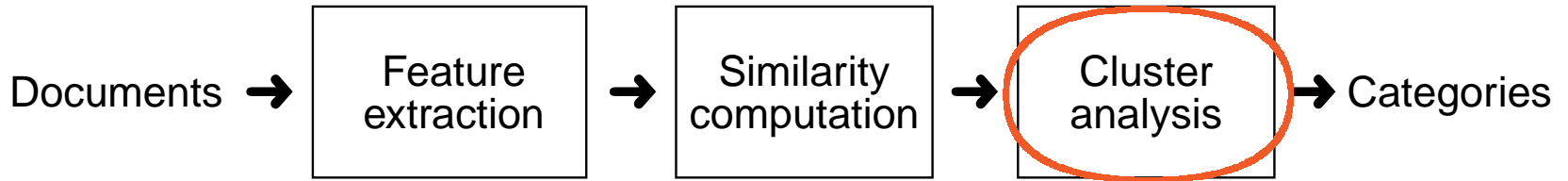
Introduction



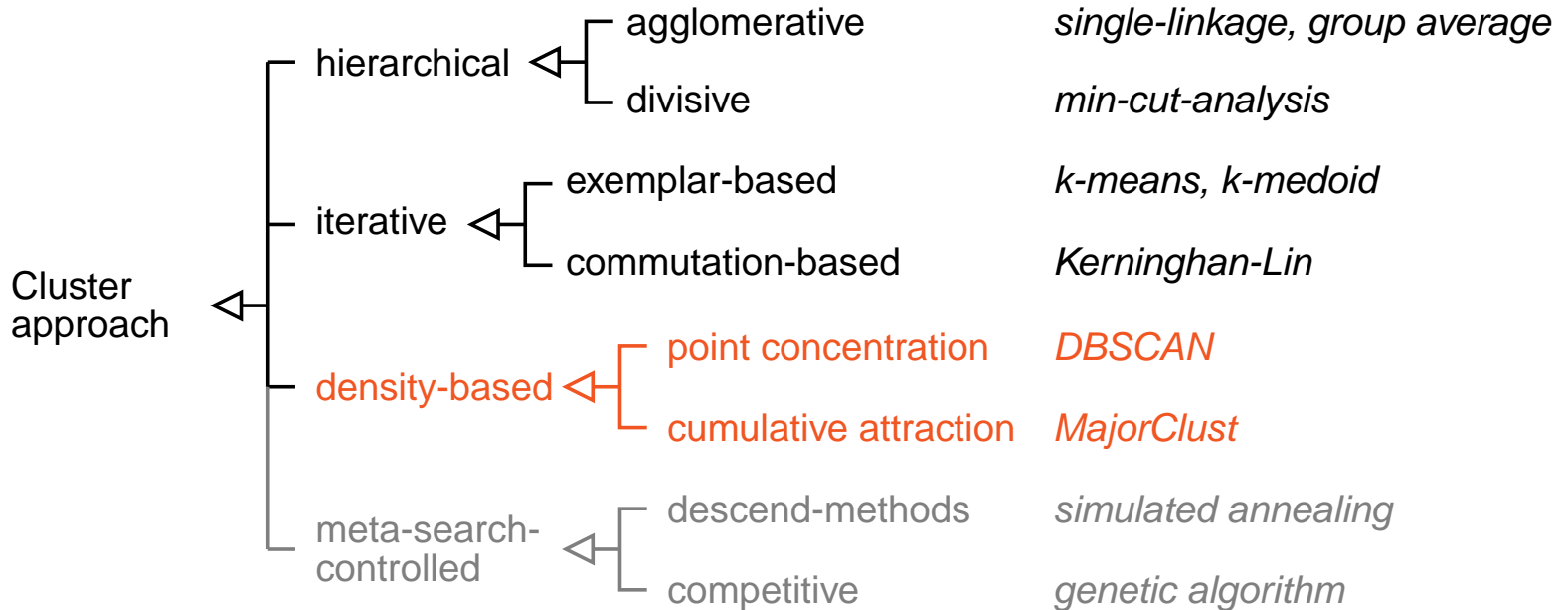
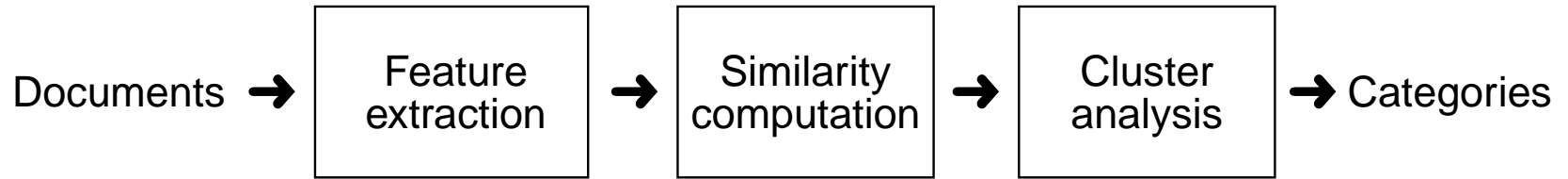
e. g. under the vector space model:



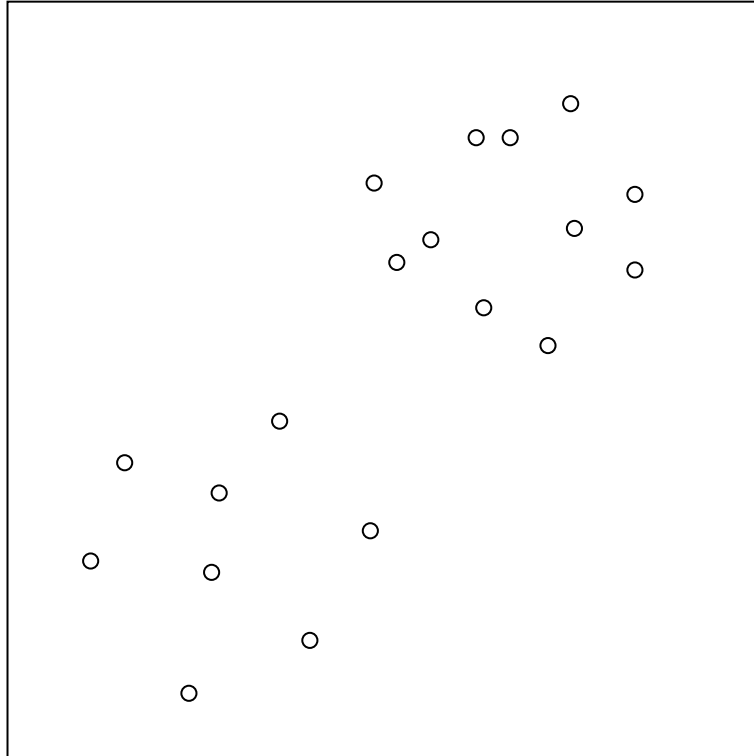
Introduction



Introduction



Hierarchical agglomerative: Single-linkage



Introduction

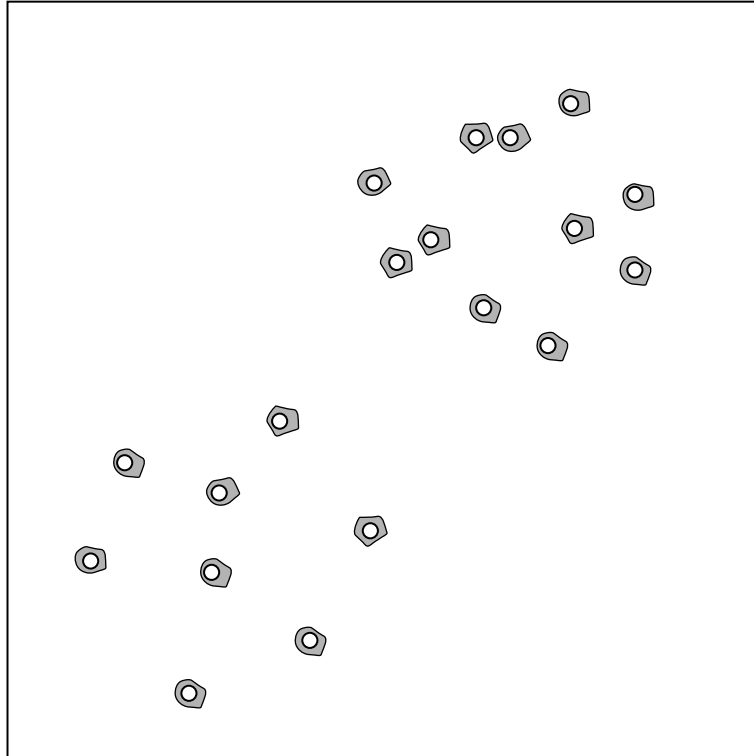
**Cluster
Algorithms**

Density-based
Algorithms

Analysis

Summary

Hierarchical agglomerative: Single-linkage



Introduction

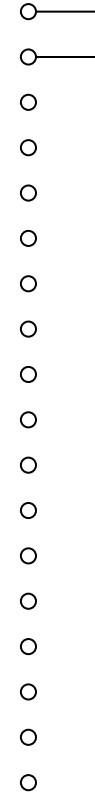
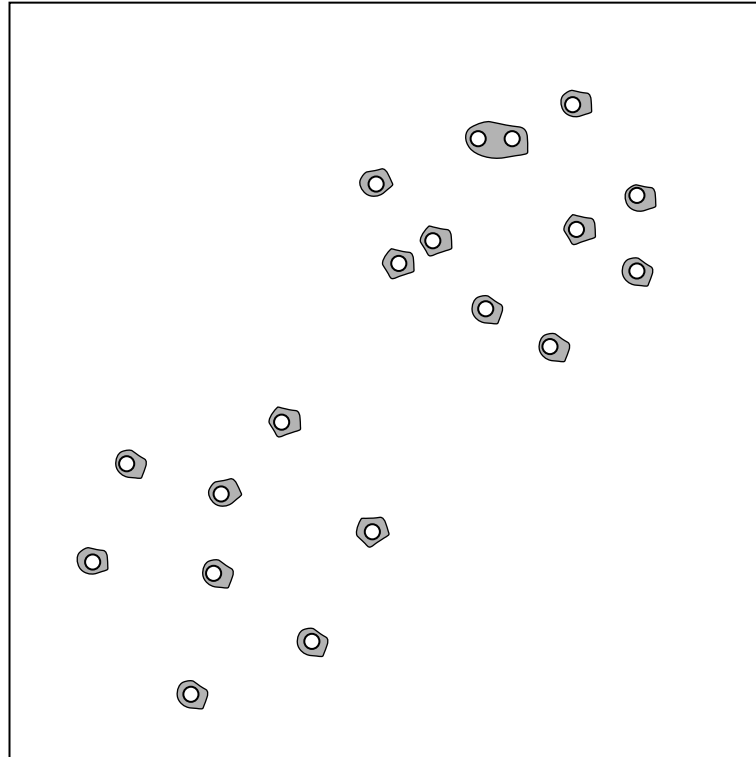
**Cluster
Algorithms**

Density-based
Algorithms

Analysis

Summary

Hierarchical agglomerative: Single-linkage



→ Distanz

Introduction

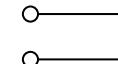
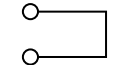
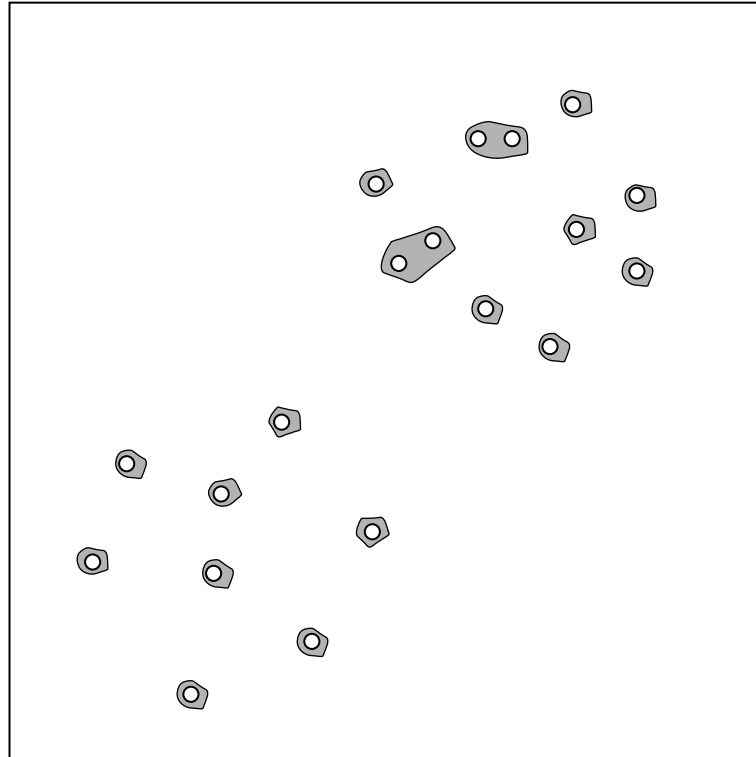
Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Hierarchical agglomerative: Single-linkage



—————→ Distanz

Introduction

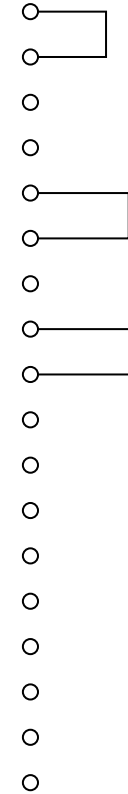
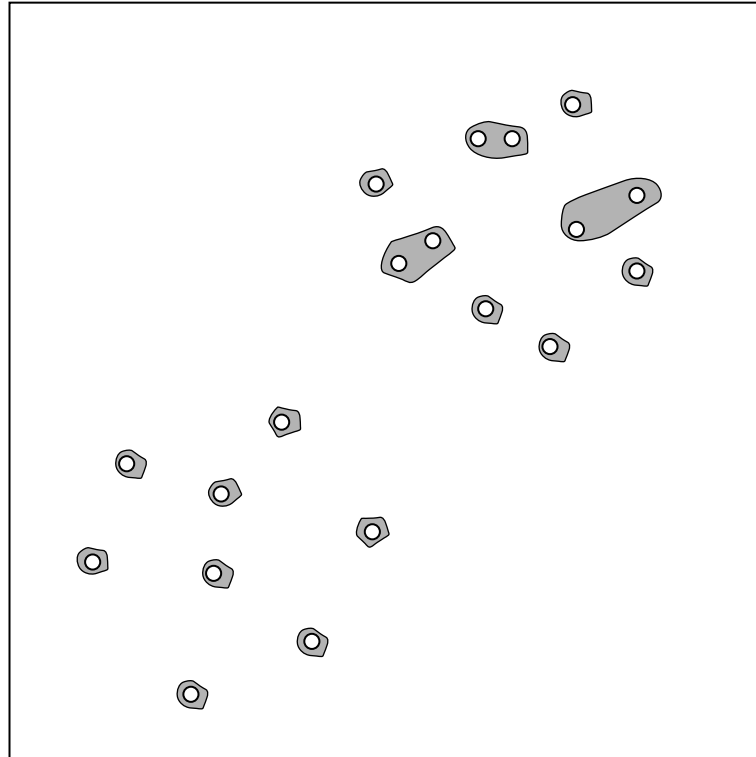
Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Hierarchical agglomerative: Single-linkage



→ Distanz

Introduction

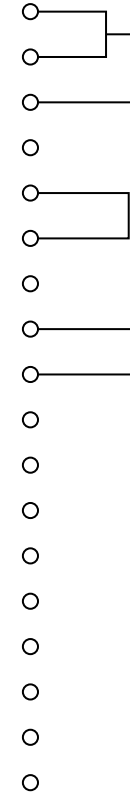
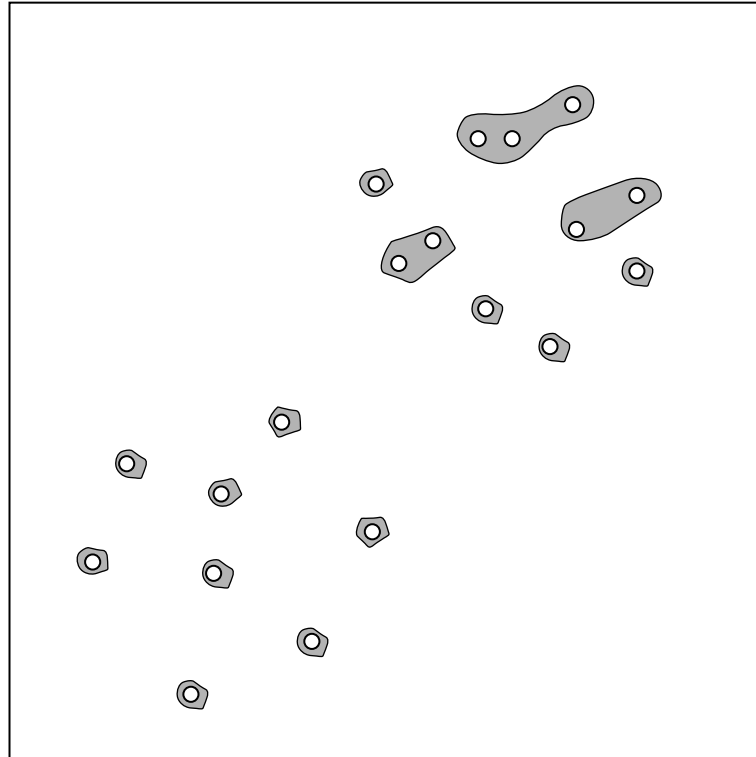
Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Hierarchical agglomerative: Single-linkage



