

Clustering and mapper

Andrew J. Blumberg (blumberg@math.utexas.edu)

June 17th, 2014

Overview

Goal of talk

Explain Mapper, which is the most widely used and most successful TDA technique. (At core of Ayasdi, TDA company founded by Gunnar Carlsson.)

Basic idea: perform clustering at different “scales”, track how clusters change as scale varies.

Motivation:

- 1 Coarser than manifold learning, but still works in nonlinear situations.
- 2 Still retains meaningful geometric information about data set.
- 3 Efficiently computable (and so can apply to very large data sets).

Overview

Goal of talk

Explain Mapper, which is the most widely used and most successful TDA technique. (At core of Ayasdi, TDA company founded by Gunnar Carlsson.)

Basic idea: perform clustering at different “scales”, track how clusters change as scale varies.

Motivation:

- 1 Coarser than manifold learning, but still works in nonlinear situations.
- 2 Still retains meaningful geometric information about data set.
- 3 Efficiently computable (and so can apply to very large data sets).

Goal of talk

Explain Mapper, which is the most widely used and most successful TDA technique. (At core of Ayasdi, TDA company founded by Gunnar Carlsson.)

Basic idea: perform clustering at different “scales”, track how clusters change as scale varies.

Motivation:

- 1 Coarser than manifold learning, but still works in nonlinear situations.
- 2 Still retains meaningful geometric information about data set.
- 3 Efficiently computable (and so can apply to very large data sets).

Goal of talk

Explain Mapper, which is the most widely used and most successful TDA technique. (At core of Ayasdi, TDA company founded by Gunnar Carlsson.)

Basic idea: perform clustering at different “scales”, track how clusters change as scale varies.

Motivation:

- 1 Coarser than manifold learning, but still works in nonlinear situations.
- 2 Still retains meaningful geometric information about data set.
- 3 Efficiently computable (and so can apply to very large data sets).

Goal of talk

Explain Mapper, which is the most widely used and most successful TDA technique. (At core of Ayasdi, TDA company founded by Gunnar Carlsson.)

Basic idea: perform clustering at different “scales”, track how clusters change as scale varies.

Motivation:

- 1 Coarser than manifold learning, but still works in nonlinear situations.
- 2 Still retains meaningful geometric information about data set.
- 3 Efficiently computable (and so can apply to very large data sets).

Goal of talk

Explain Mapper, which is the most widely used and most successful TDA technique. (At core of Ayasdi, TDA company founded by Gunnar Carlsson.)

Basic idea: perform clustering at different “scales”, track how clusters change as scale varies.

Motivation:

- 1 Coarser than manifold learning, but still works in nonlinear situations.
- 2 Still retains meaningful geometric information about data set.
- 3 Efficiently computable (and so can apply to very large data sets).

Basic idea

Describe topology of a smooth manifold M using levelsets of a suitable function $h: M \rightarrow \mathbb{R}$.

- We recover M by looking at $h^{-1}((-\infty, t])$, as t scans over the range of h .
- Topology of M changes at critical points of h .

Basic idea

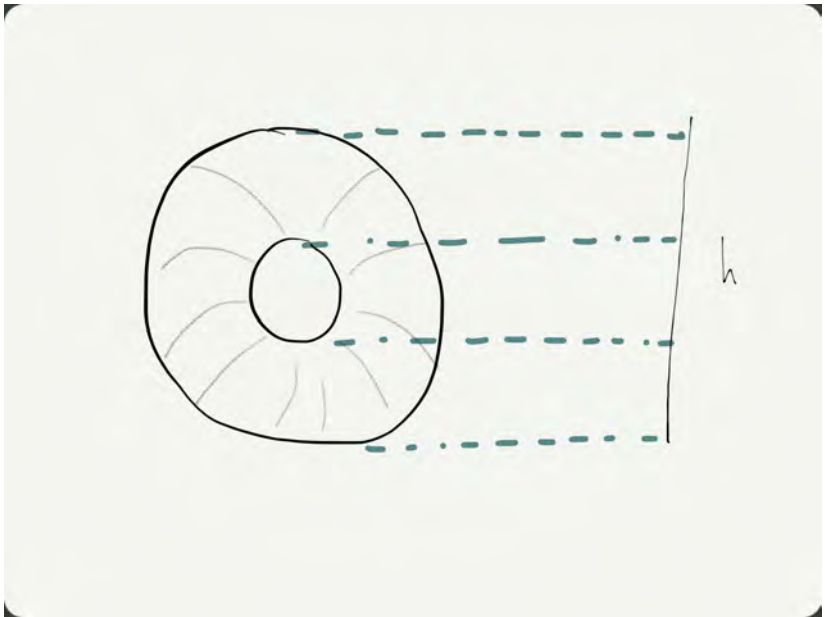
Describe topology of a smooth manifold M using levelsets of a suitable function $h: M \rightarrow \mathbb{R}$.

