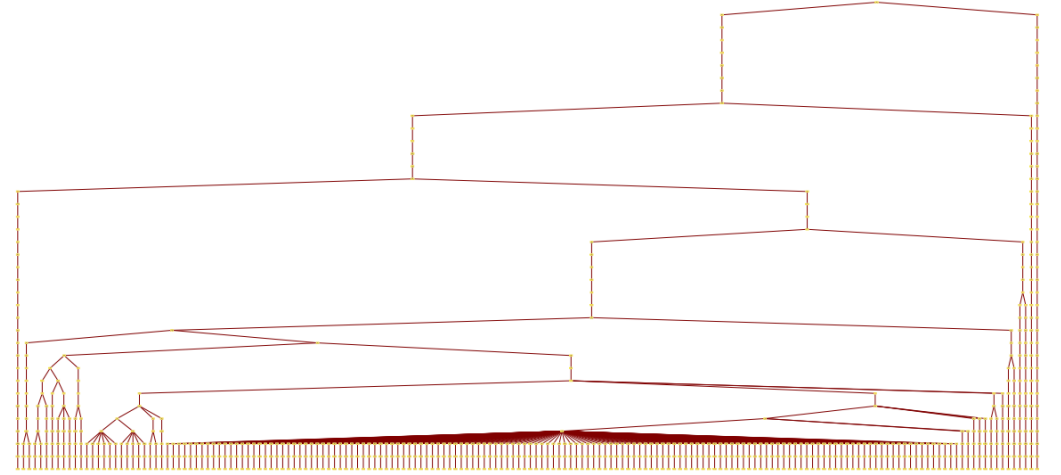
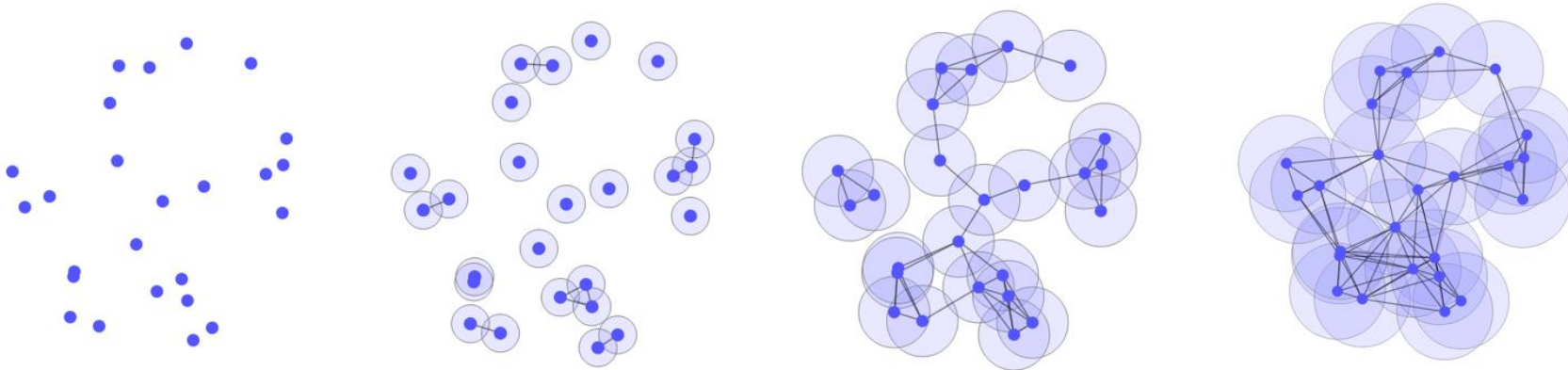


Where can we use graph algorithms we learned in ML?



Hierarchical Clustering



Hierarchical Clustering

Hierarchical clustering is a family of clustering algorithms that build a tree clusters. It is usually done by merging or splitting the clusters successively.

Hierarchical Clustering

Hierarchical clustering is a family of clustering algorithms that build a tree clusters. It is usually done by merging or splitting the clusters successively.

Hierarchical clustering is usually represented by a tree called the dendrogram that represents the clusterings at all levels.