→ Distanz

Introduction

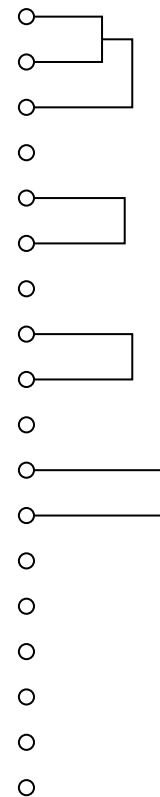
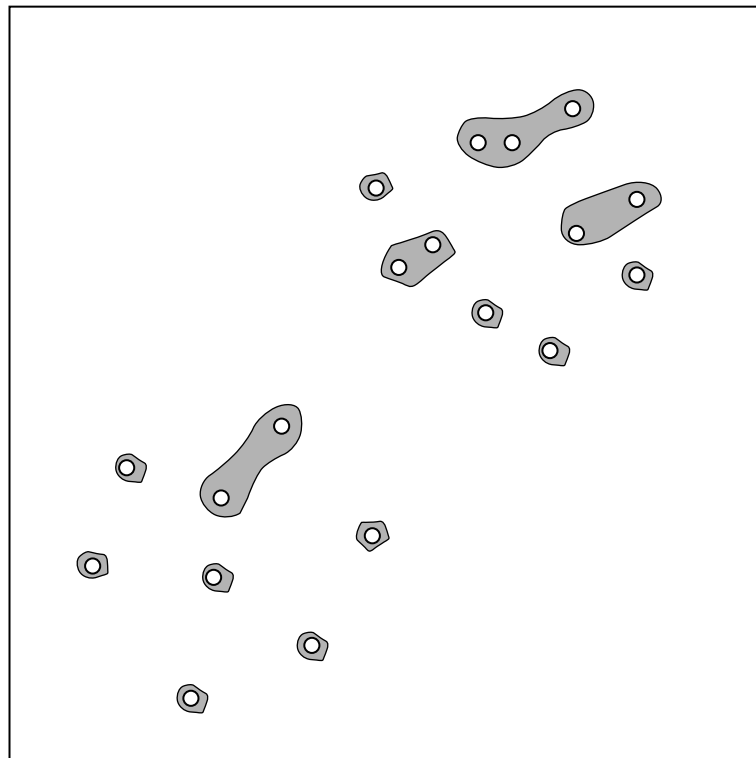
Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Hierarchical agglomerative: Single-linkage



→ Distanz

Introduction

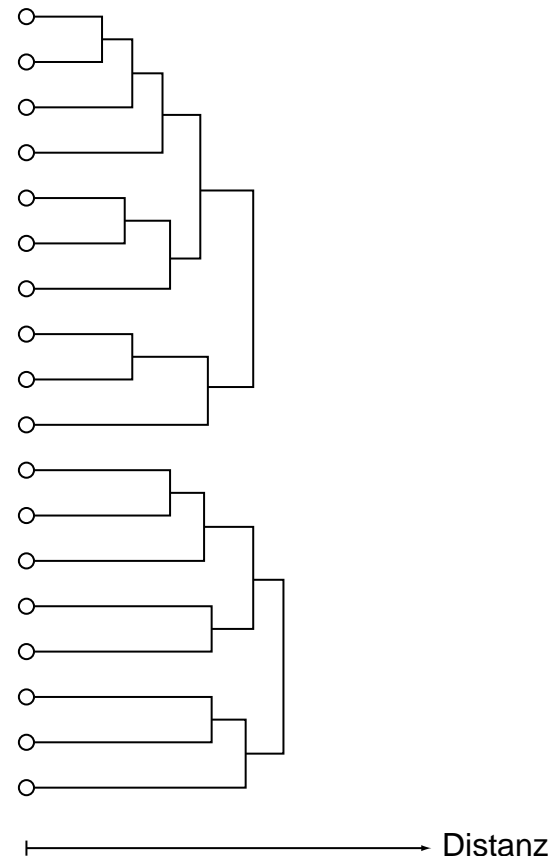
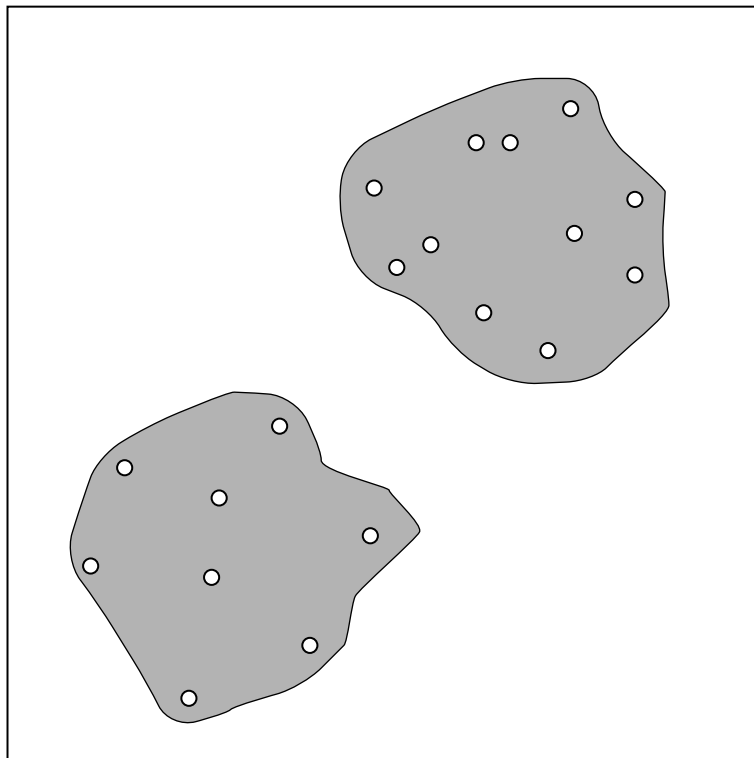
Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Hierarchical agglomerative: Single-linkage



Introduction

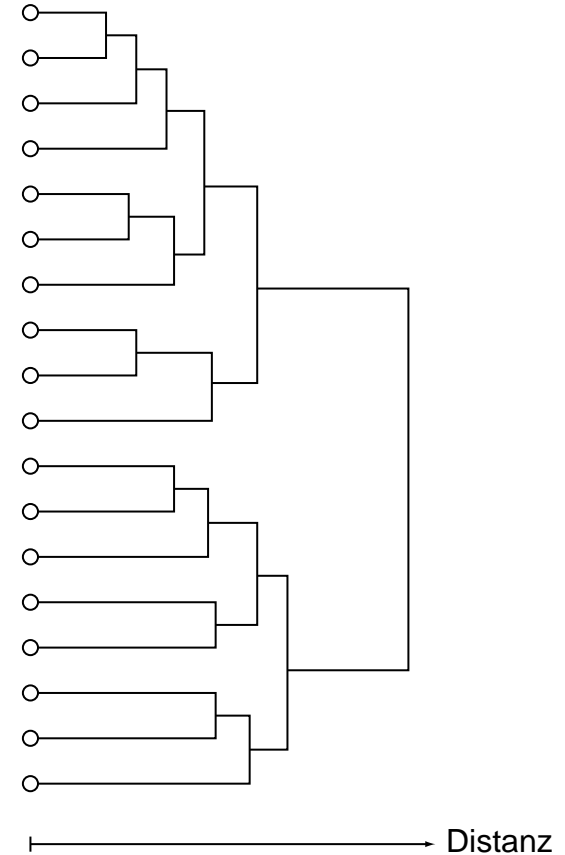
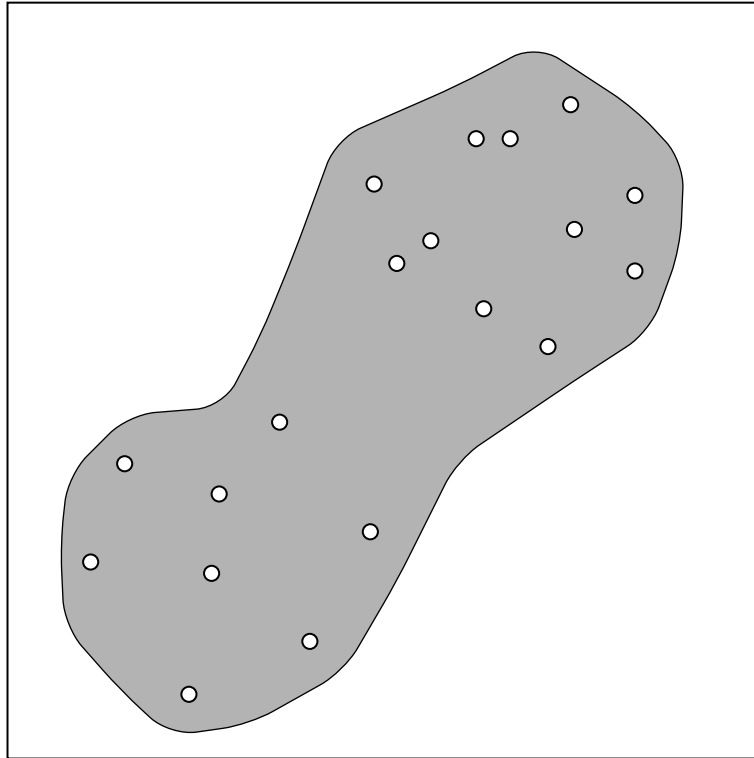
Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Hierarchical agglomerative: Single-linkage



Introduction

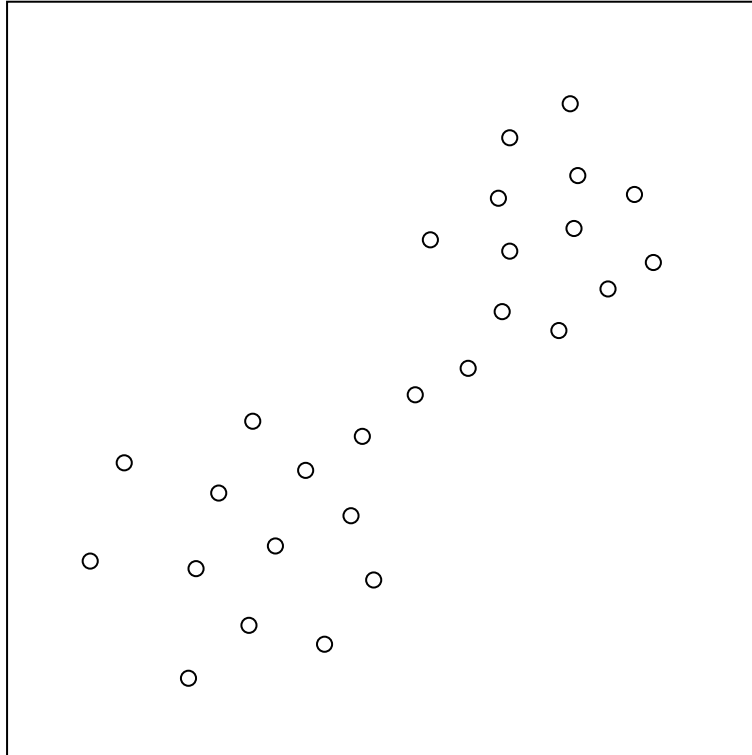
Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Single-linkage: Chaining Problem