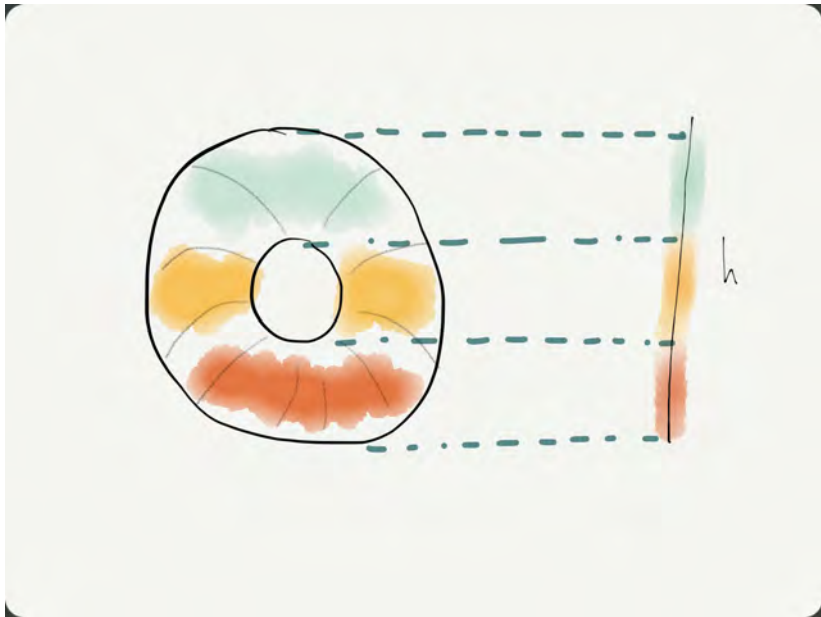
- We recover M by looking at $h^{-1}((-\infty, t])$, as t scans over the range of h .
- Topology of M changes at critical points of h .

Basic idea

Describe topology of a smooth manifold M using levelsets of a suitable function $h: M \rightarrow \mathbb{R}$.

- We recover M by looking at $h^{-1}((-\infty, t])$, as t scans over the range of h .
- Topology of M changes at critical points of h .





Convenient simplification:

- 1 For each $t \in \mathbb{R}$, contract each component of $f^{-1}(t)$ to a point.
- 2 Resulting structure is a graph.

Convenient simplification:

- 1 For each $t \in \mathbb{R}$, contract each component of $f^{-1}(t)$ to a point.
- 2 Resulting structure is a graph.

Convenient simplification:

- 1 For each $t \in \mathbb{R}$, contract each component of $f^{-1}(t)$ to a point.
- 2 Resulting structure is a graph.



The mapper algorithm is a generalization of this procedure.
[Singh-Memoli-Carlsson]

Assume given a data set X .

- ① Choose a filter function $f: X \rightarrow \mathbb{R}$.
- ② Choose a cover U_α of X .
- ③ Cluster each inverse image $f^{-1}(U_\alpha)$.
- ④ Form a graph where:
 - ① Clusters are vertices.
 - ② An edge connects two clusters C and C' if both $U_\alpha \cap U_{\alpha'} \neq \emptyset$ and $C \cap C' \neq \emptyset$.
- ⑤ Color vertices according to average value of f in the cluster.

The mapper algorithm is a generalization of this procedure.
[Singh-Memoli-Carlsson]

Assume given a data set X .

- ① Choose a filter function $f: X \rightarrow \mathbb{R}$.
- ② Choose a cover U_α of X .
- ③ Cluster each inverse image $f^{-1}(U_\alpha)$.
- ④ Form a graph where:
 - ① Clusters are vertices.
 - ② An edge connects two clusters C and C' if both $U_\alpha \cap U_{\alpha'} \neq \emptyset$ and $C \cap C' \neq \emptyset$.
- ⑤ Color vertices according to average value of f in the cluster.