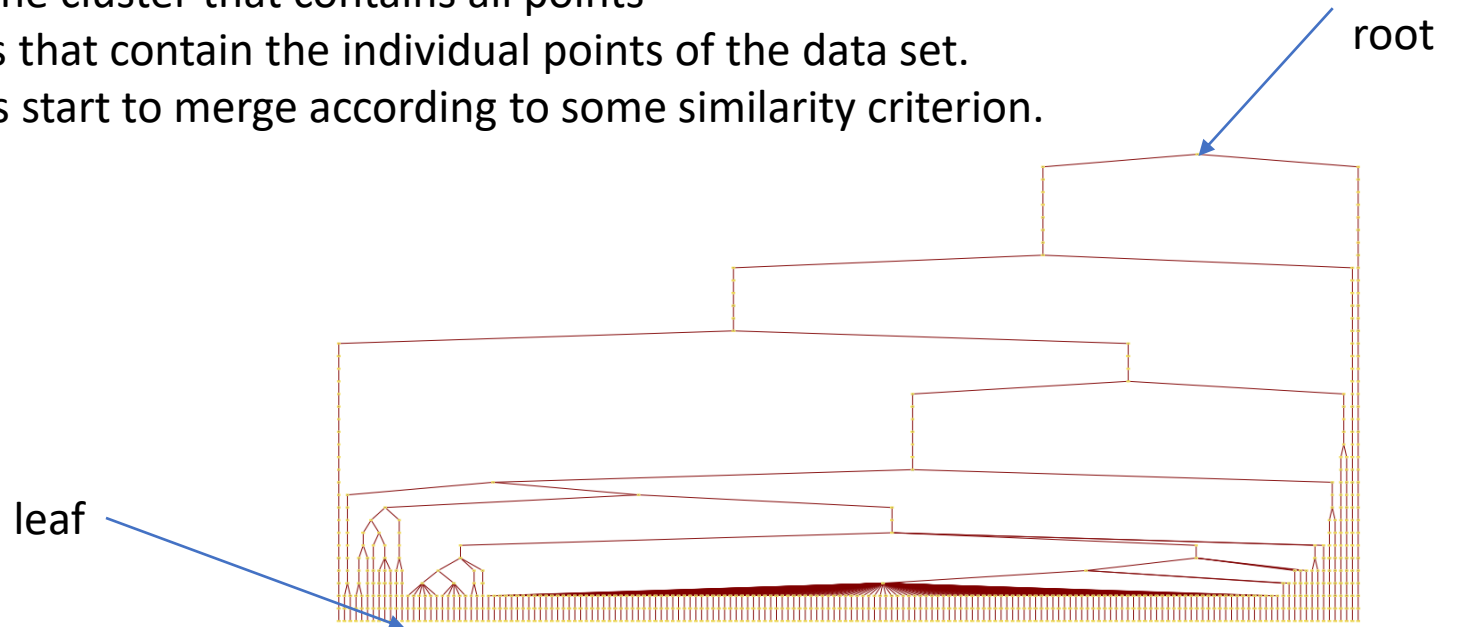
Each node in the tree represents a cluster

- In particular the root of the tree represents the cluster that contains all points
- The leaves of the tree represents the clusters that contain the individual points of the data set.
- As we go from the leaves to the root, clusters start to merge according to some similarity criterion.

There are two types of Hierarchical clustering

Agglomerative : bottom up (Merge).

Divisive : Top down (Split).



General Steps

General Steps in a standard hierarchical agglomerative clustering algorithm:

1. Compute the distance matrix, or the dissimilarity, between all points in the input data. Choosing the metric will impact the results largely.

General Steps

General Steps in a standard hierarchical agglomerative clustering algorithm:

1. Compute the distance matrix, or the dissimilarity, between all points in the input data. Choosing the metric will impact the results largely.
2. Initialize every point in the dataset to be its own cluster

General Steps

General Steps in a standard hierarchical agglomerative clustering algorithm:

1. Compute the distance matrix, or the dissimilarity, between all points in the input data. Choosing the metric will impact the results largely.
2. Initialize every point in the dataset to be its own cluster
3. Compute the distance between all clusters

General Steps

General Steps in a standard hierarchical agglomerative clustering algorithm:

1. Compute the distance matrix, or the dissimilarity, between all points in the input data. Choosing the metric will impact the results largely.
2. Initialize every point in the dataset to be its own cluster
3. Compute the distance between all clusters
4. Combine the closest two clusters: the two clusters c_i and c_j with $\min_{i,j} D(c_i, c_j)$

General Steps

General Steps in a standard hierarchical agglomerative clustering algorithm:

1. Compute the distance matrix, or the dissimilarity, between all points in the input data. Choosing the metric will impact the results largely.
2. Initialize every point in the dataset to be its own cluster
3. Compute the distance between all clusters
4. Combine the closest two clusters: the two clusters c_i and c_j with $\min_{i,j} D(c_i, c_j)$
5. Remove the clusters c_i and c_j and add the cluster $c_i + c_j$.

General Steps

General Steps in a standard hierarchical agglomerative clustering algorithm:

1. Compute the distance matrix, or the dissimilarity, between all points in the input data. Choosing the metric will impact the results largely.
2. Initialize every point in the dataset to be its own cluster
3. Compute the distance between all clusters
4. Combine the closest two clusters: the two clusters c_i and c_j with $\min_{i,j} D(c_i, c_j)$
5. Remove the clusters c_i and c_j and add the cluster $c_i + c_j$.
6. Go back to 3 and repeat until we have a single cluster.