Introduction

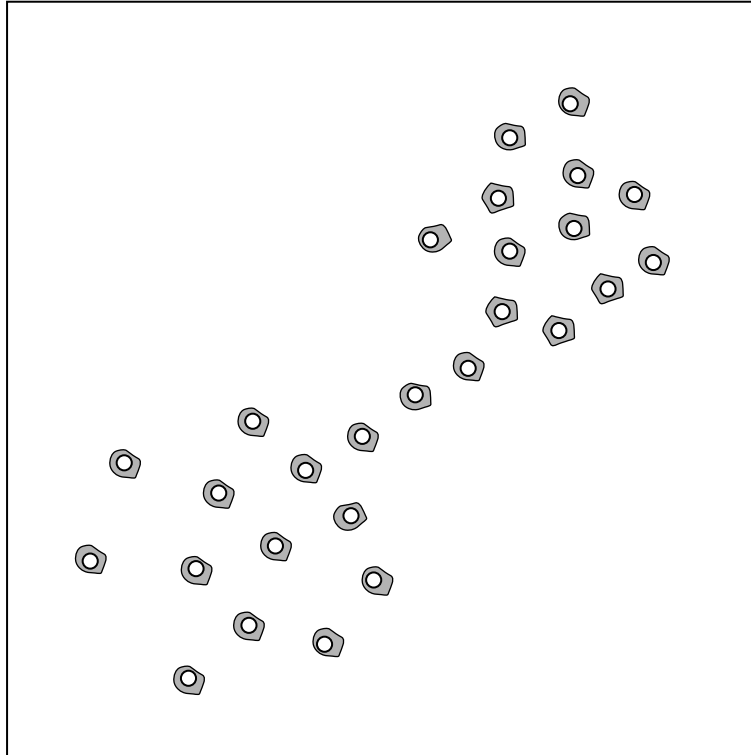
**Cluster
Algorithms**

Density-based
Algorithms

Analysis

Summary

Single-linkage: Chaining Problem



Introduction

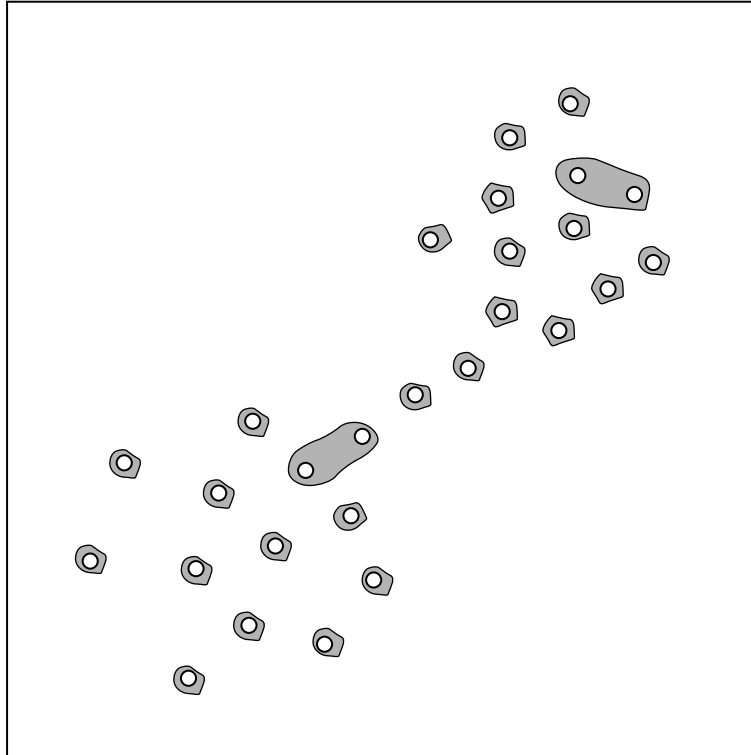
**Cluster
Algorithms**

Density-based
Algorithms

Analysis

Summary

Single-linkage: Chaining Problem



Introduction

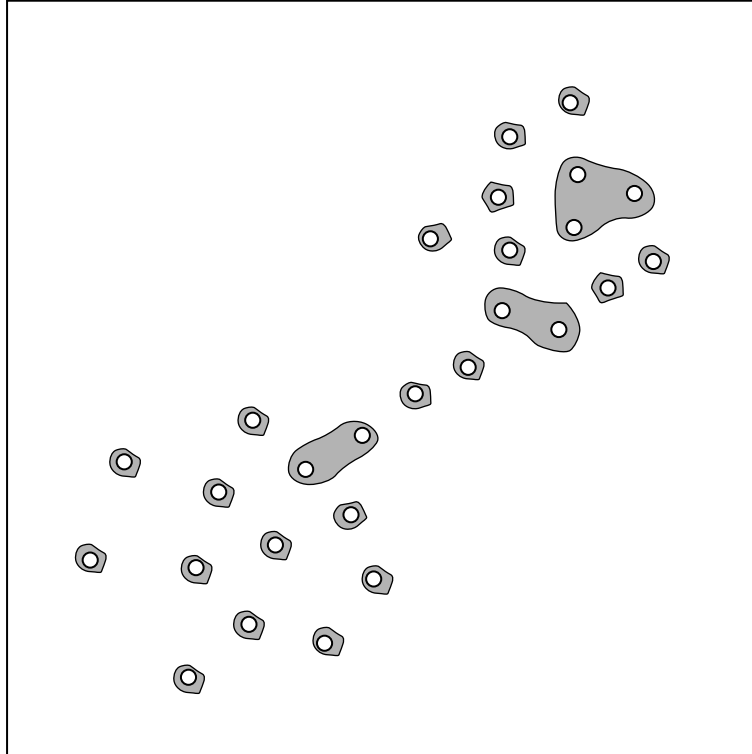
**Cluster
Algorithms**

Density-based
Algorithms

Analysis

Summary

Single-linkage: Chaining Problem



Introduction

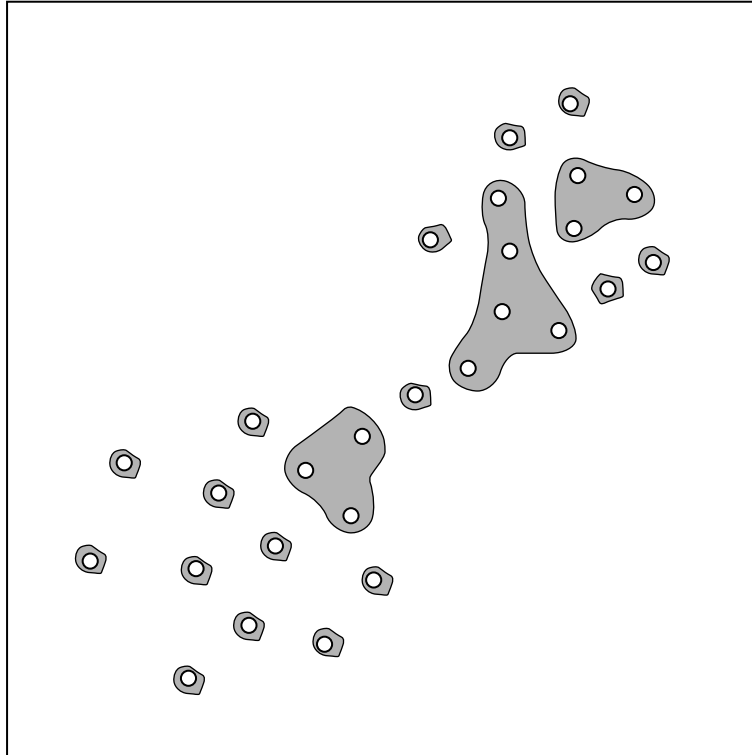
**Cluster
Algorithms**

Density-based
Algorithms

Analysis

Summary

Single-linkage: Chaining Problem



Introduction

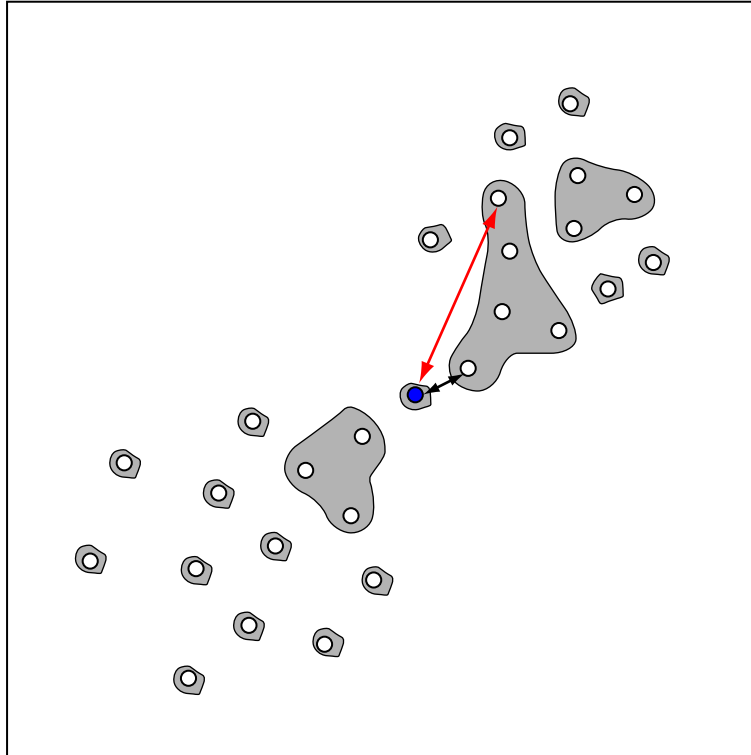
**Cluster
Algorithms**

Density-based
Algorithms

Analysis

Summary

Single-linkage: Chaining Problem



Introduction

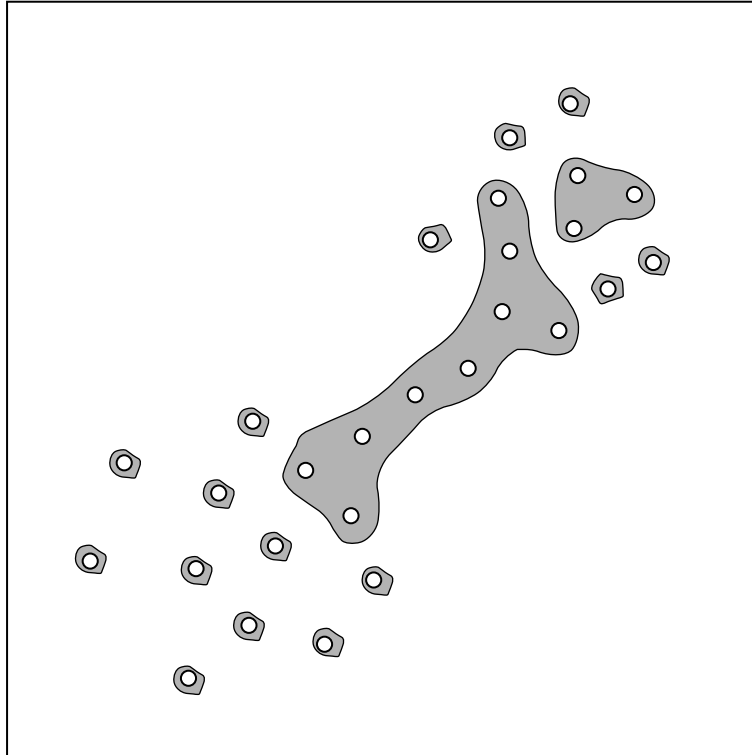
**Cluster
Algorithms**

Density-based
Algorithms

Analysis

Summary

Single-linkage: Chaining Problem



Introduction

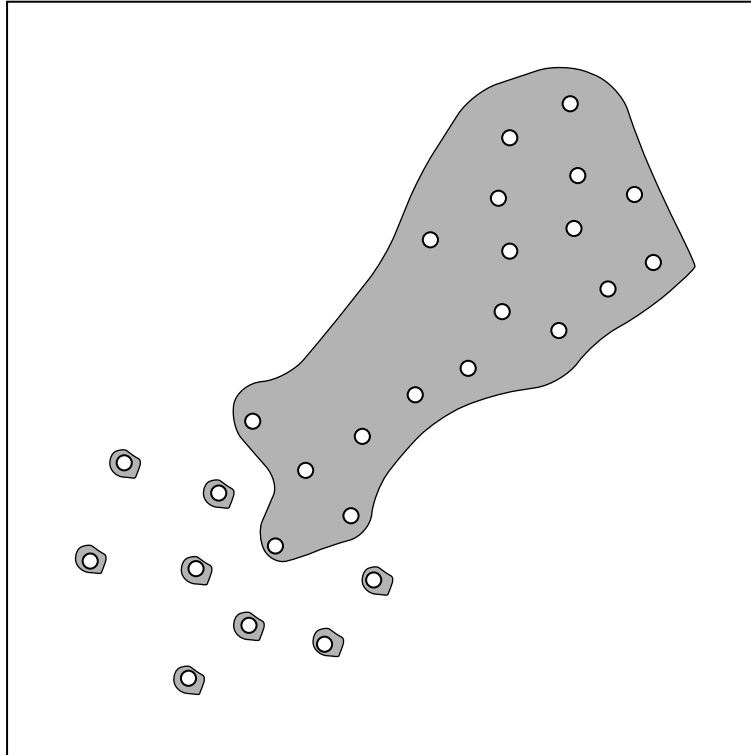
**Cluster
Algorithms**

Density-based
Algorithms

Analysis

Summary

Single-linkage: Chaining Problem



Introduction

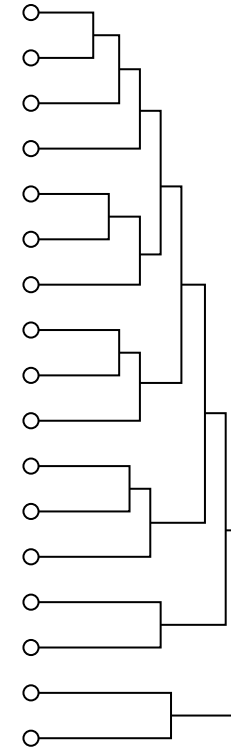
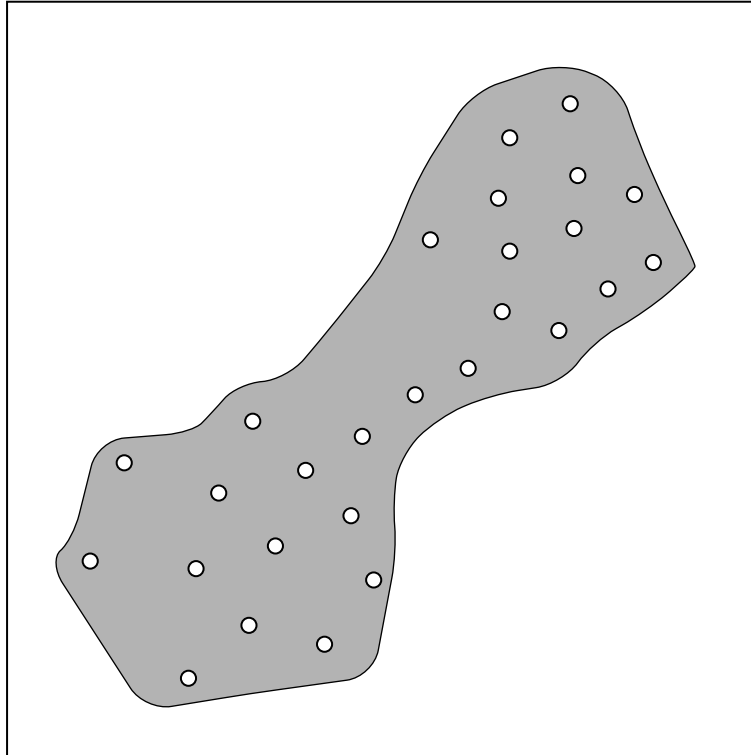
**Cluster
Algorithms**

Density-based
Algorithms

Analysis

Summary

Single-linkage: Chaining Problem



→ Distanz

Introduction

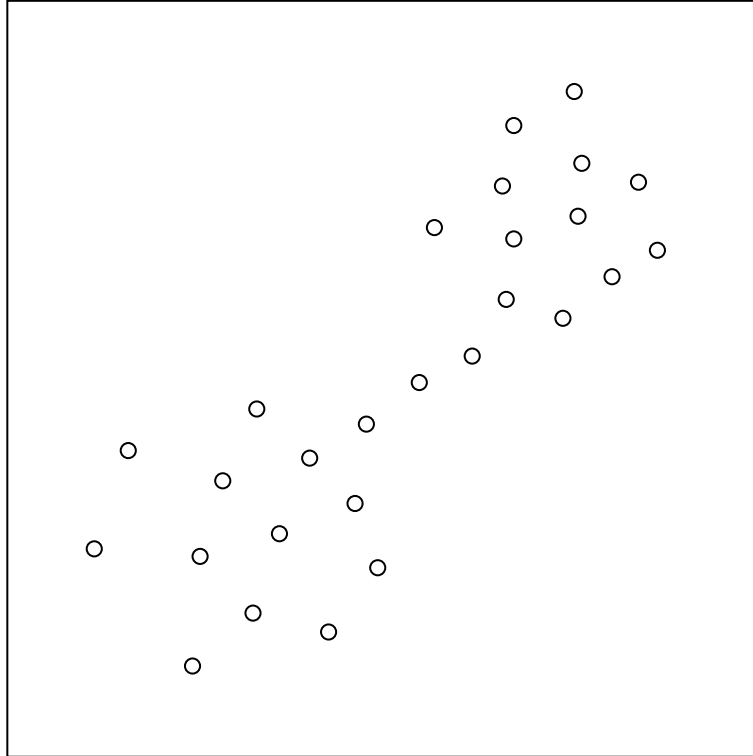
Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Exemplar-based algorithm: k -Means



Introduction

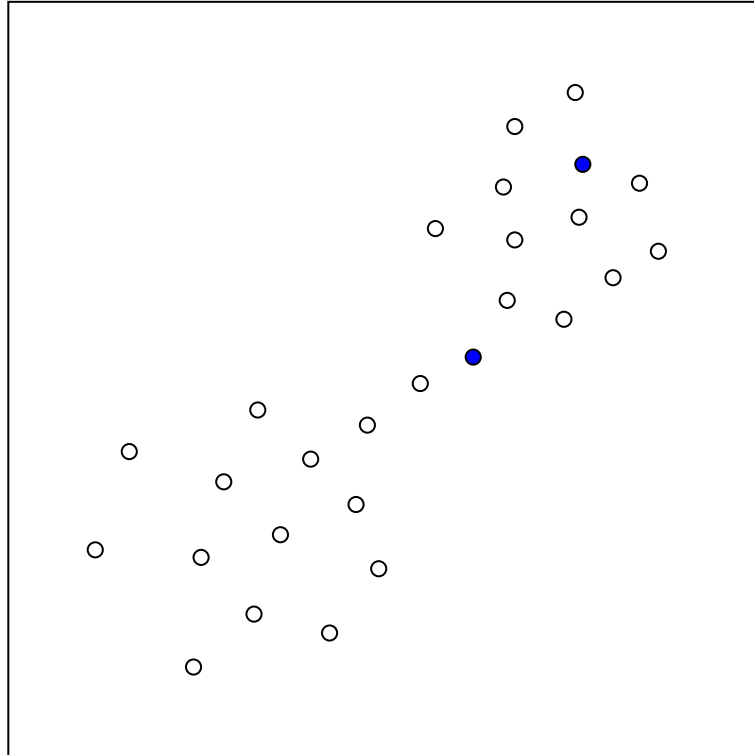
**Cluster
Algorithms**

Density-based
Algorithms

Analysis

Summary

Exemplar-based algorithm: k -Means



Introduction

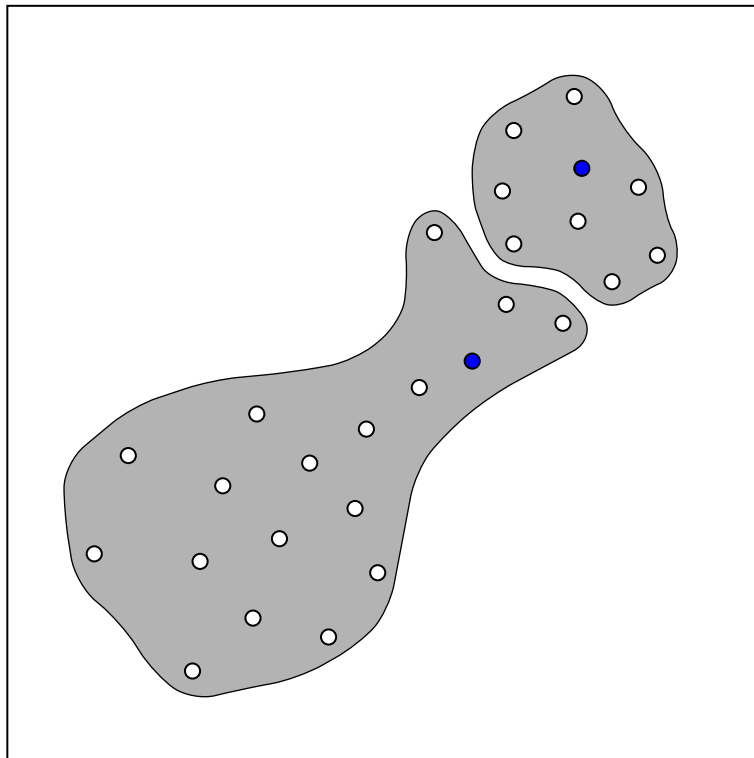
Cluster
Algorithms

Density-based
Algorithms

Analysis

Summary

Exemplar-based algorithm: k -Means



Introduction

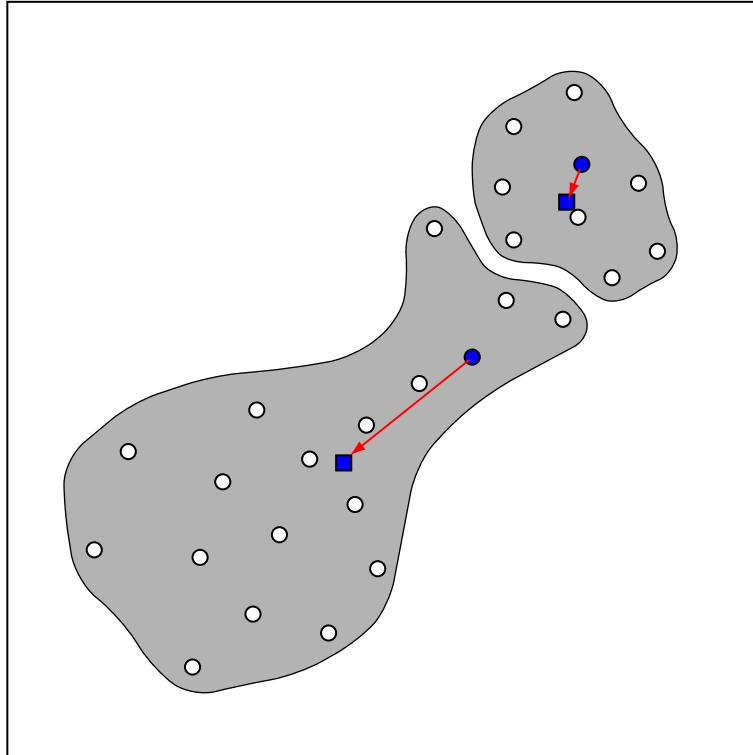
Cluster
Algorithms

Density-based
Algorithms

Analysis

Summary

Exemplar-based algorithm: k -Means



Introduction

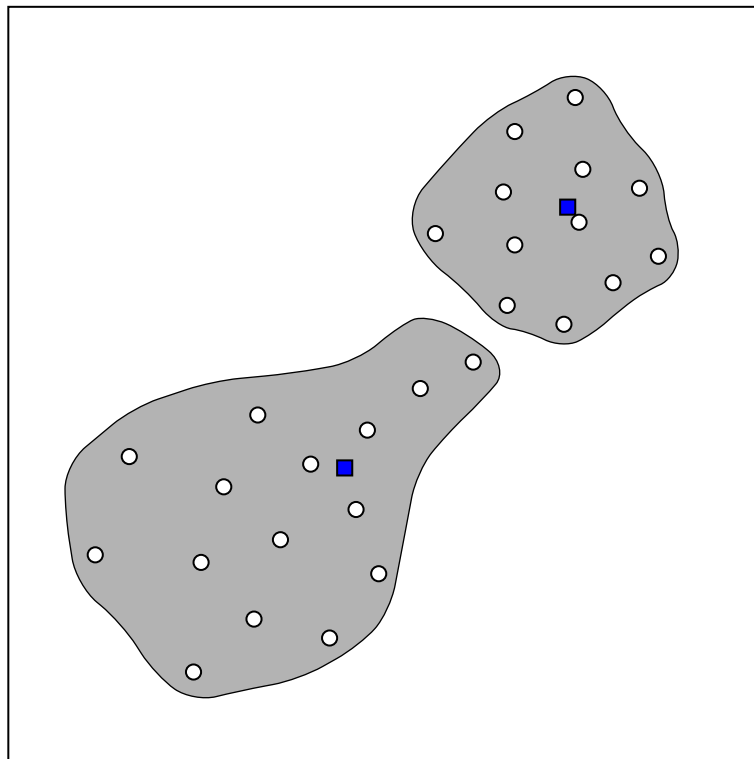
Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Exemplar-based algorithm: k -Means