The mapper algorithm is a generalization of this procedure.
[Singh-Memoli-Carlsson]

Assume given a data set X .

- ① Choose a filter function $f: X \rightarrow \mathbb{R}$.
- ② Choose a cover U_α of X .
- ③ Cluster each inverse image $f^{-1}(U_\alpha)$.
- ④ Form a graph where:
 - ① Clusters are vertices.
 - ② An edge connects two clusters C and C' if both $U_\alpha \cap U_{\alpha'} \neq \emptyset$ and $C \cap C' \neq \emptyset$.
- ⑤ Color vertices according to average value of f in the cluster.

The mapper algorithm is a generalization of this procedure.
[Singh-Memoli-Carlsson]

Assume given a data set X .

- ① Choose a filter function $f: X \rightarrow \mathbb{R}$.
- ② Choose a cover U_α of X .
- ③ Cluster each inverse image $f^{-1}(U_\alpha)$.
- ④ Form a graph where:
 - ① Clusters are vertices.
 - ② An edge connects two clusters C and C' if both $U_\alpha \cap U_{\alpha'} \neq \emptyset$ and $C \cap C' \neq \emptyset$.
- ⑤ Color vertices according to average value of f in the cluster.

The mapper algorithm is a generalization of this procedure.
[Singh-Memoli-Carlsson]

Assume given a data set X .

- ① Choose a filter function $f: X \rightarrow \mathbb{R}$.
- ② Choose a cover U_α of X .
- ③ Cluster each inverse image $f^{-1}(U_\alpha)$.
- ④ Form a graph where:
 - ① Clusters are vertices.
 - ② An edge connects two clusters C and C' if both $U_\alpha \cap U_{\alpha'} \neq \emptyset$ and $C \cap C' \neq \emptyset$.
- ⑤ Color vertices according to average value of f in the cluster.

Mapper

The mapper algorithm is a generalization of this procedure.
[Singh-Memoli-Carlsson]

Assume given a data set X .

- ① Choose a filter function $f: X \rightarrow \mathbb{R}$.
- ② Choose a cover U_α of X .
- ③ Cluster each inverse image $f^{-1}(U_\alpha)$.
- ④ Form a graph where:
 - ① Clusters are vertices.
 - ② An edge connects two clusters C and C' if both $U_\alpha \cap U_{\alpha'} \neq \emptyset$ and $C \cap C' \neq \emptyset$.
- ⑤ Color vertices according to average value of f in the cluster.

The mapper algorithm is a generalization of this procedure.
[Singh-Memoli-Carlsson]

Assume given a data set X .

- ① Choose a filter function $f: X \rightarrow \mathbb{R}$.
- ② Choose a cover U_α of X .
- ③ Cluster each inverse image $f^{-1}(U_\alpha)$.
- ④ Form a graph where:
 - ① Clusters are vertices.
 - ② An edge connects two clusters C and C' if both $U_\alpha \cap U_{\alpha'} \neq \emptyset$ and $C \cap C' \neq \emptyset$.
- ⑤ Color vertices according to average value of f in the cluster.

Mapper

The mapper algorithm is a generalization of this procedure.
[Singh-Memoli-Carlsson]

Assume given a data set X .

- ① Choose a filter function $f: X \rightarrow \mathbb{R}$.
- ② Choose a cover U_α of X .
- ③ Cluster each inverse image $f^{-1}(U_\alpha)$.
- ④ Form a graph where:
 - ① Clusters are vertices.
 - ② An edge connects two clusters C and C' if both $U_\alpha \cap U_{\alpha'} \neq \emptyset$ and $C \cap C' \neq \emptyset$.
- ⑤ Color vertices according to average value of f in the cluster.

The mapper algorithm is a generalization of this procedure.
[Singh-Memoli-Carlsson]

Assume given a data set X .

- ① Choose a filter function $f: X \rightarrow \mathbb{R}$.
- ② Choose a cover U_α of X .
- ③ Cluster each inverse image $f^{-1}(U_\alpha)$.
- ④ Form a graph where:
 - ① Clusters are vertices.
 - ② An edge connects two clusters C and C' if both $U_\alpha \cap U_{\alpha'} \neq \emptyset$ and $C \cap C' \neq \emptyset$.
- ⑤ Color vertices according to average value of f in the cluster.