General Steps

General Steps in a standard hierarchical agglomerative clustering algorithm:

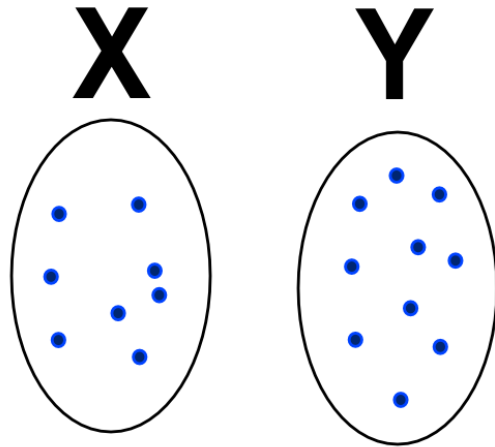
1. Compute the distance matrix, or the dissimilarity, between all points in the input data. Choosing the metric will impact the results largely.
2. Initialize every point in the dataset to be its own cluster
3. Compute the distance between all clusters
4. Combine the closest two clusters: the two clusters c_i and c_j with $\min_{i,j} D(c_i, c_j)$
5. Remove the clusters c_i and c_j and add the cluster $c_i + c_j$.
6. Go back to 3 and repeat until we have a single cluster.

Question : how do we measure the distance between two clusters ?

This is important because step 3 we need to measure the distance between clusters rather than points.

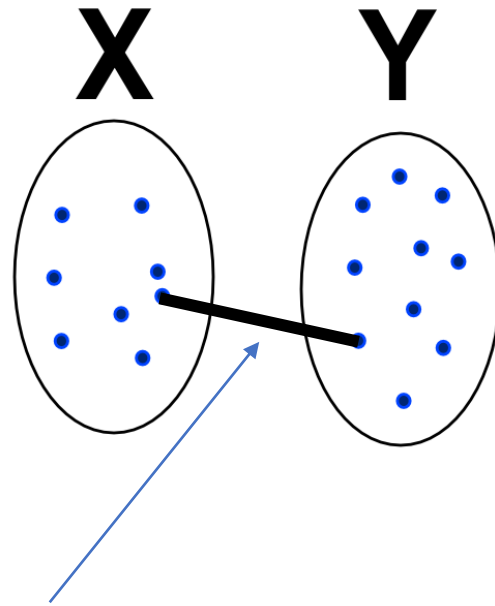
Distance between clusters

Given two clusters X and Y. How do we measure the distance between them ?



Distance between clusters

Given two clusters X and Y. How do we measure the distance between them ?

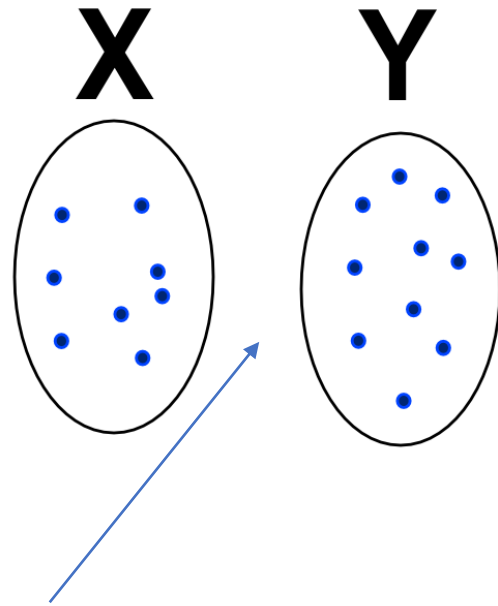


$$D(X, Y) := \min_{x \in X, y \in Y} d(x, y)$$

One way is to measure the minimal distance between all points of X and Y.
This distance induce [single linkage clustering](#).

Distance between clusters

Given two clusters X and Y. How do we measure the distance between them ?