Introduction

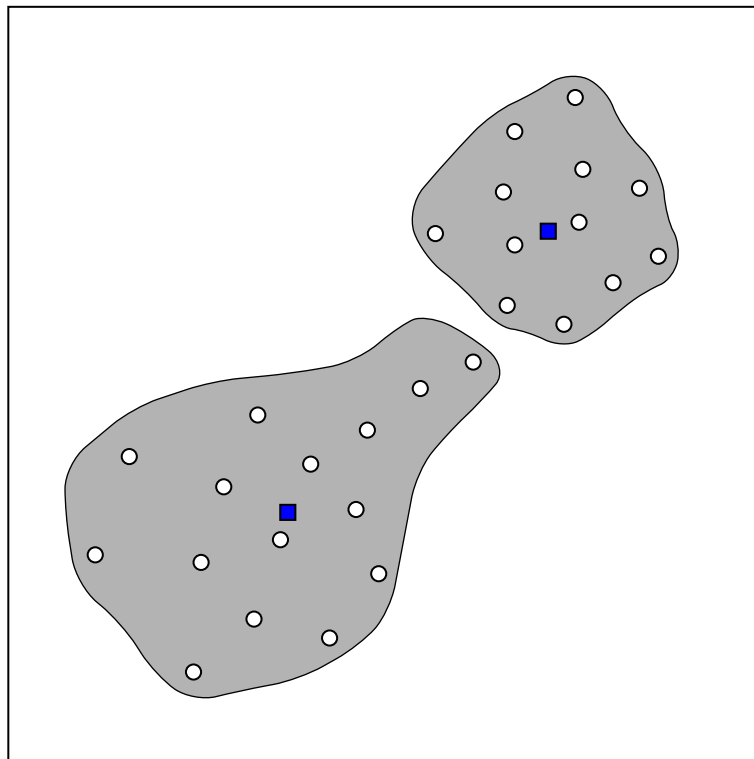
Cluster
Algorithms

Density-based
Algorithms

Analysis

Summary

Exemplar-based algorithm: k -Means



Introduction

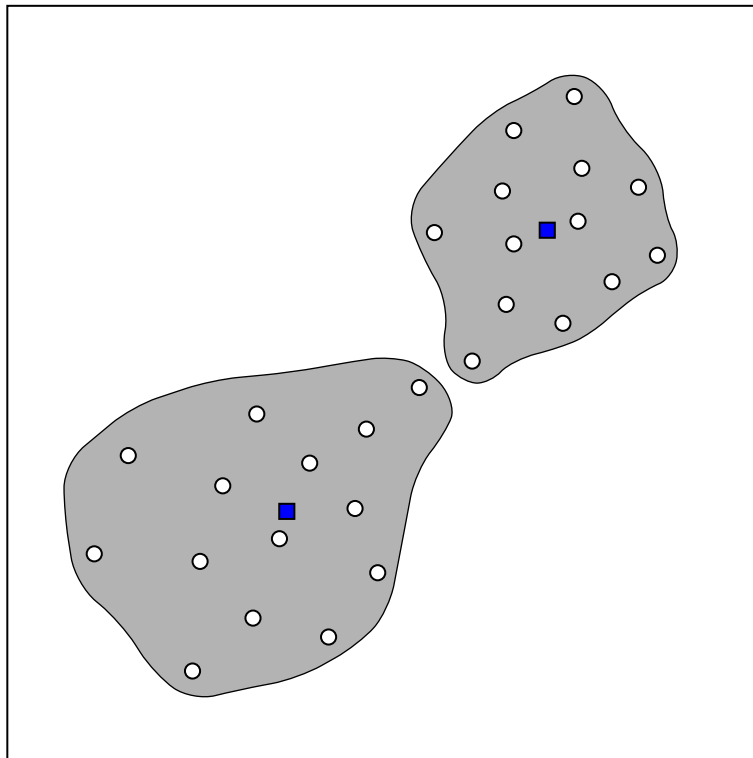
Cluster
Algorithms

Density-based
Algorithms

Analysis

Summary

Exemplar-based algorithm: k -Means



Introduction

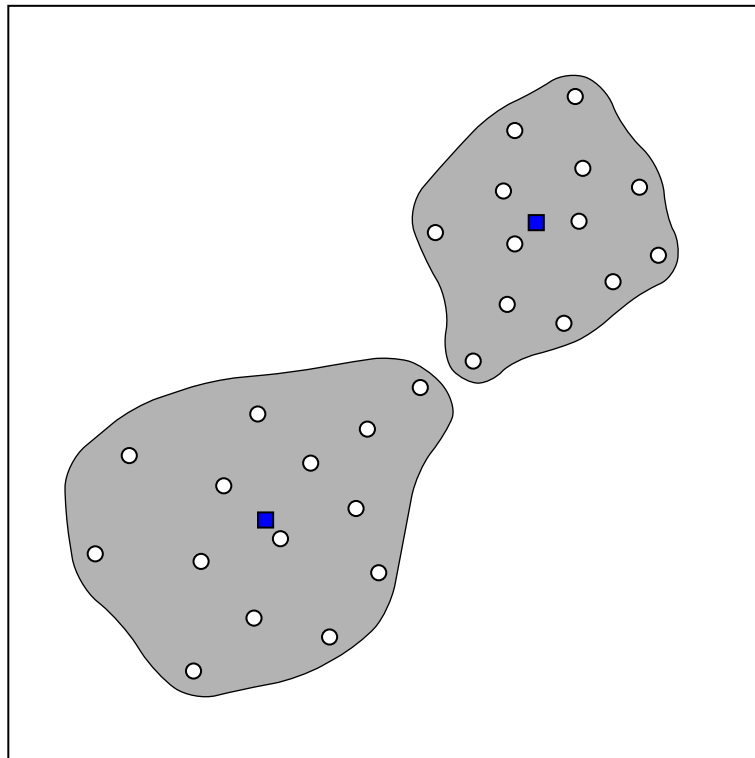
Cluster
Algorithms

Density-based
Algorithms

Analysis

Summary

Exemplar-based algorithm: k -Means



Introduction

Cluster Algorithms

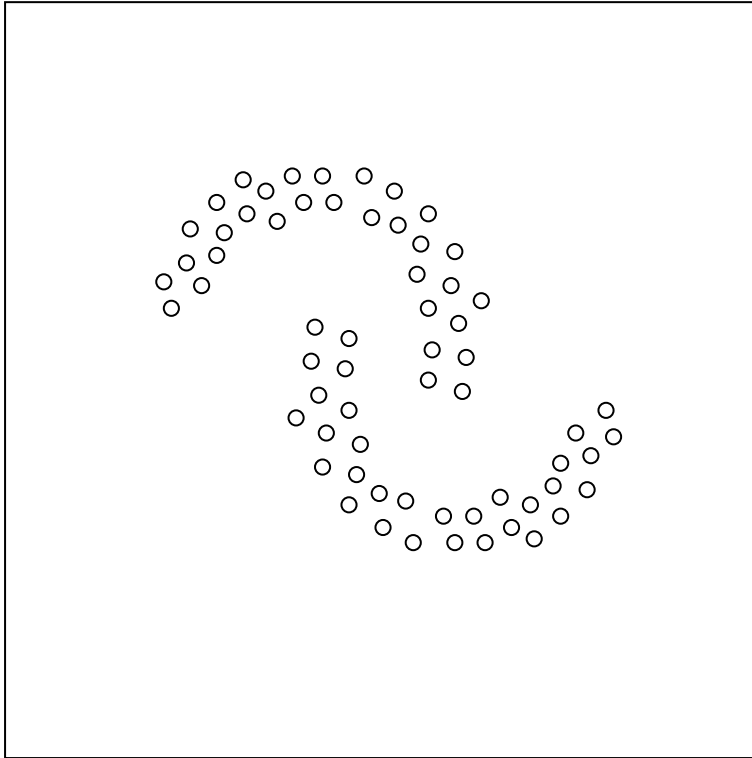
Density-based Algorithms

Analysis

Summary

Exemplar-based versus Linkage

Exemplar-based algorithms fail with large differences in size.



Introduction

**Cluster
Algorithms**

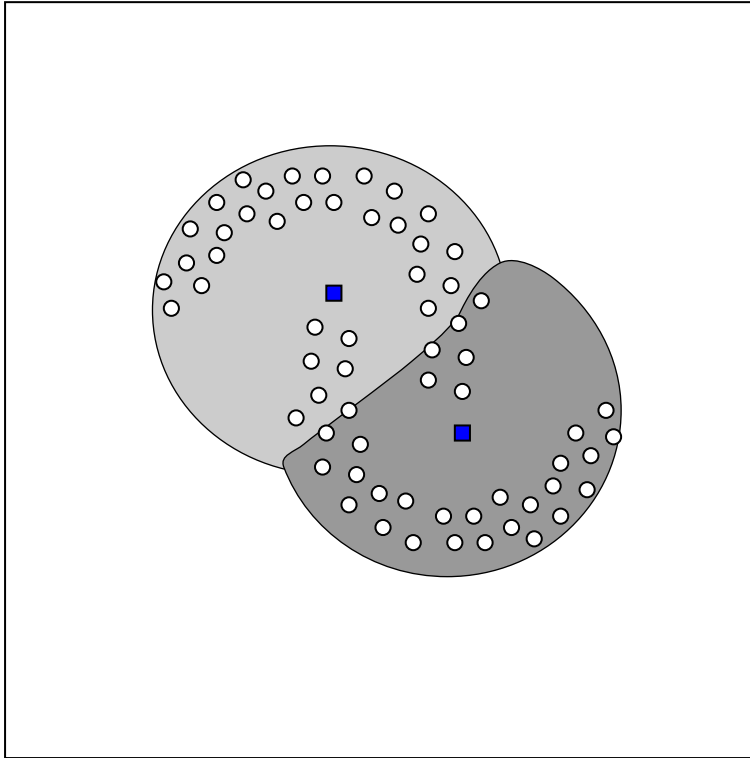
Density-based
Algorithms

Analysis

Summary

Exemplar-based versus Linkage

Exemplar-based algorithms fail with entwined clusters.



Introduction

Cluster
Algorithms

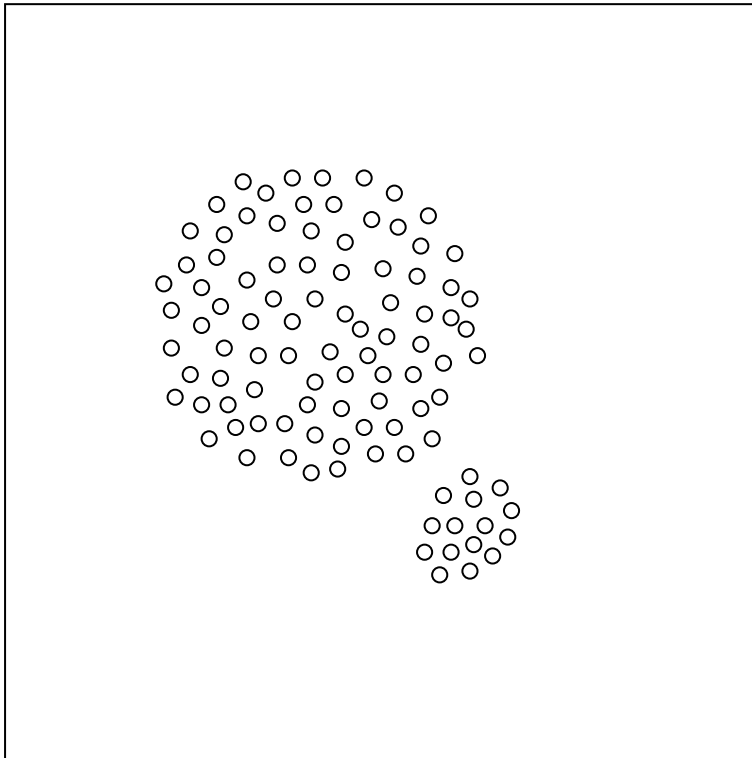
Density-based
Algorithms

Analysis

Summary

Exemplar-based versus Linkage

Exemplar-based algorithms fail with entwined clusters.



Introduction

**Cluster
Algorithms**

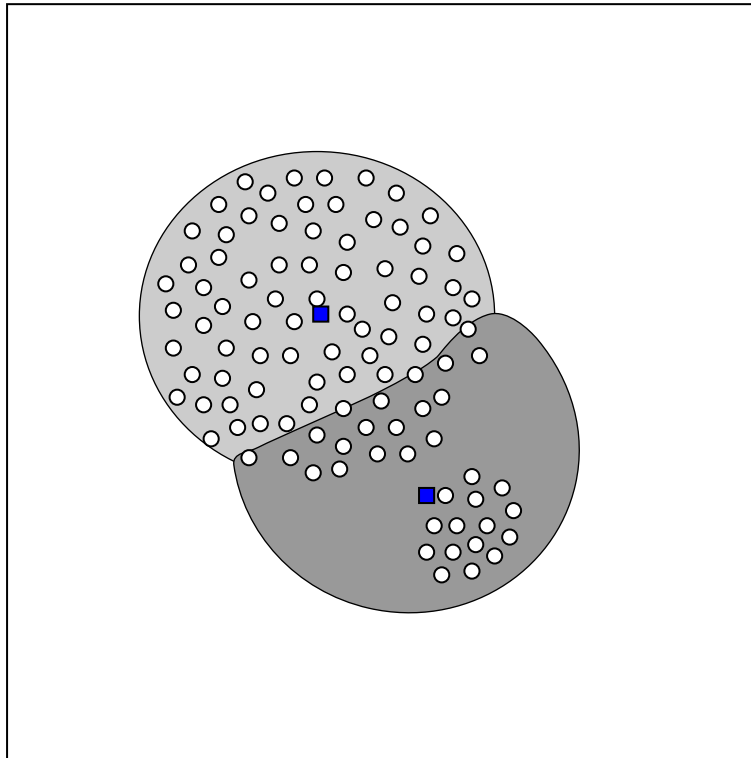
Density-based
Algorithms

Analysis

Summary

Exemplar-based versus Linkage

Exemplar-based algorithms fail with entwined clusters.



Introduction

Cluster
Algorithms

Density-based
Algorithms

Analysis

Summary

Density-based Cluster Analysis

Density-based algorithms try to separate the set D into subsets of similar densities.

Density estimation can happen

- parameter-based: the underlying distribution is known
- parameter-less: histogram, kernel function
(construct bar charts, superimpose continuous functions)

Introduction

Cluster
Algorithms

Density-based
Algorithms

Analysis

Summary

Density-based Cluster Analysis

Density-based algorithms try to separate the set D into subsets of similar densities.

Density estimation can happen

- parameter-based: the underlying distribution is known
- parameter-less: histogram, kernel function
(construct barcharts, superimpose continuous functions)

Example (Carribean Islands):



Introduction

Cluster Algorithms

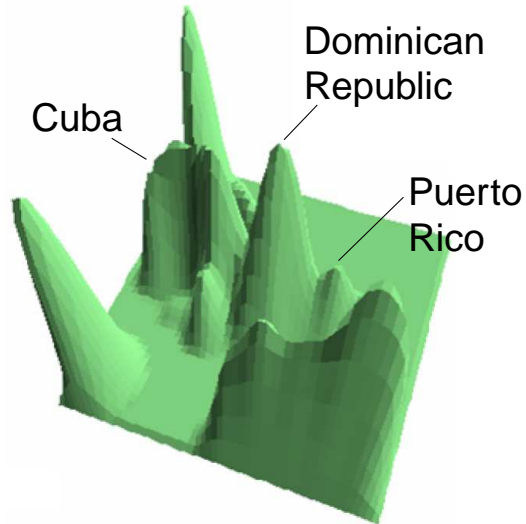
Density-based Algorithms

Analysis

Summary

Density-based Cluster Analysis

Density estimation with Gaussian Kernel for the example.



Introduction

Cluster
Algorithms

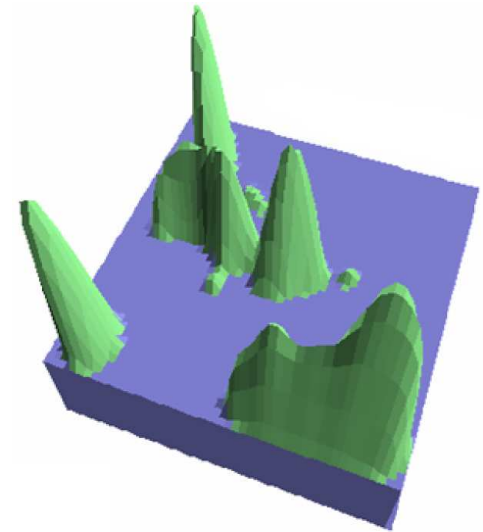
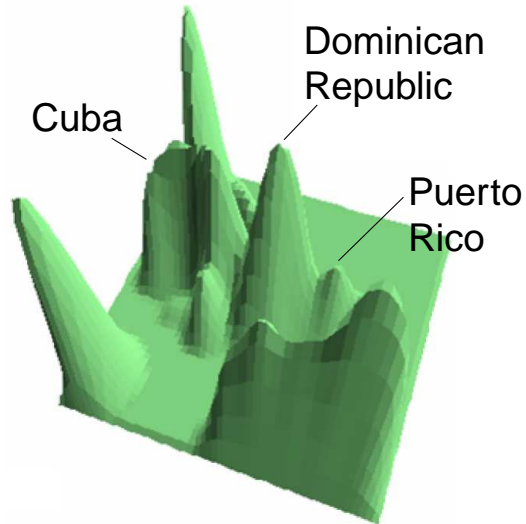
**Density-based
Algorithms**

Analysis

Summary

Density-based Cluster Analysis

Density estimation with Gaussian Kernel for the example.