Clearly, choice of filter function is essential.

- Some kind of density measure.
- A score measure difference (distance) from some baseline.
- An eccentricity measure.

Clearly, choice of filter function is essential.

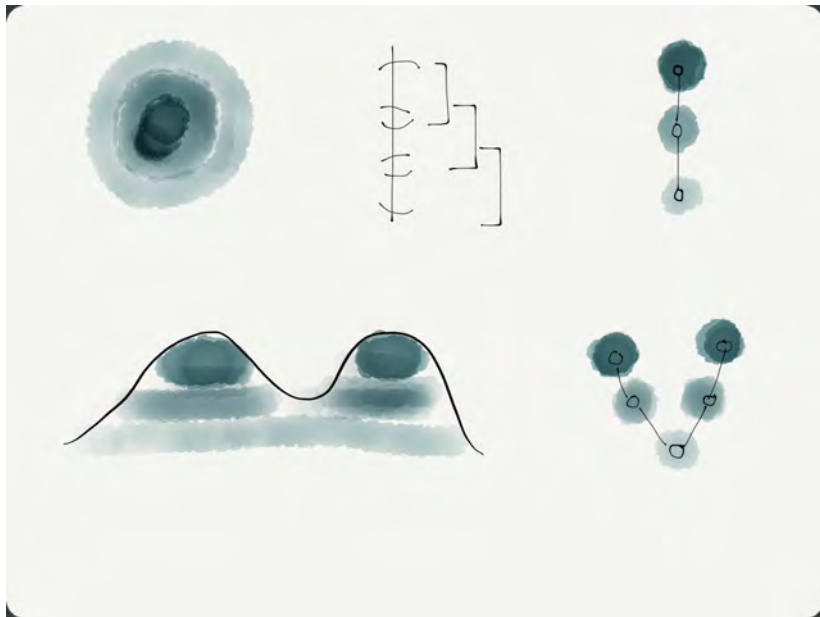
- Some kind of density measure.
- A score measure difference (distance) from some baseline.
- An eccentricity measure.

Clearly, choice of filter function is essential.

- Some kind of density measure.
- A score measure difference (distance) from some baseline.
- An eccentricity measure.

Clearly, choice of filter function is essential.

- Some kind of density measure.
- A score measure difference (distance) from some baseline.
- An eccentricity measure.



Breast cancer example

Highly successful example of real data analysis. [Nicolau, Carlsson, Levine]

- Working with vectors of gene expression data.
- Distance metric is correlation.
- Filter is a measure of (unsigned) deviation of expression from normal tissue.

Results identify previously unknown c-MYB+ region, which are very different from normal tissue but have very high survival rates.

Breast cancer example

Highly successful example of real data analysis. [Nicolau, Carlsson, Levine]

- Working with vectors of gene expression data.
- Distance metric is correlation.
- Filter is a measure of (unsigned) deviation of expression from normal tissue.

Results identify previously unknown c-MYB+ region, which are very different from normal tissue but have very high survival rates.

Breast cancer example

Highly successful example of real data analysis. [Nicolau, Carlsson, Levine]

- Working with vectors of gene expression data.
- Distance metric is correlation.
- Filter is a measure of (unsigned) deviation of expression from normal tissue.

Results identify previously unknown c-MYB+ region, which are very different from normal tissue but have very high survival rates.

Breast cancer example

Highly successful example of real data analysis. [Nicolau, Carlsson, Levine]

- Working with vectors of gene expression data.
- Distance metric is correlation.
- Filter is a measure of (unsigned) deviation of expression from normal tissue.

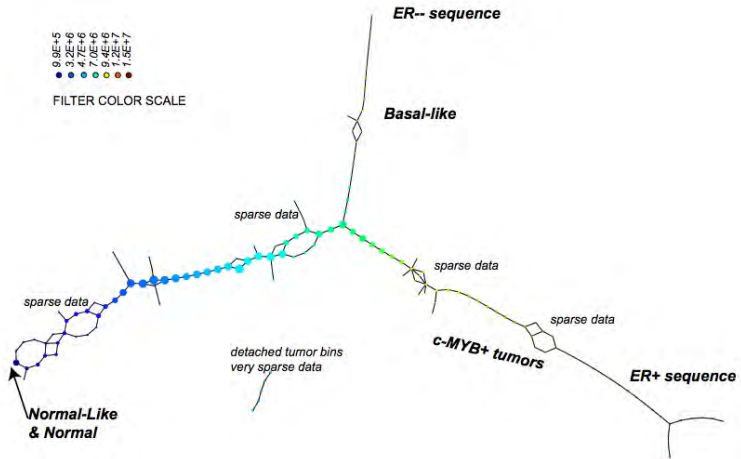
Results identify previously unknown c-MYB+ region, which are very different from normal tissue but have very high survival rates.