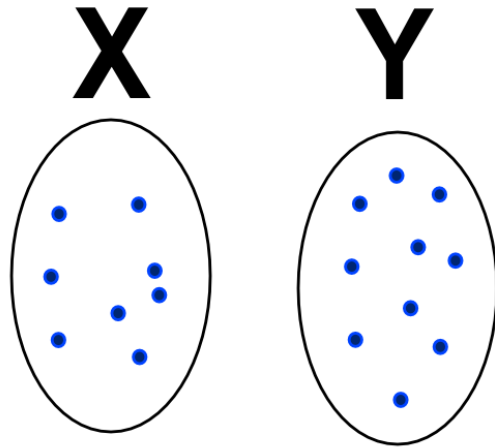


$$D(X, Y) := \frac{1}{|X||Y|} \sum_{x \in X, y \in Y}^n d(x, y)$$

We could also consider the mean distance between the points of the clusters

Distance between clusters

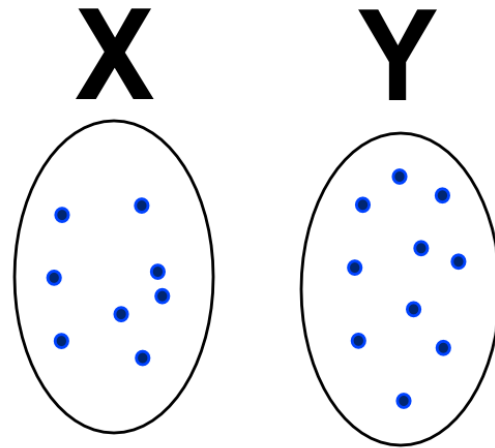
Given two clusters X and Y. How do we measure the distance between them ?



There are other measures as well : [The minimal energy criterion](#) , the distance between the centroids of the clusters

Distance between clusters

Given two clusters X and Y. How do we measure the distance between them ?



For efficient calculations we usually require the following condition:
Knowing the distance $D(A,C)$, $D(B,C)$
Implies we can calculate in constant time the distance $D(A+B,C)$

There are other measures as well : [The minimal energy criterion](#) , the distance between the centroids of the clusters

Single-linkage clustering algorithm

The input of the algorithm is a distance matrix D contains all distances $d(i,j)$ between the points. At the beginning each point is its own cluster.

Single-linkage clustering algorithm

1-Find the most similar pair of clusters in the current clustering, say pair $(r), (s)$, according to $d[(r),(s)] = \min d[(i),(j)]$ where the minimum is over all pairs of clusters in the current clustering.

Single-linkage clustering algorithm

- 1-Find the most similar pair of clusters in the current clustering, say pair $(r), (s)$, according to $d[(r),(s)] = \min d[(i),(j)]$ where the minimum is over all pairs of clusters in the current clustering.
- 2- Merge clusters (r) and (s) into a single cluster to form the next clustering m .

Single-linkage clustering algorithm

1-Find the most similar pair of clusters in the current clustering, say pair (r), (s), according to $d[(r),(s)] = \min d[(i),(j)]$ where the minimum is over all pairs of clusters in the current clustering.

2- Merge clusters (r) and (s) into a single cluster to form the next clustering m.

3-Update the distance matrix, D, by deleting the rows and columns corresponding to clusters (r) and (s) and adding a row and column corresponding to the newly formed cluster.

Single-linkage clustering algorithm

1-Find the most similar pair of clusters in the current clustering, say pair (r), (s), according to $d[(r),(s)] = \min d[(i),(j)]$ where the minimum is over all pairs of clusters in the current clustering.

2- Merge clusters (r) and (s) into a single cluster to form the next clustering m.

3-Update the distance matrix, D, by deleting the rows and columns corresponding to clusters (r) and (s) and adding a row and column corresponding to the newly formed cluster.

4- The distance between the new cluster, denoted (r,s) and old cluster (k) is defined as $d[(k), (r,s)] = \min d[(k),(r)], d[(k),(s)]$.

Single-linkage clustering algorithm

1-Find the most similar pair of clusters in the current clustering, say pair (r), (s), according to $d[(r),(s)] = \min d[(i),(j)]$ where the minimum is over all pairs of clusters in the current clustering.

2- Merge clusters (r) and (s) into a single cluster to form the next clustering m.

3-Update the distance matrix, D, by deleting the rows and columns corresponding to clusters (r) and (s) and adding a row and column corresponding to the newly formed cluster.

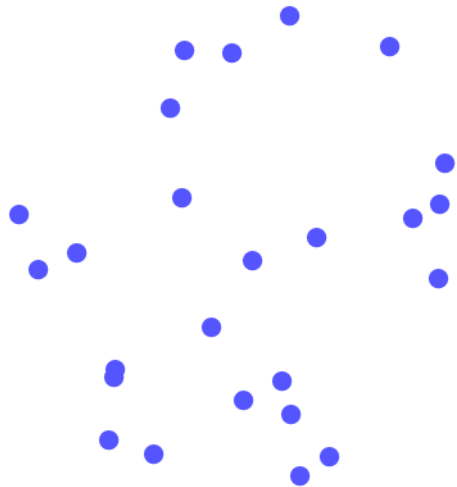
4- The distance between the new cluster, denoted (r,s) and old cluster (k) is defined as $d[(k), (r,s)] = \min d[(k),(r)], d[(k),(s)]$.

5- If all objects are in one cluster, stop. Else, go to step 1.