Introduction

Cluster Algorithms

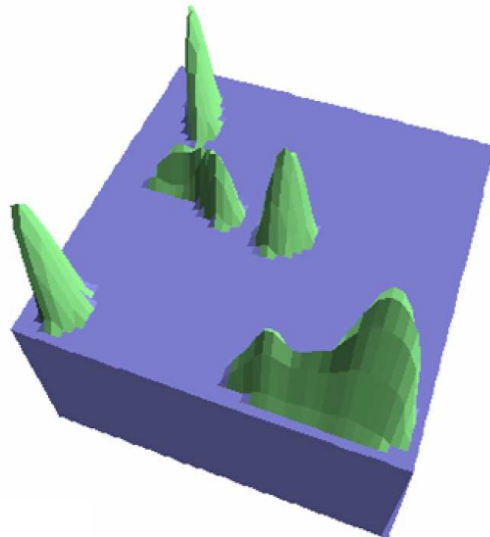
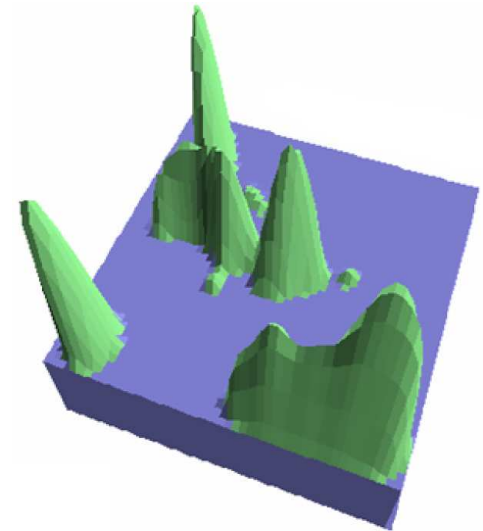
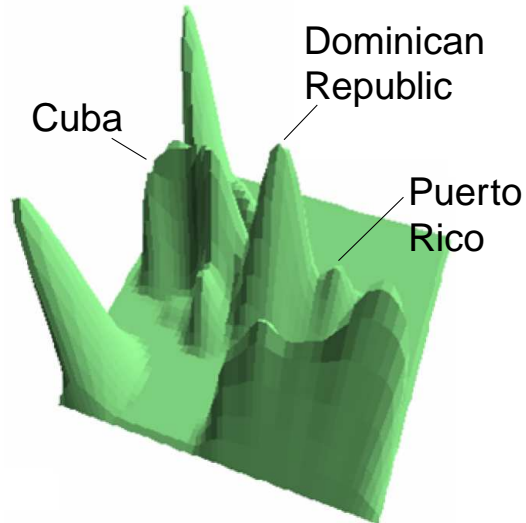
Density-based Algorithms

Analysis

Summary

Density-based Cluster Analysis

Density estimation with Gaussian Kernel for the example.



Introduction

Cluster Algorithms

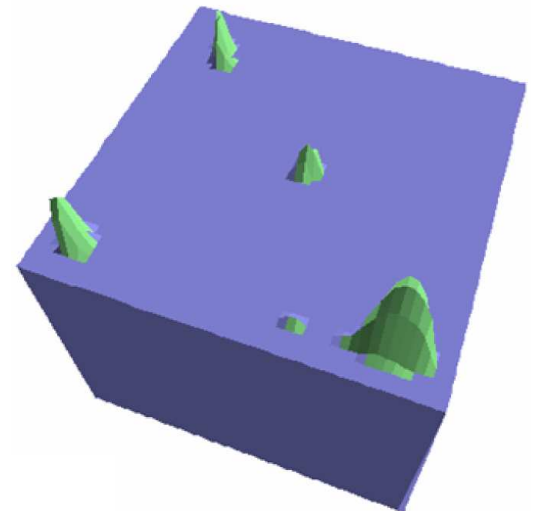
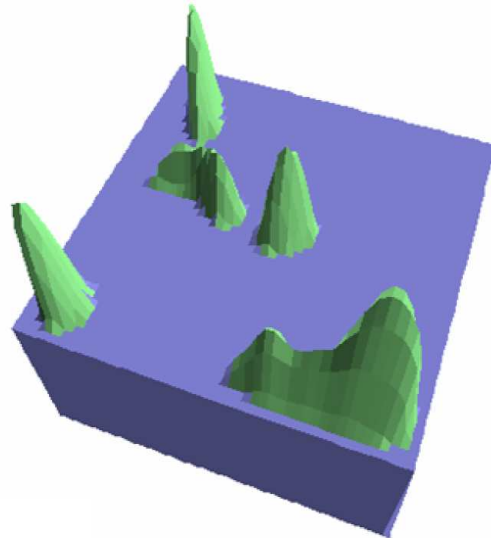
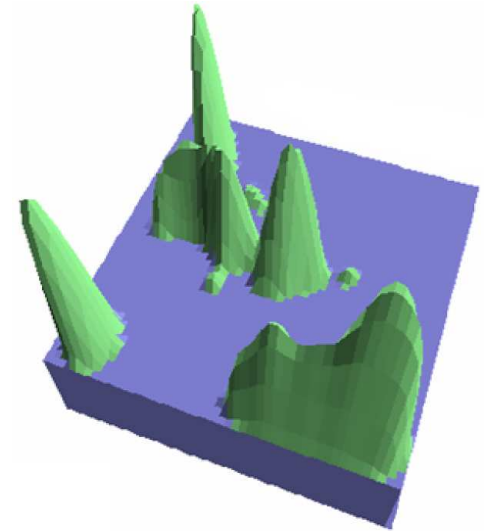
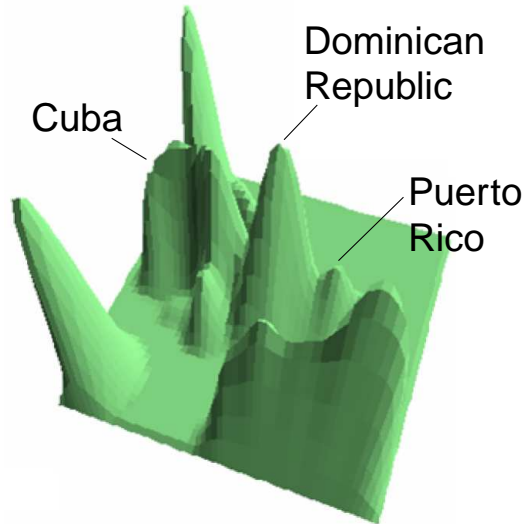
Density-based Algorithms

Analysis

Summary

Density-based Cluster Analysis

Density estimation with Gaussian Kernel for the example.



Introduction

Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Density-based Algorithm: DBSCAN

[Ester et al. 1996]

Introduction

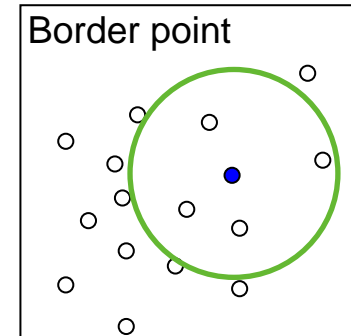
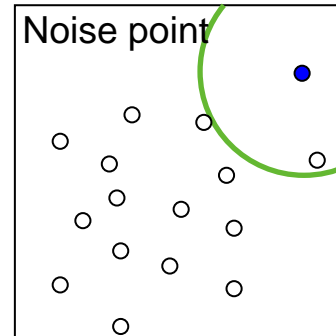
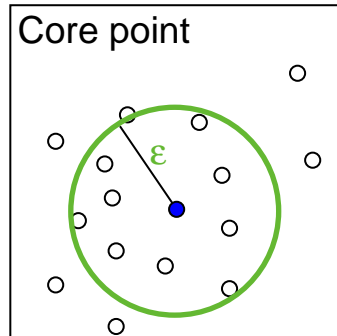
Cluster
Algorithms

**Density-based
Algorithms**

Analysis

Summary

Density-based Algorithm: DBSCAN [Ester et al. 1996]



p is core point:

$$|N_{\varepsilon}(p)| \geq \text{MinPts}.$$

p is noise point:

p is not **density-reachable** from a core point.

p is border point:

otherwise.

Introduction

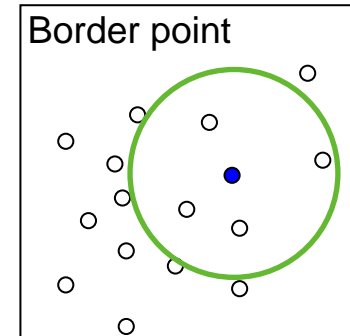
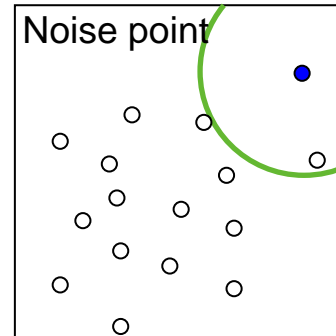
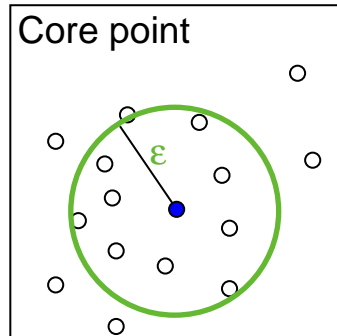
Cluster
Algorithms

Density-based
Algorithms

Analysis

Summary

Density-based Algorithm: DBSCAN



p is core point: $|N_\varepsilon(p)| \geq \text{MinPts}$.

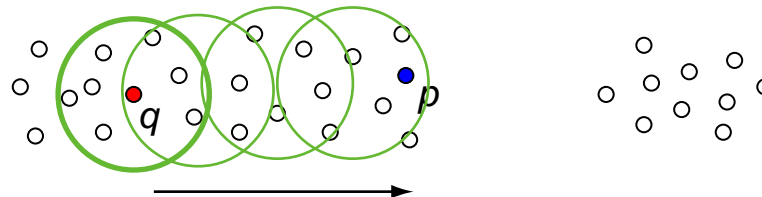
p is noise point: p is not **density-reachable** from a core point.

p is border point: otherwise.

p is density-reachable from q :

(a) $p \in |N_\varepsilon(q)|$, where q is a corepoint

(b) transitive application of condition (a):



Introduction

Cluster Algorithms

Density-based Algorithms

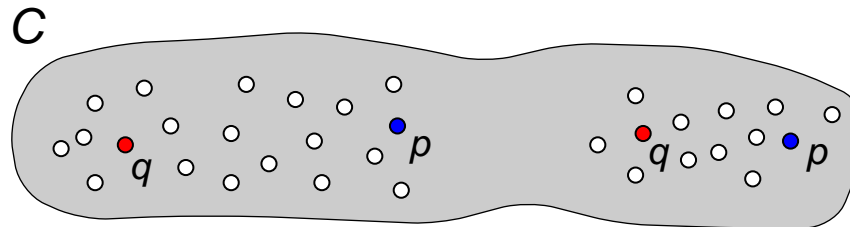
Analysis

Summary

Density-based Algorithm: DBSCAN

A cluster $C \subseteq D$ satisfies the following conditions:

1. $\forall p, q : \text{If } p \in C \text{ and } q \text{ is density-reachable from } p \text{ then } q \in C.$



Maximality
condition

Introduction

Cluster
Algorithms

Density-based
Algorithms

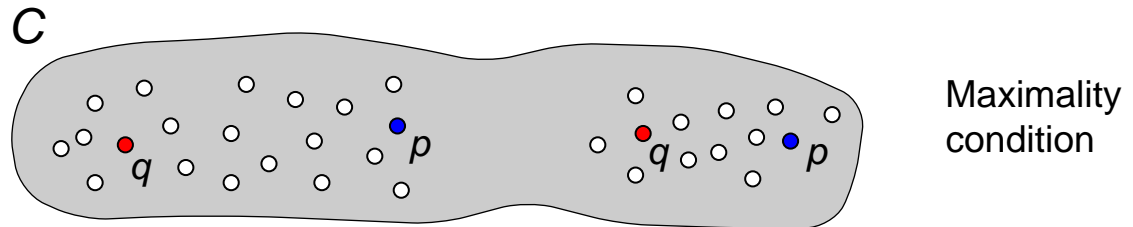
Analysis

Summary

Density-based Algorithm: DBSCAN

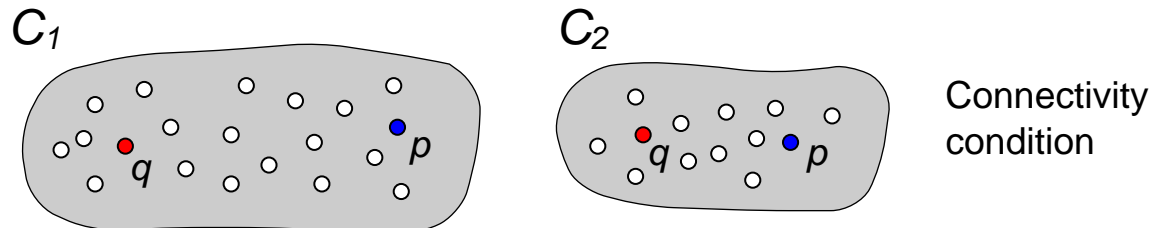
A cluster $C \subseteq D$ satisfies the following conditions:

1. $\forall p, q$: If $p \in C$ and q is density-reachable from p then $q \in C$.



2. $\forall p, q$: p is **density-connected** to q .

There is a point o such that both, p and q are density-reachable from o .



Introduction

Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Density-based Algorithm: DBSCAN

Overall cluster procedure:

1. Select unclassified point $p \in D$.
2. Construct ε -neighborhood $N_\varepsilon(p)$.
3. **If** p is a core point
Then Insert $N_\varepsilon(p)$ into new cluster C .
Recursively analyze the ε -neighborhoods of $q \in N_\varepsilon(p)$
and insert all density-reachable points into C .
Else Classify p as noise.

Introduction

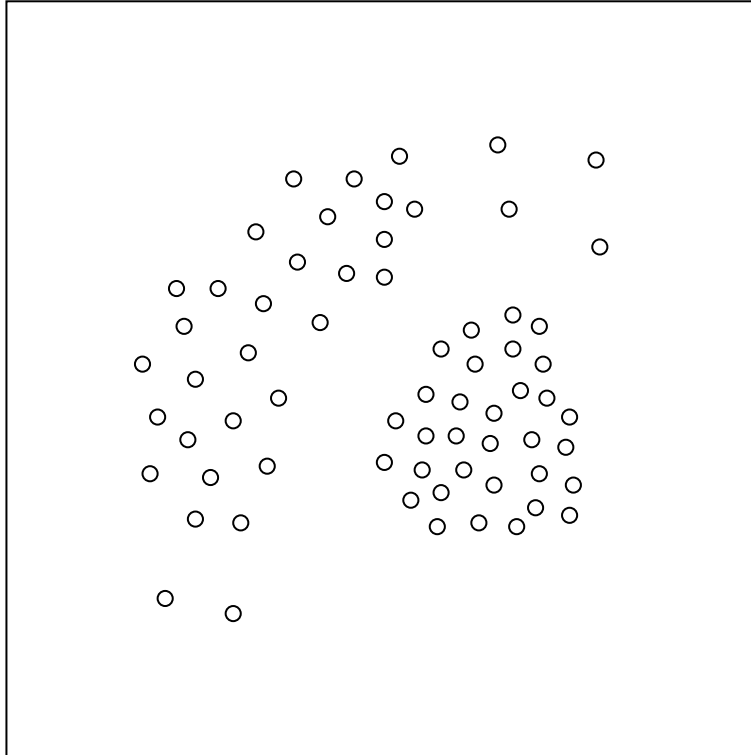
Cluster
Algorithms

Density-based
Algorithms

Analysis

Summary

Density-based Algorithm: DBSCAN



Introduction

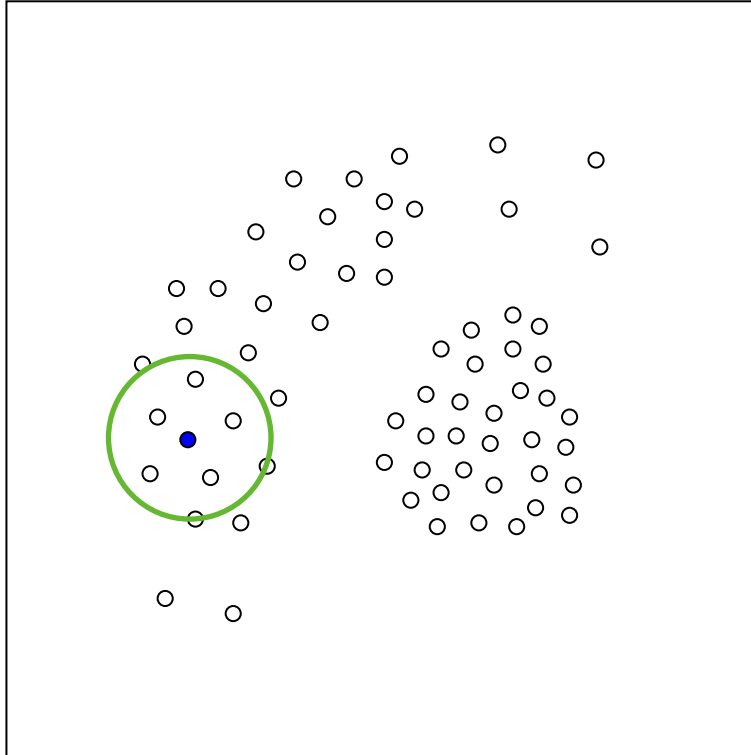
Cluster
Algorithms

**Density-based
Algorithms**

Analysis

Summary

Density-based Algorithm: DBSCAN



Introduction

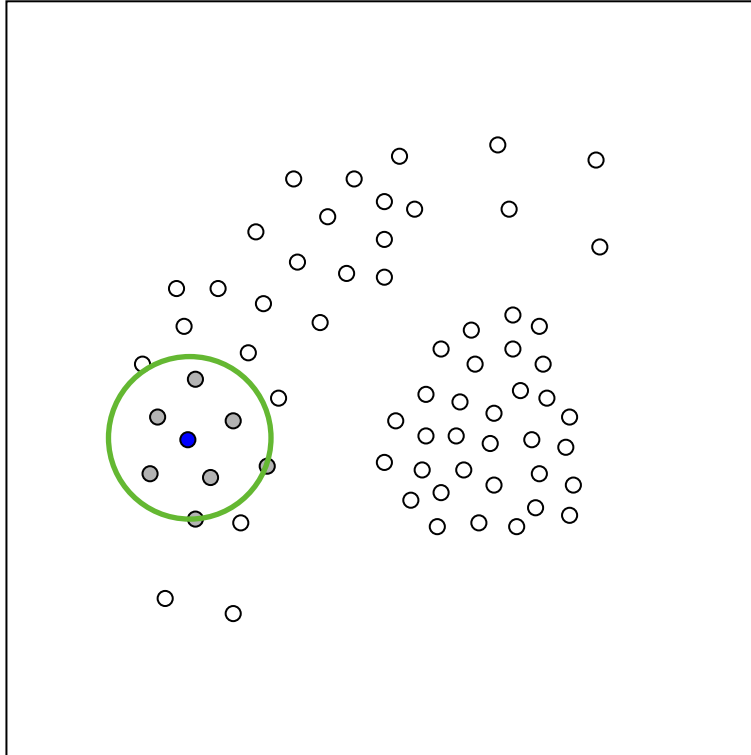
Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Density-based Algorithm: DBSCAN