Breast cancer example

Highly successful example of real data analysis. [Nicolau, Carlsson, Levine]

- Working with vectors of gene expression data.
- Distance metric is correlation.
- Filter is a measure of (unsigned) deviation of expression from normal tissue.

Results identify previously unknown c-MYB⁺ region, which are very different from normal tissue but have very high survival rates.



Clever example of application to sports analytics. [Alagappan]

- Data set consists of vectors of statistics (points scored, rebounds, etc.).
- Distance metric is Euclidean.
- Filter is points per minute.

Results identify many “new” positions.

Clever example of application to sports analytics. [Alagappan]

- Data set consists of vectors of statistics (points scored, rebounds, etc.).
- Distance metric is Euclidean.
- Filter is points per minute.

Results identify many “new” positions.

Clever example of application to sports analytics. [Alagappan]

- Data set consists of vectors of statistics (points scored, rebounds, etc.).
- Distance metric is Euclidean.
- Filter is points per minute.

Results identify many “new” positions.

Clever example of application to sports analytics. [Alagappan]

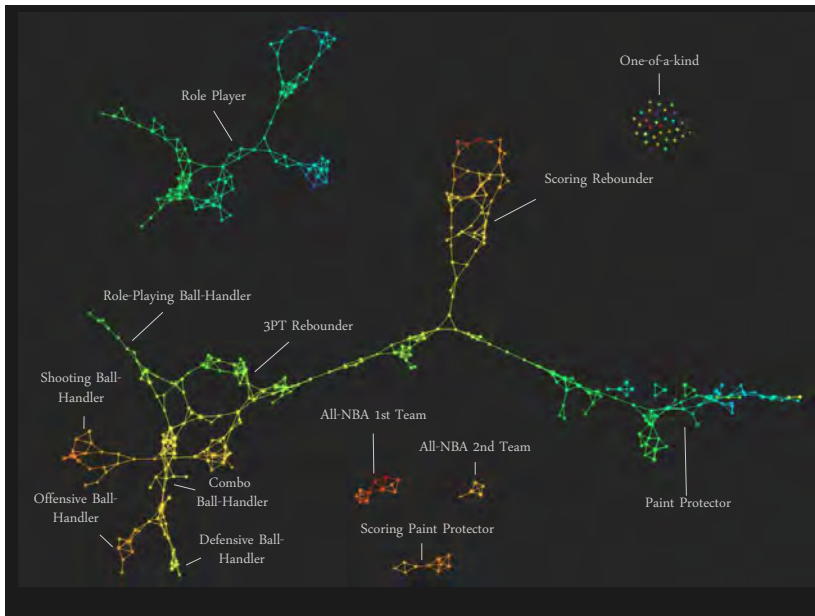
- Data set consists of vectors of statistics (points scored, rebounds, etc.).
- Distance metric is Euclidean.
- Filter is points per minute.

Results identify many “new” positions.

Clever example of application to sports analytics. [Alagappan]

- Data set consists of vectors of statistics (points scored, rebounds, etc.).
- Distance metric is Euclidean.
- Filter is points per minute.

Results identify many “new” positions.



Claim

Mapper can be successfully applied to analysis of geometric structures in large data sets from a wide variety of domains.

- Key idea: clustering across “scales”, represent relationships between clusters as scale varies
- Choice of filter function(s) is critical to successful application.

Claim

Mapper can be successfully applied to analysis of geometric structures in large data sets from a wide variety of domains.

- Key idea: clustering across “scales”, represent relationships between clusters as scale varies
- Choice of filter function(s) is critical to successful application.

Claim

Mapper can be successfully applied to analysis of geometric structures in large data sets from a wide variety of domains.

- Key idea: clustering across “scales”, represent relationships between clusters as scale varies
- Choice of filter function(s) is critical to successful application.