Single Linkage Hierarchical Clustering and the ε - Neighborhood Graph

Suppose that we are given a set of points $X = \{p_1, p_2, \dots, p_n\}$ in R^d with a distance function d defined on them.

Consider the connected components of the ε -neighborhood graph as we continuously increase ε from zero to infinity.

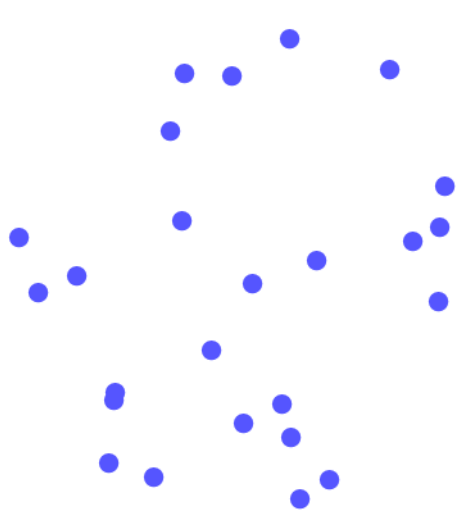


Every point is a connected component

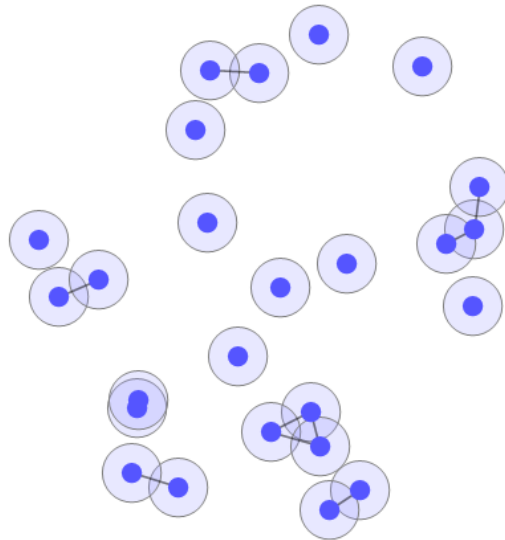
Single Linkage Hierarchical Clustering and the ε - Neighborhood Graph

Suppose that we are given a set of points $X = \{p_1, p_2, \dots, p_n\}$ in R^d with a distance function d defined on them.

Consider the connected components of the ε -neighborhood graph as we continuously increase ε from zero to infinity.



Every point is a connected component

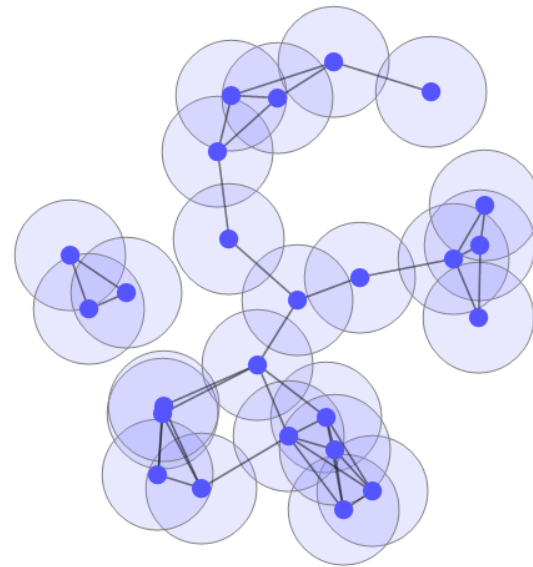
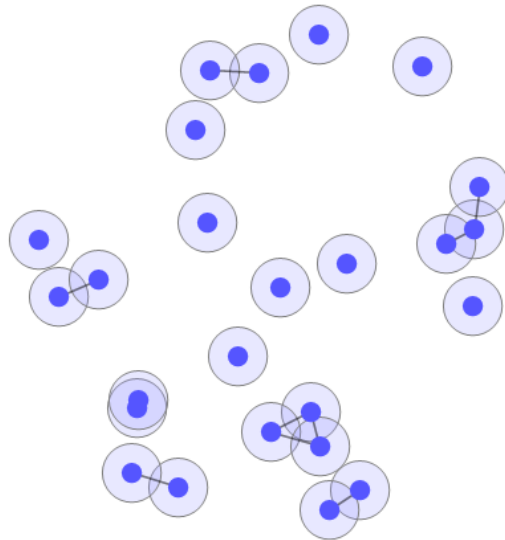
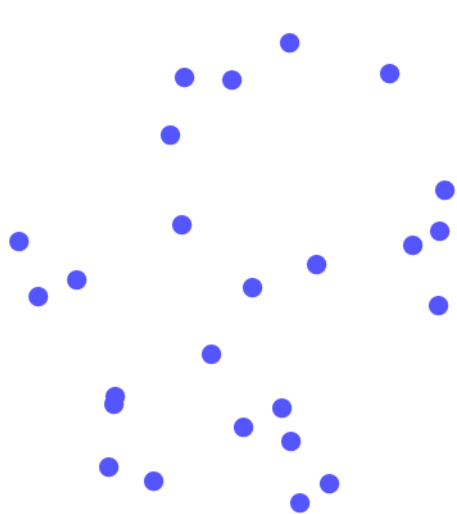