Introduction

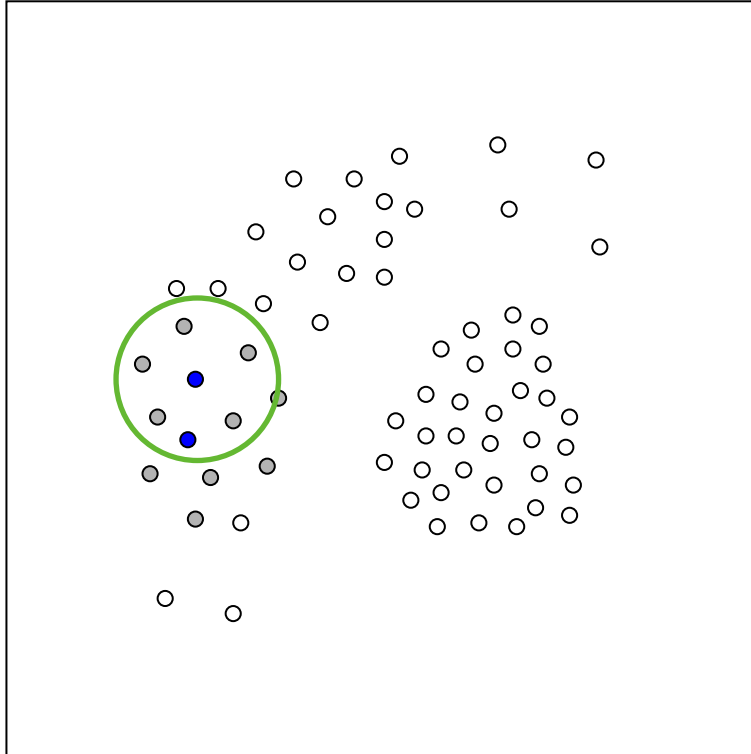
Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Density-based Algorithm: DBSCAN



Introduction

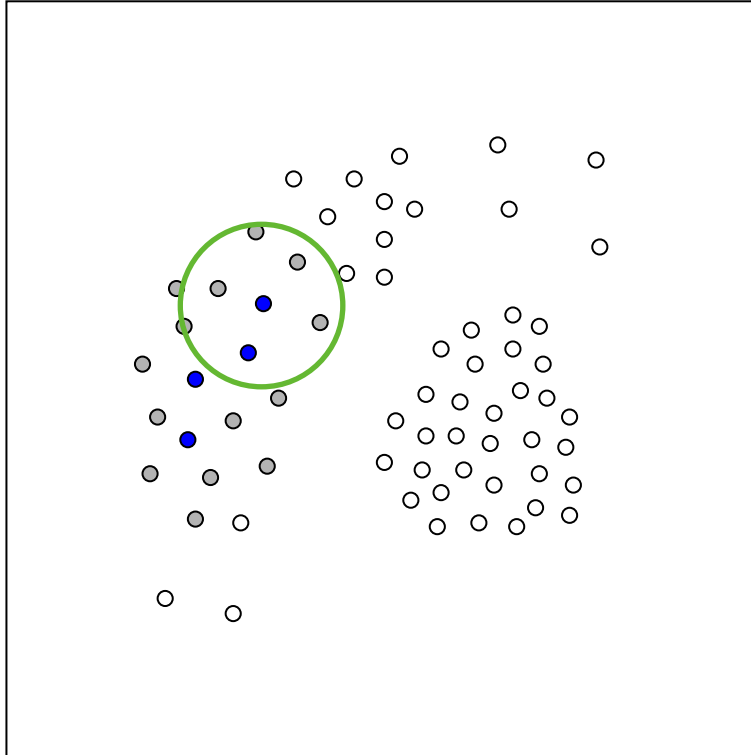
Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Density-based Algorithm: DBSCAN



Introduction

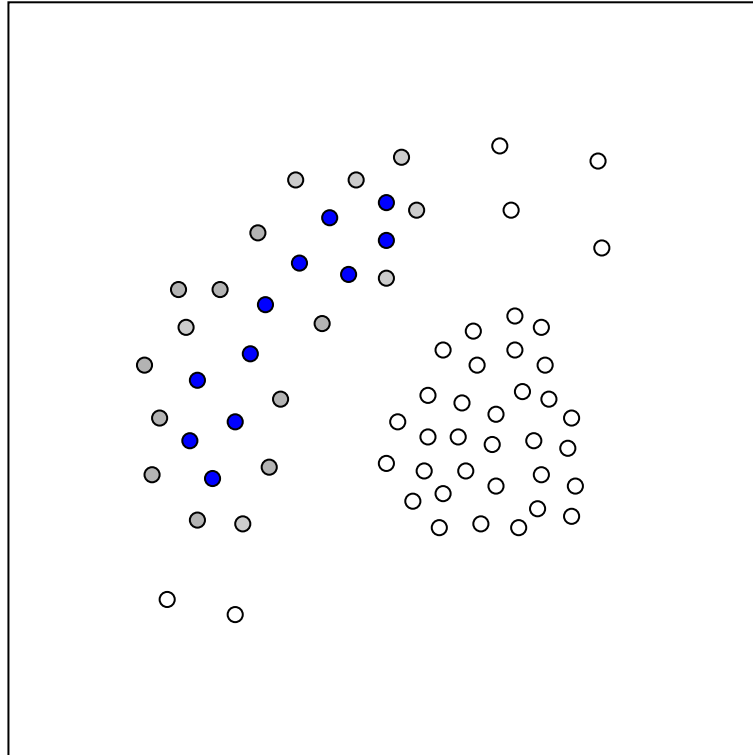
Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Density-based Algorithm: DBSCAN



- Core point
- Border point

Introduction

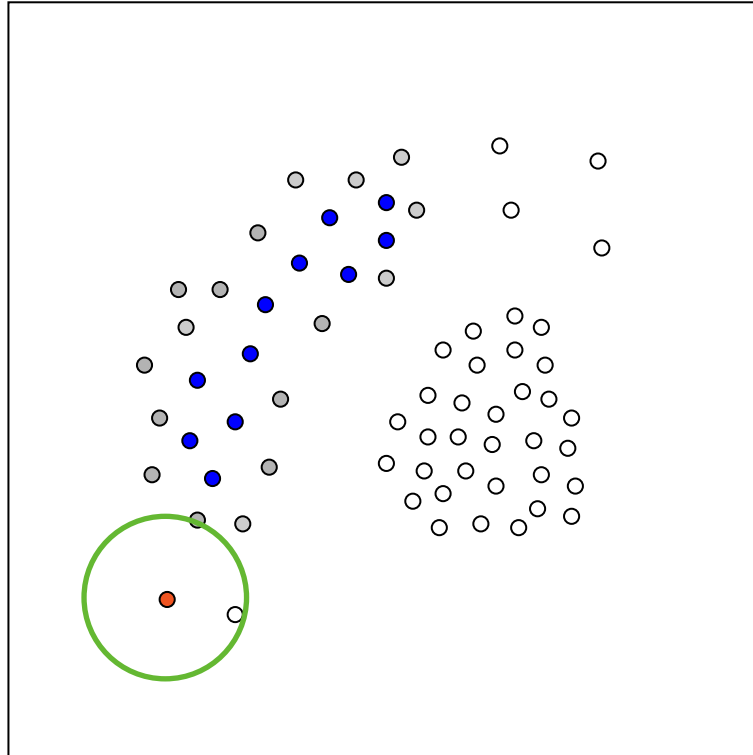
Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Density-based Algorithm: DBSCAN



- Core point
- Border point

Introduction

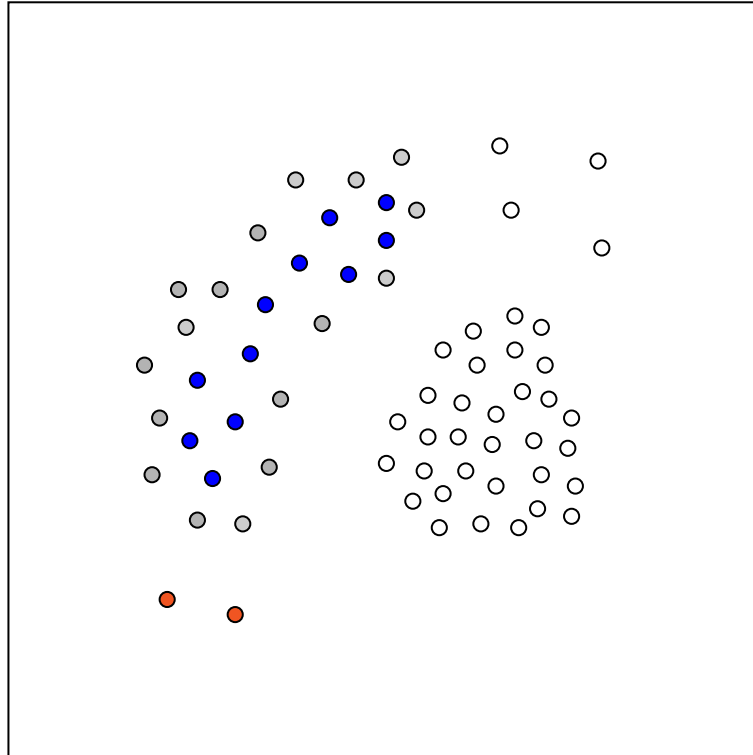
Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Density-based Algorithm: DBSCAN



- Core point
- Border point
- Noise point

Introduction

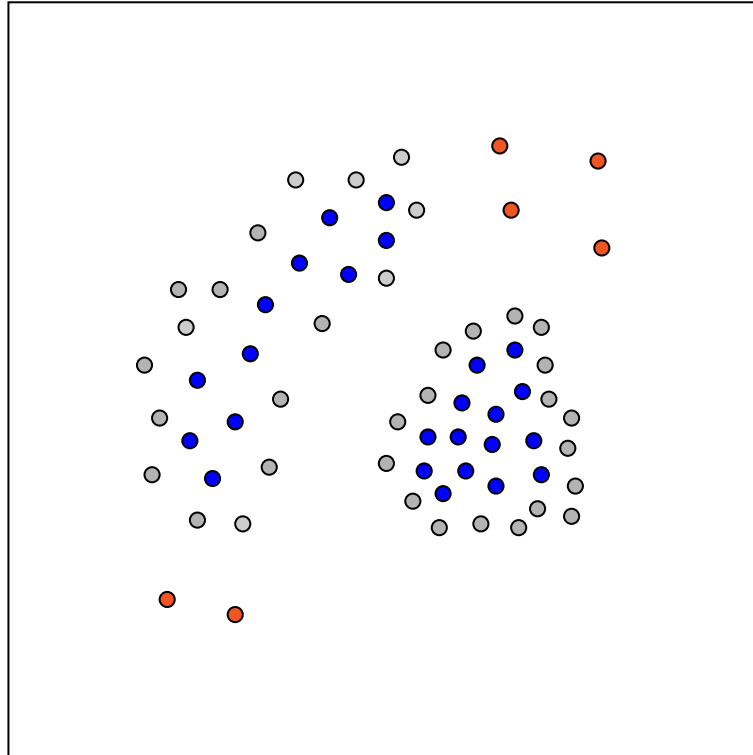
Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Density-based Algorithm: DBSCAN



- Core point
- Border point
- Noise point

Introduction

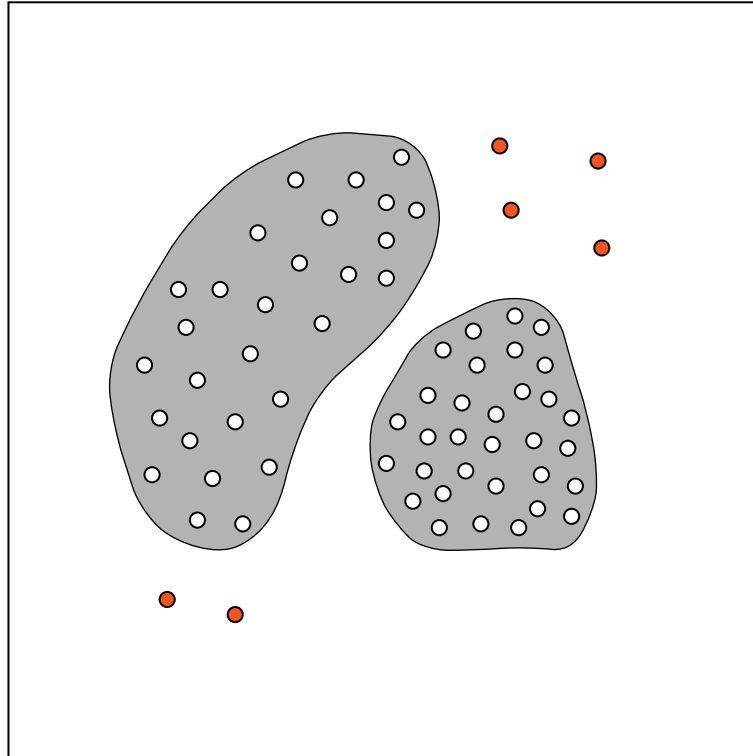
Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Density-based Algorithm: DBSCAN



● Noise point

Introduction

Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Density-based Algorithm: MajorClust

Introduction

Cluster
Algorithms

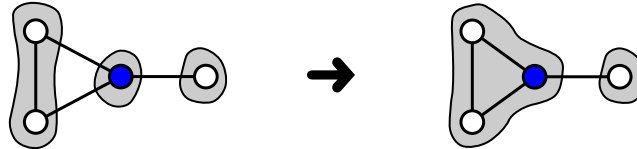
**Density-based
Algorithms**

Analysis

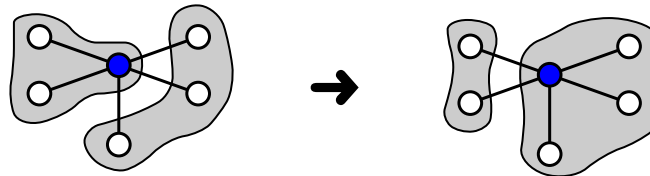
Summary

Density-based Algorithm: MajorClust

- Definite majority decision (agglomeration):



- Definite majority decision (node changes cluster):



Introduction

Cluster
Algorithms

Density-based
Algorithms

Analysis

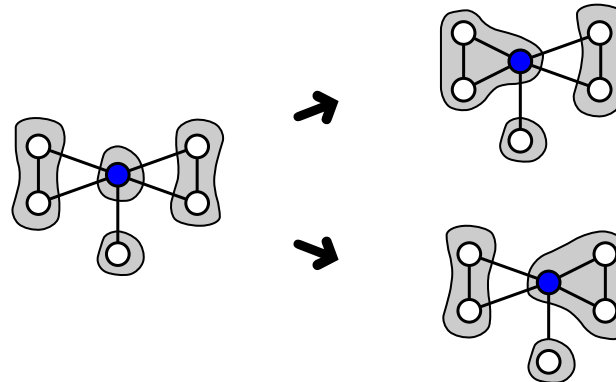
Summary

Density-based Algorithm: MajorClust

- Definite majority decision (agglomeration):



- Indefinite majority decision:



Introduction

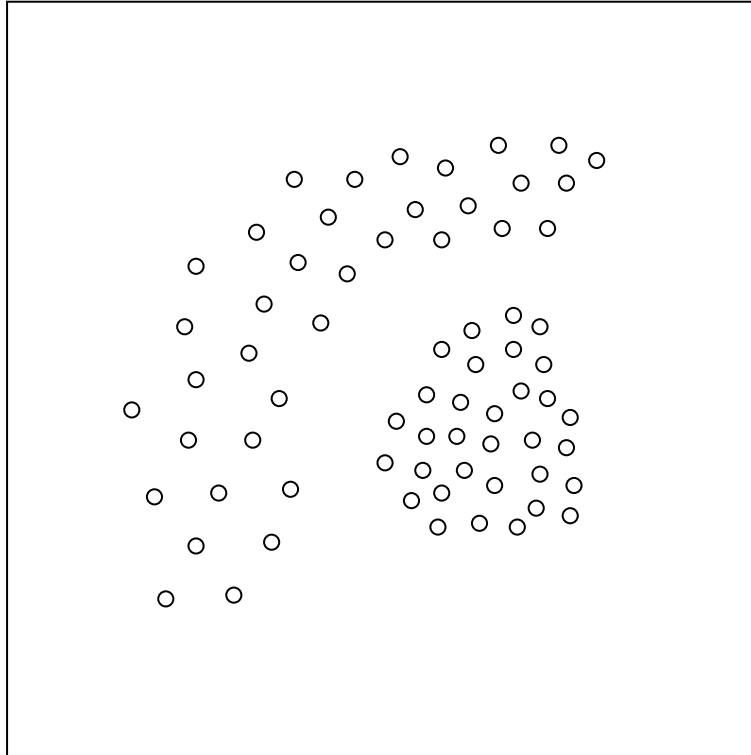
Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Density-based Algorithm: MajorClust



Introduction

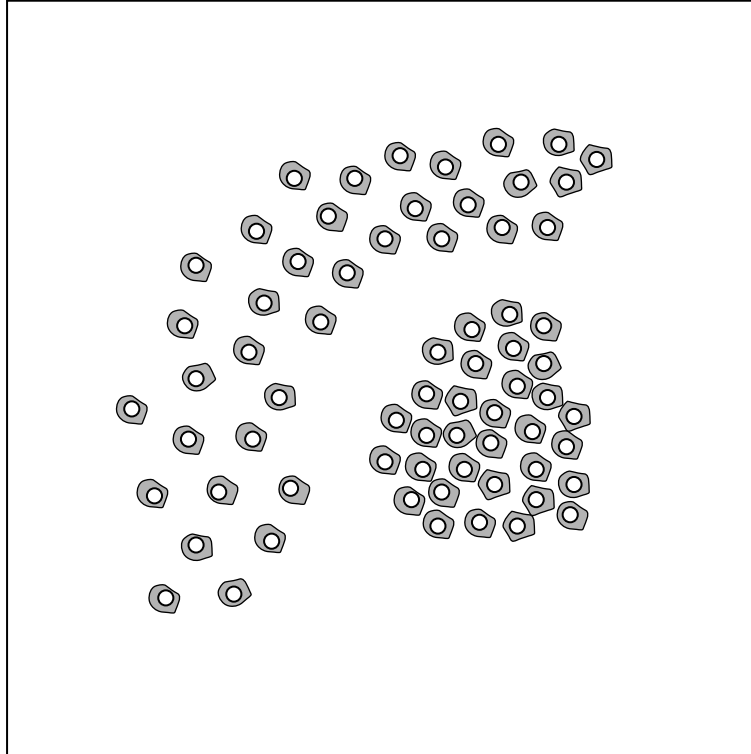
Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Density-based Algorithm: MajorClust