When ε is a little larger we start
some clusters starts to get form

Single Linkage Hierarchical Clustering and the ε - Neighborhood Graph

Suppose that we are given a set of points $X = \{p_1, p_2, \dots, p_n\}$ in R^d with a distance function d defined on them.

Consider the connected components of the ε -neighborhood graph as we continuously increase ε from zero to infinity.



Every point is a connected component

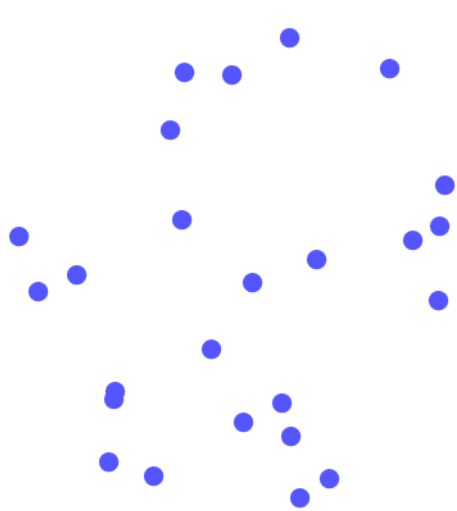
When ε is a little larger we start some clusters starts to get form

When ε is even larger we have few clusters
As the clusters get larger and larger

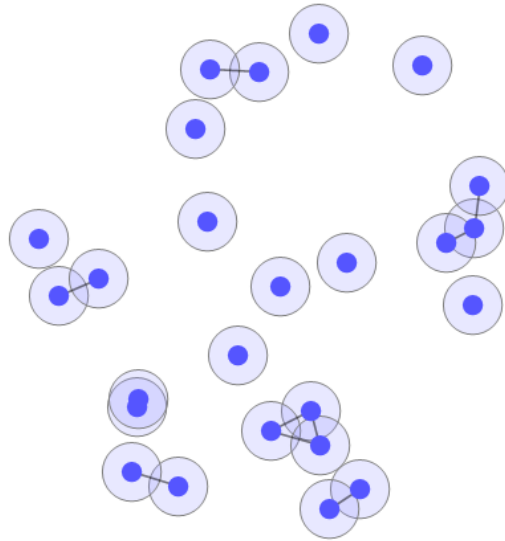
Single Linkage Hierarchical Clustering and the ε - Neighborhood Graph

Suppose that we are given a set of points $X = \{p_1, p_2, \dots, p_n\}$ in R^d with a distance function d defined on them.

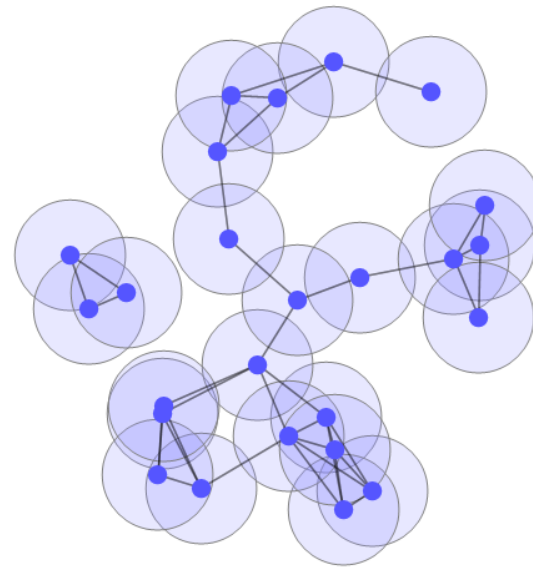
Consider the connected components of the ε -neighborhood graph as we continuously increase ε from zero to infinity.