Introduction

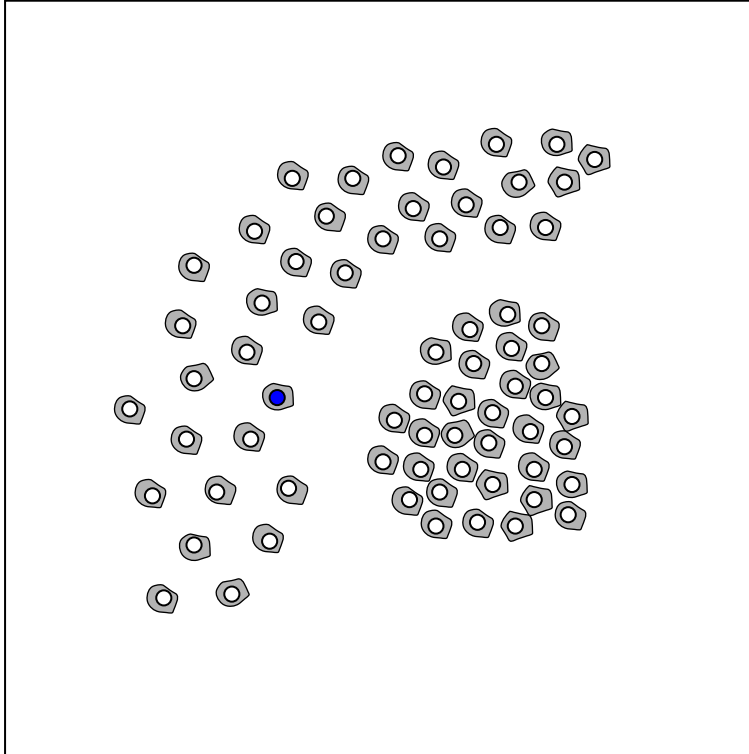
Cluster
Algorithms

**Density-based
Algorithms**

Analysis

Summary

Density-based Algorithm: MajorClust



Introduction

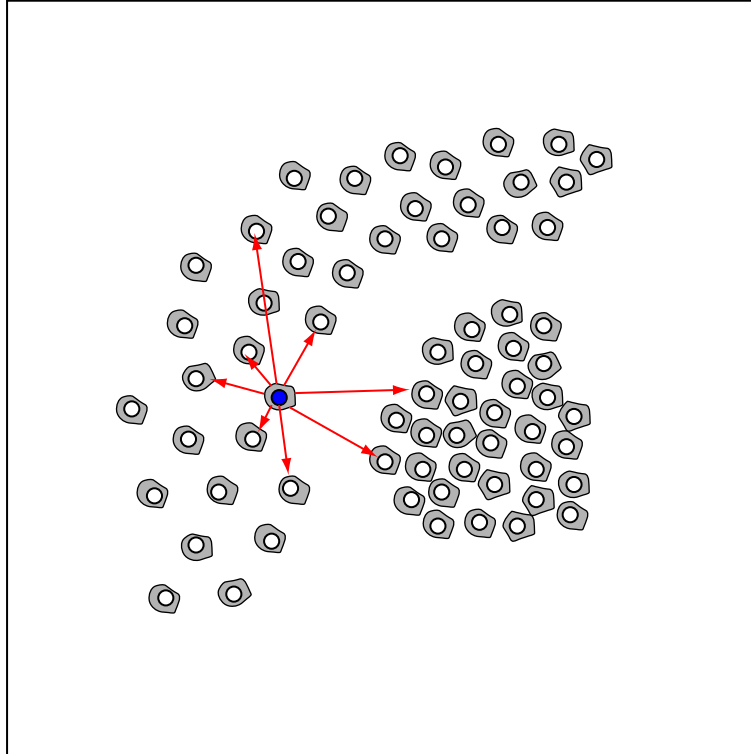
Cluster
Algorithms

**Density-based
Algorithms**

Analysis

Summary

Density-based Algorithm: MajorClust



Introduction

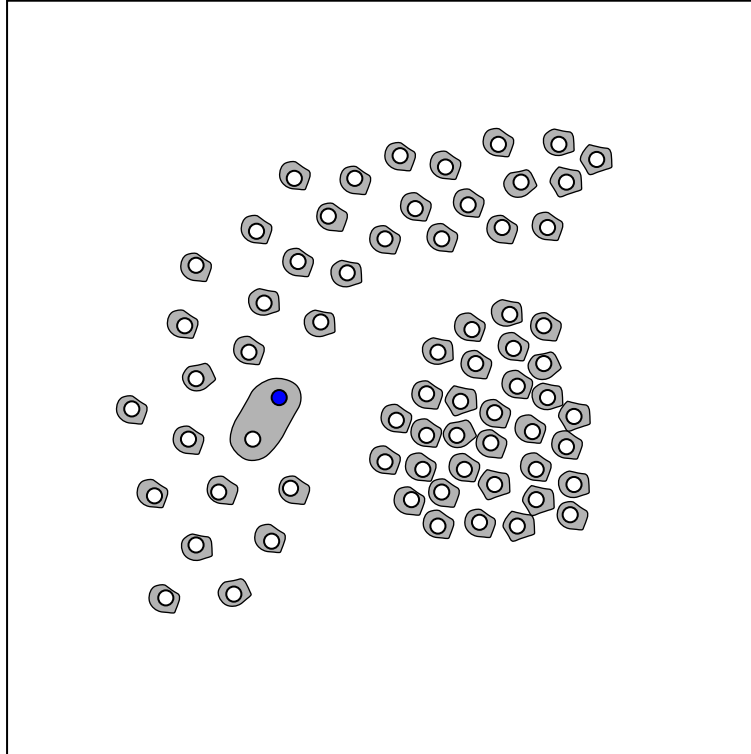
Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Density-based Algorithm: MajorClust



Introduction

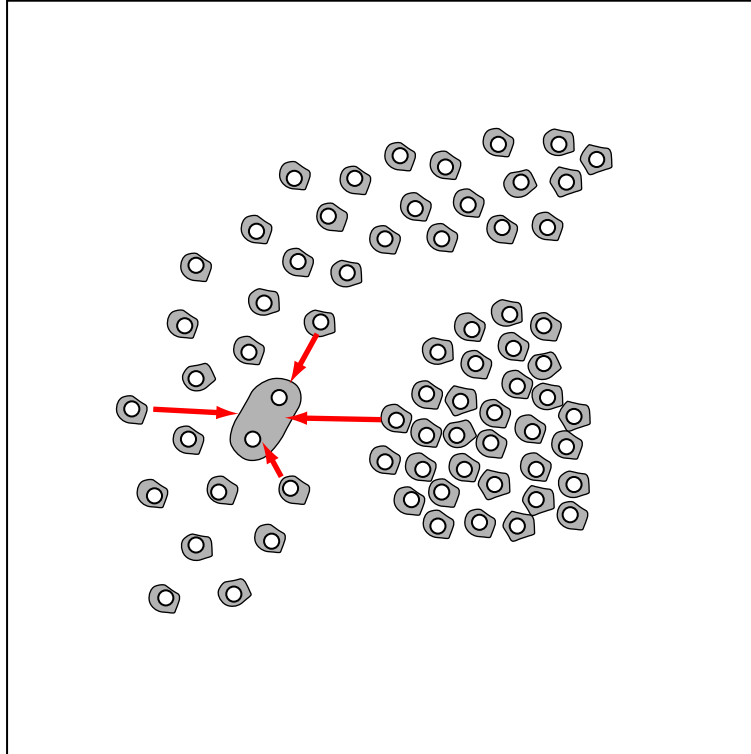
Cluster
Algorithms

**Density-based
Algorithms**

Analysis

Summary

Density-based Algorithm: MajorClust



Introduction

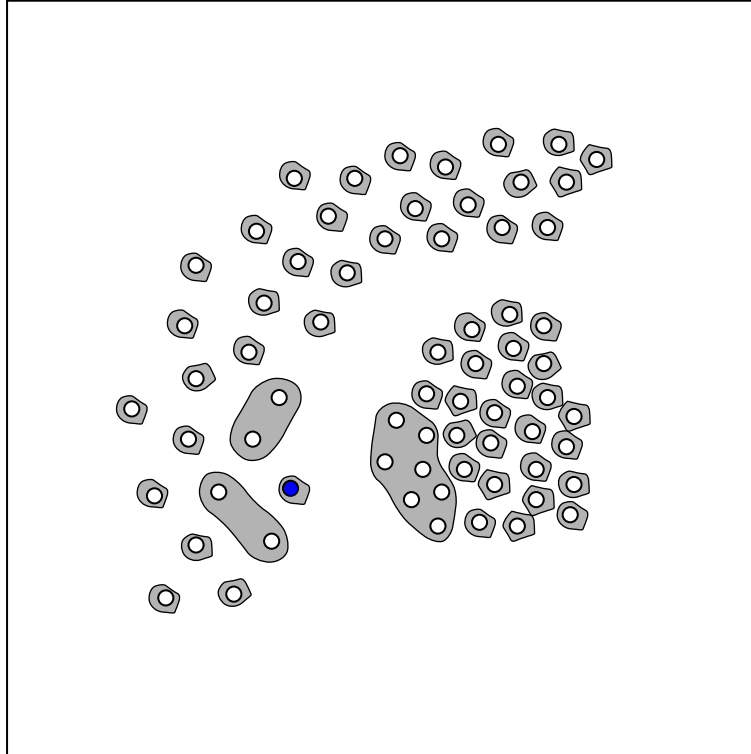
Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Density-based Algorithm: MajorClust



Introduction

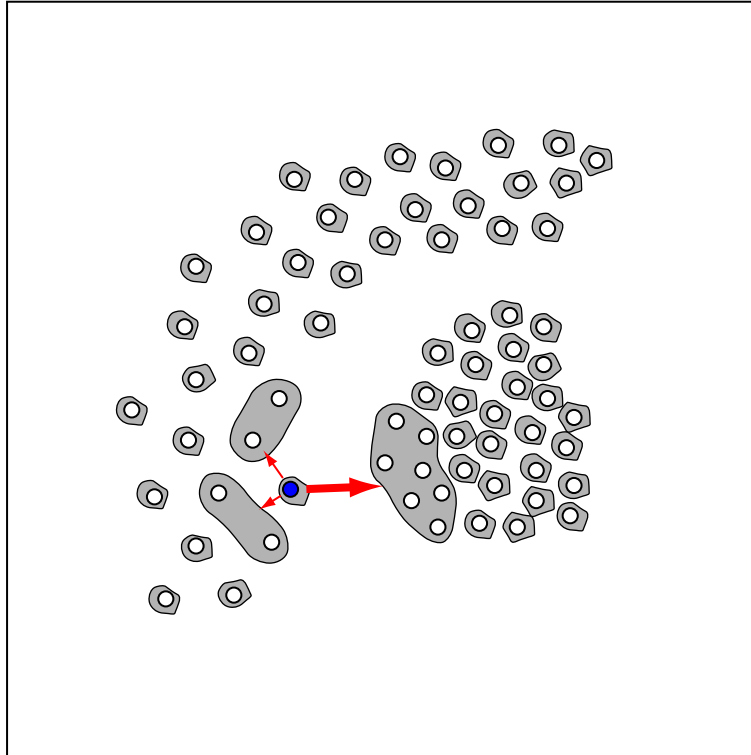
Cluster
Algorithms

**Density-based
Algorithms**

Analysis

Summary

Density-based Algorithm: MajorClust



Introduction

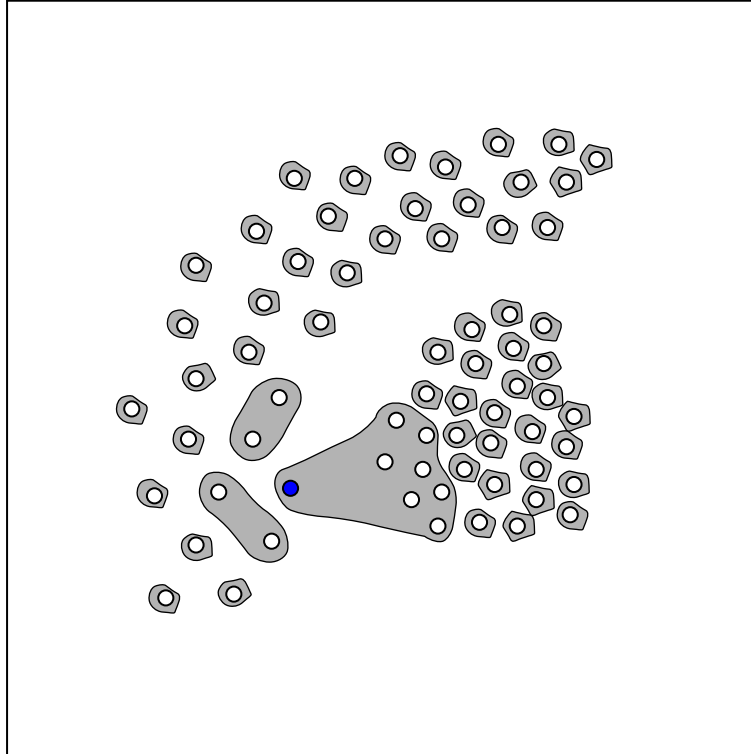
Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Density-based Algorithm: MajorClust



Introduction

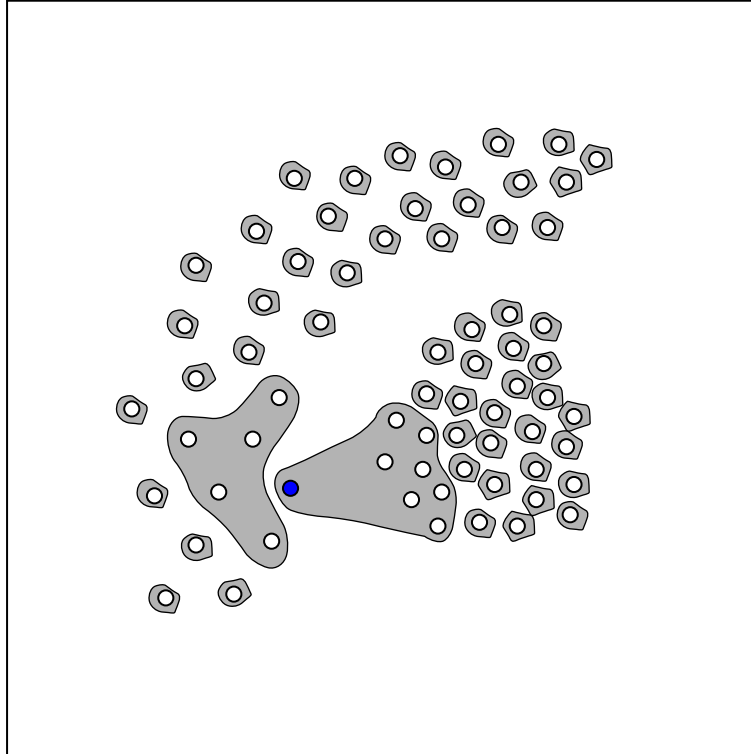
Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Density-based Algorithm: MajorClust



Introduction

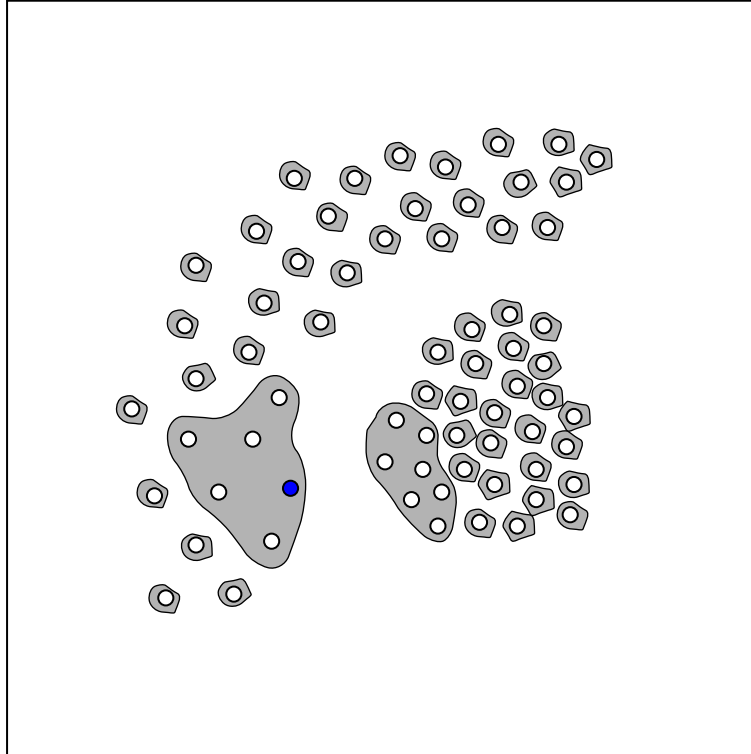
Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Density-based Algorithm: MajorClust



Introduction

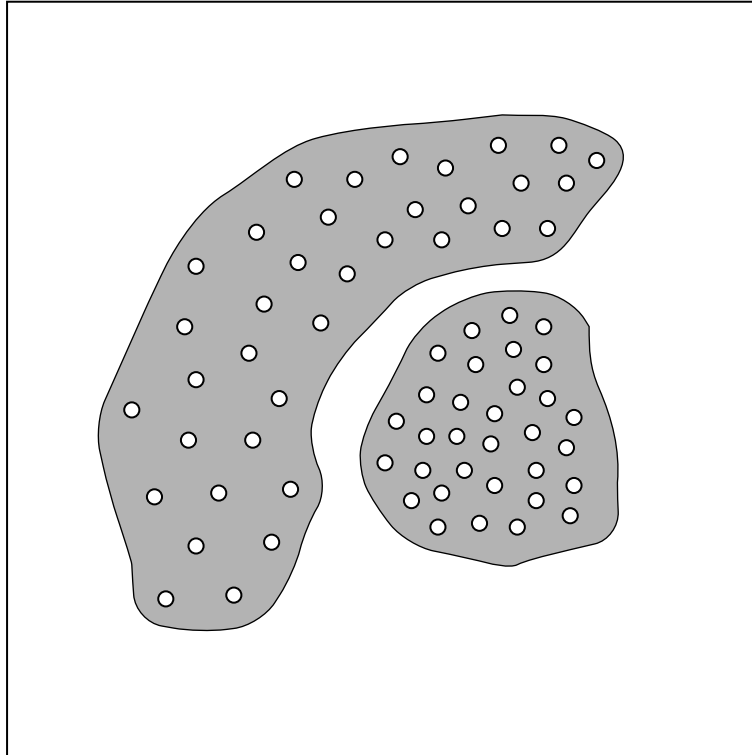
Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Density-based Algorithm: MajorClust