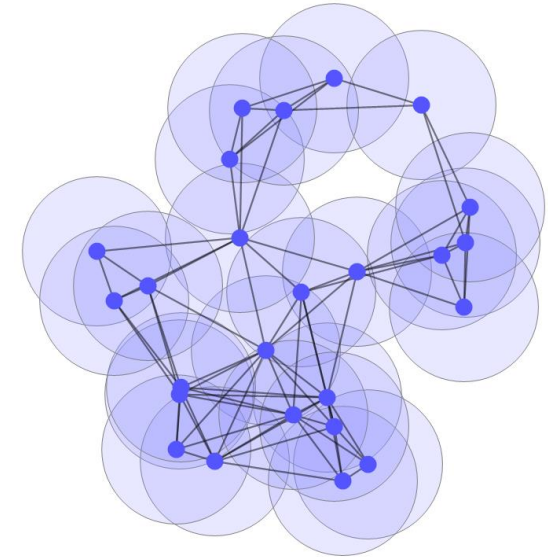
Every point is a connected component



When ε is a little larger we start some clusters starts to get form

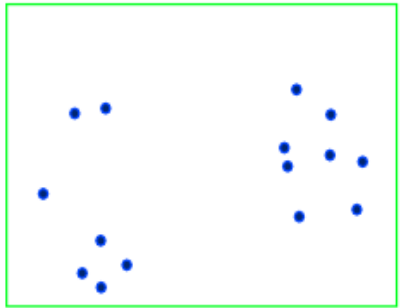


When ε is even larger we have few clusters
As the clusters get larger and larger

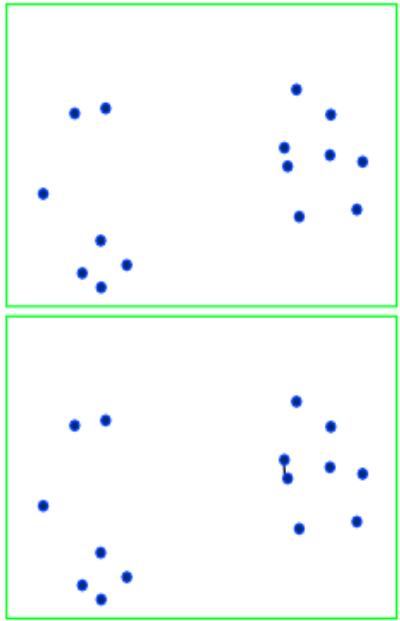


At some point all points become a part of a single cluster

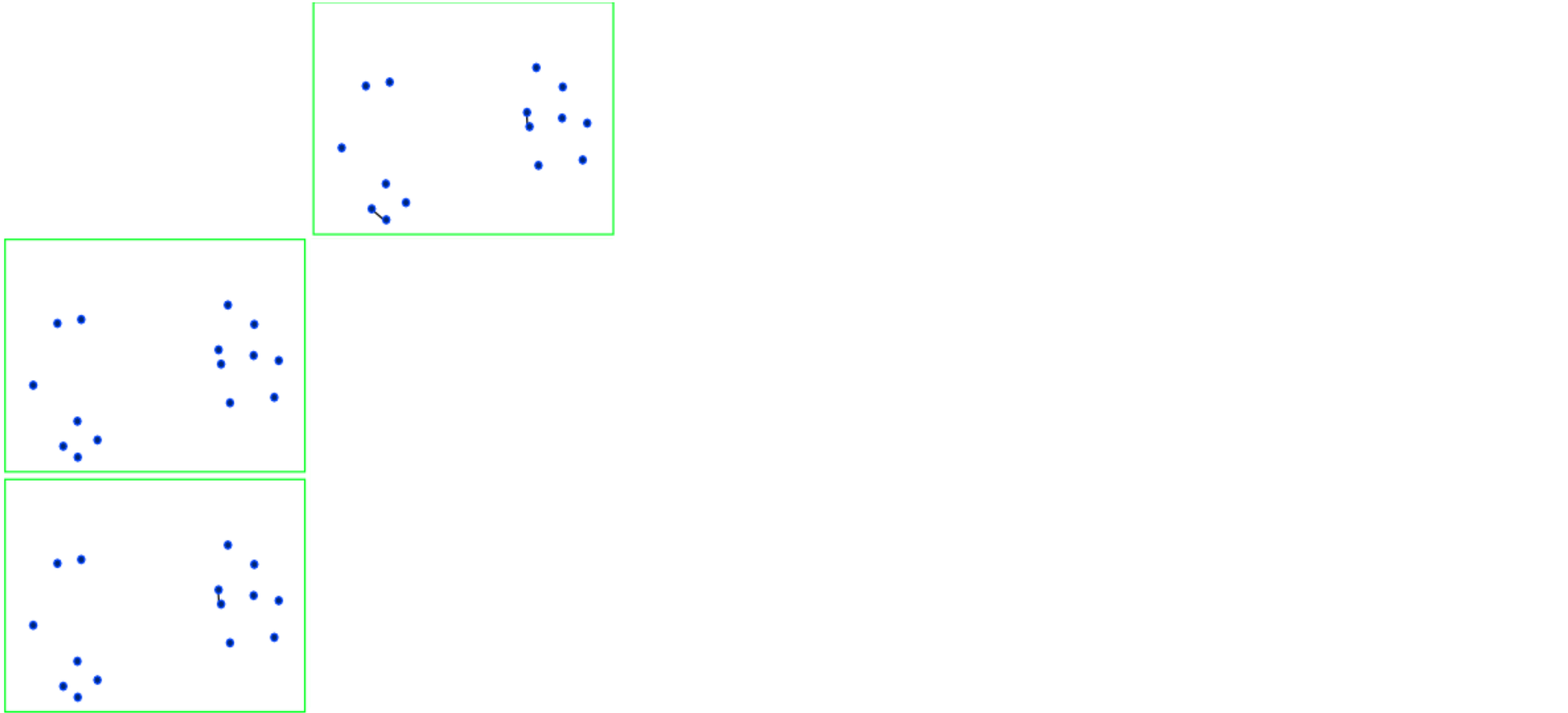
Single Linkage Hierarchical Clustering and the and Kruskal's algorithm