Introduction

Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Analysis I (low-dimensional)

Geometrical Data—map of the Caribbean Islands (approx. 20,000 points) :

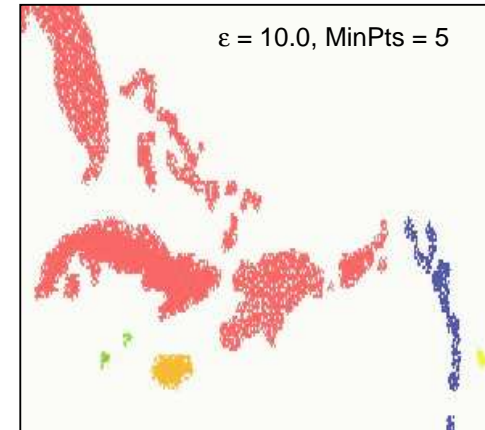
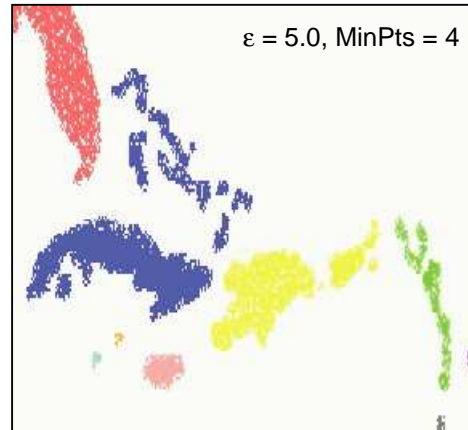
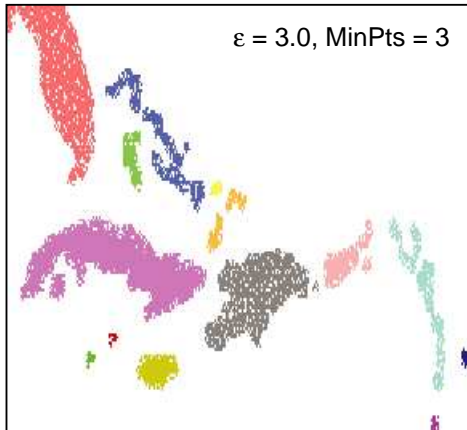


Analysis I (low-dimensional)

Geometrical Data—map of the Caribbean Islands (approx. 20,000 points) :

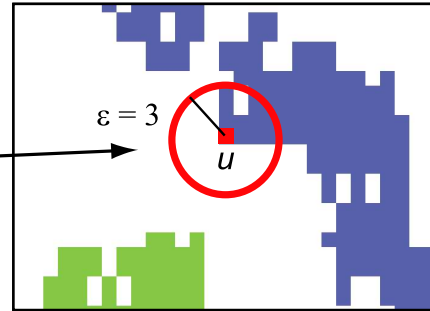
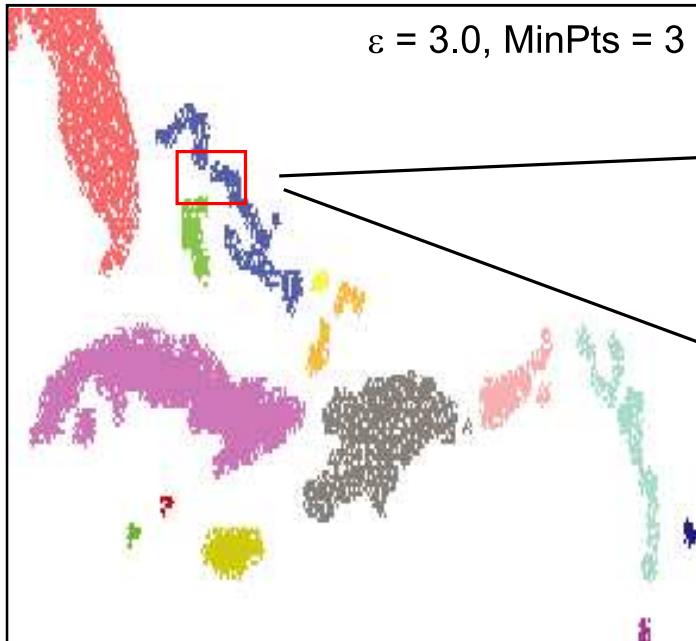


DBSCAN:

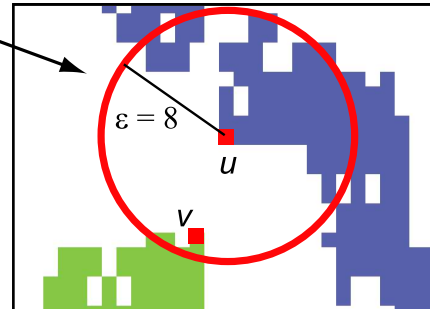


Analysis I (low-dimensional)

The problem of choosing a good ϵ -value in DBSCAN.



Two separate clusters are found.



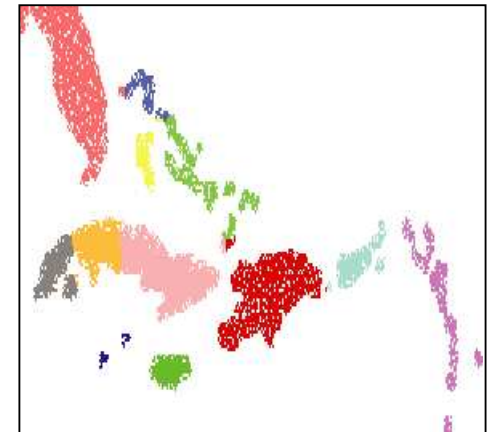
Clusters are merged.

Analysis I (low-dimensional)

Geometrical Data—map of the Caribbean Islands (approx. 20,000 points) :

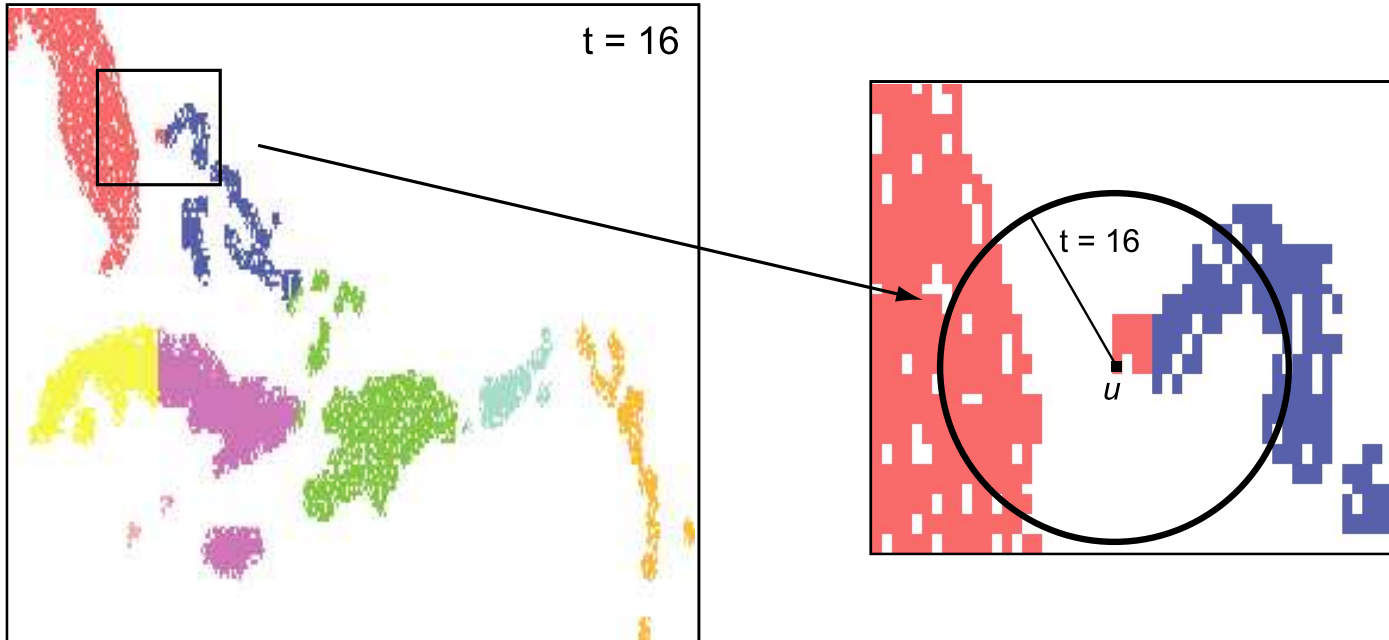


MajorClust:



Analysis I (low-dimensional)

The problem of a global analysis (no ε -neighborhood restriction) in MajorClust.



Analysis II (high-dimensional)

Document categorization with the Reuters corpus.



- ❑ 1000 documents
- ❑ 10 categories: politics, culture, economics, etc.
- ❑ uniformly distributed, exclusive membership
- ❑ **> 10,000 dimensions**

Introduction

Cluster Algorithms

Density-based Algorithms

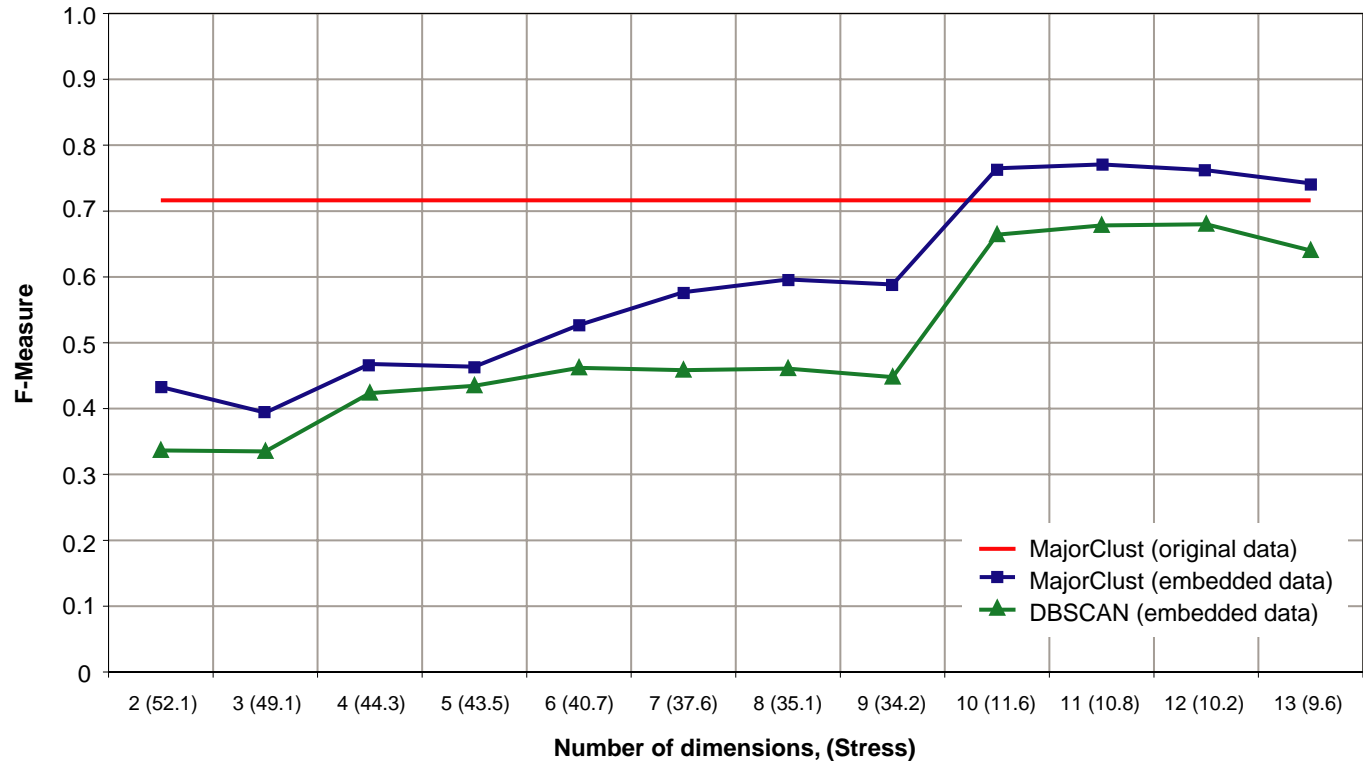
Analysis

Summary

Analysis II (high-dimensional)

DBSCAN requires embedding of data in low-dimensional space.

Classification results (F -Measure) over dimensionality:



Introduction

Cluster Algorithms

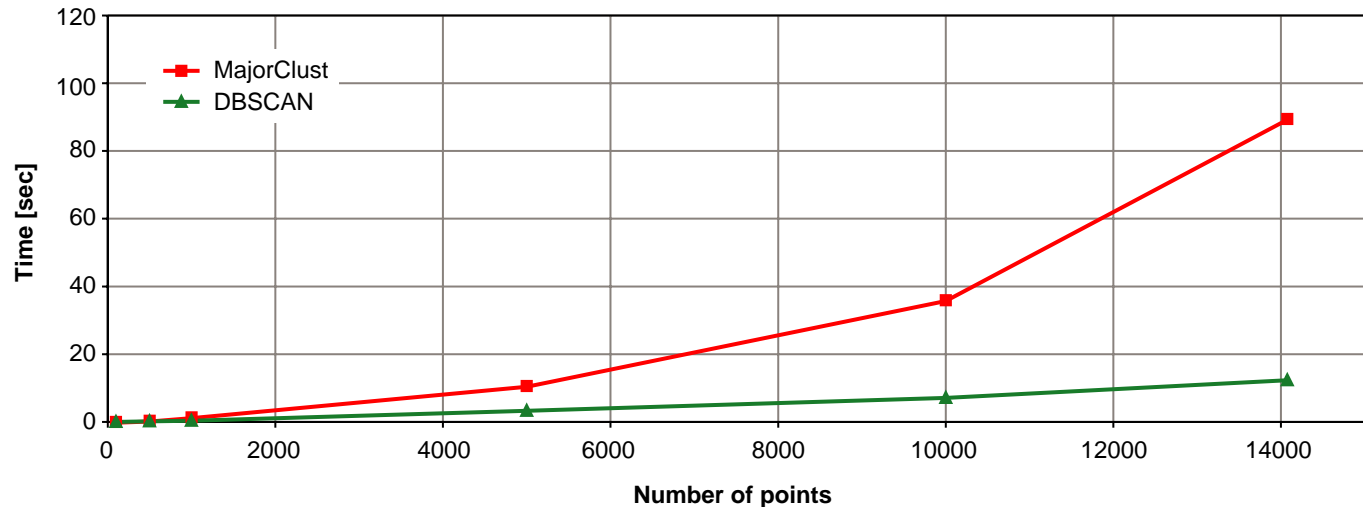
Density-based Algorithms

Analysis

Summary

Analysis (runtime)

Runtime-behavior on the geometrical data:



Note:

The embedding of data in a low-dimensional space (MDS) is computationally very expensive:

I. e., most cluster algorithms will be faster than DBSCAN + MDS.

Introduction

Cluster Algorithms

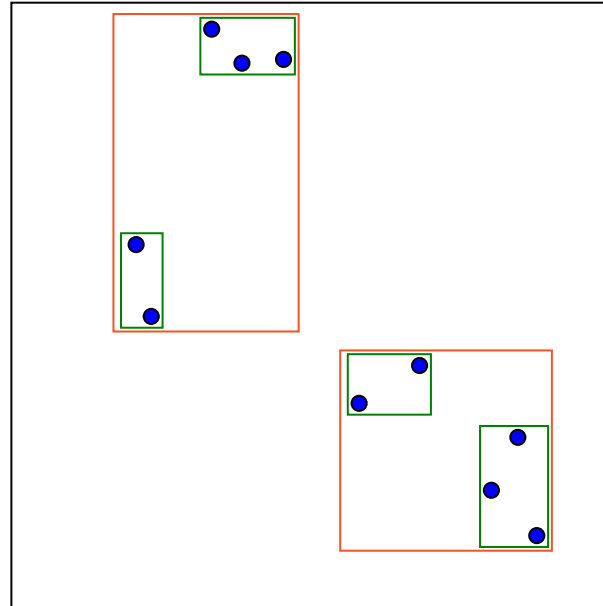
Density-based Algorithms

Analysis

Summary

Analysis (runtime)

DBSCAN employs the R -tree data structure for region queries, which constructs minimum bounding regions for inserted points:



“Existing methods are outperformed on on average by a simple sequential scan, if the number of dimensions exceeds around 10.”

[Weber 99, Gionis/Indyk/Motwani 99-04]

Introduction

Cluster
Algorithms

Density-based
Algorithms

Analysis

Summary

Summary

An alternative categorization scheme :

	Cluster approach			
	hierarchical	iterative	density-based	meta-search controlled
Analysis strategy	relative comparison based on two items	absolute comparison based on k items	relative comparison based on k items	absolute comparison based on all items
Recovery characteristics	irrevocable	revocable	revocable	revocable

Orthogonal to this scheme is the concept for similarity computation:

- distance (neighborhood) analysis in low-dimensional space
- similarity predicate in arbitrary (high-dimensional) space

Introduction

Cluster Algorithms

Density-based Algorithms

Analysis

Summary

Summary

The strengths and weaknesses of density-based cluster algorithms can be explained with the dimensionality of the data.

- ❑ DBSCAN usually outperforms other cluster algorithms on low-dimensional data.
- ❑ MajorClust usually outperforms other cluster algorithms on high-dimensional data, in particular in the document categorization field.

Current work:

How fingerprints can be utilized for efficient region queries in high-dimensional spaces.

Introduction

Cluster
Algorithms

Density-based
Algorithms

Analysis

Summary