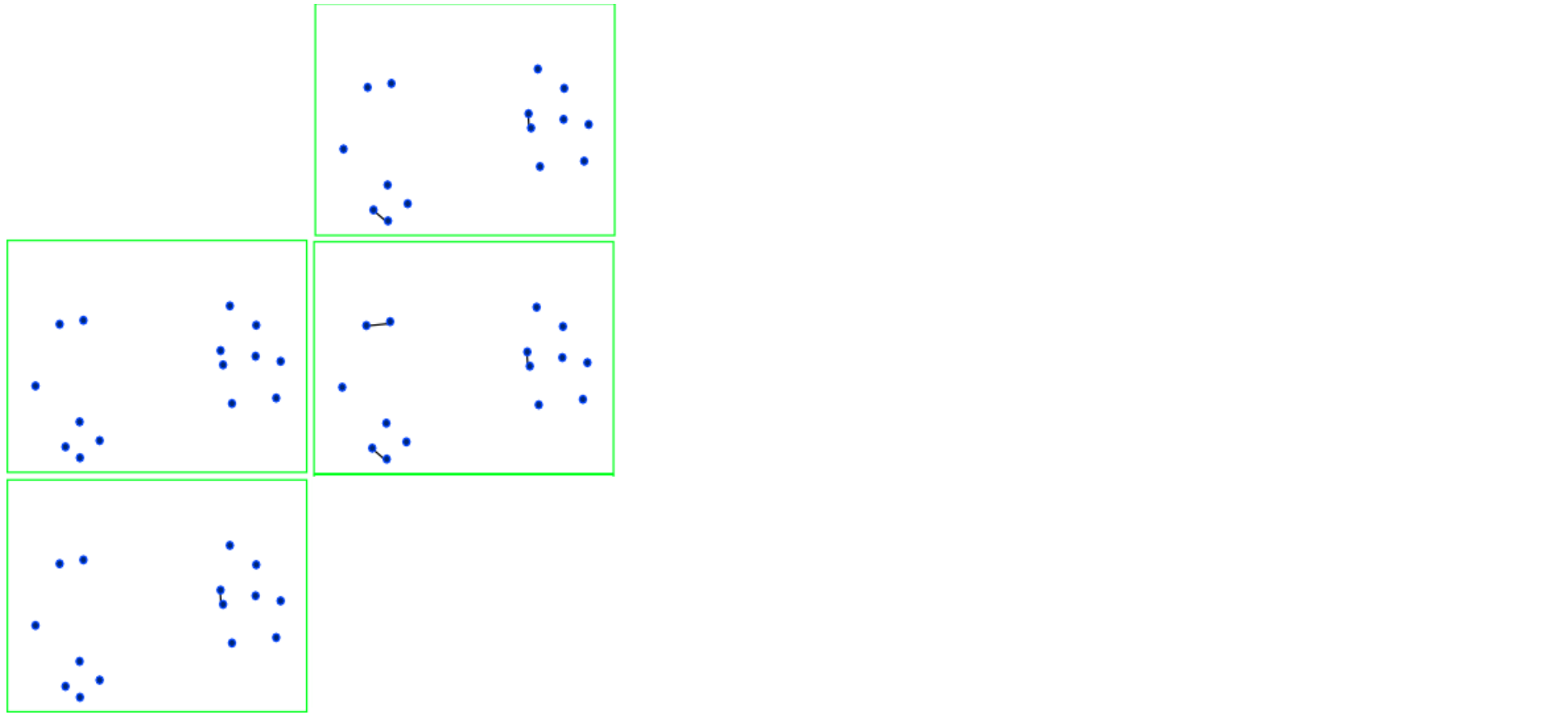
Single Linkage Hierarchical Clustering and the and Kruskal's algorithm



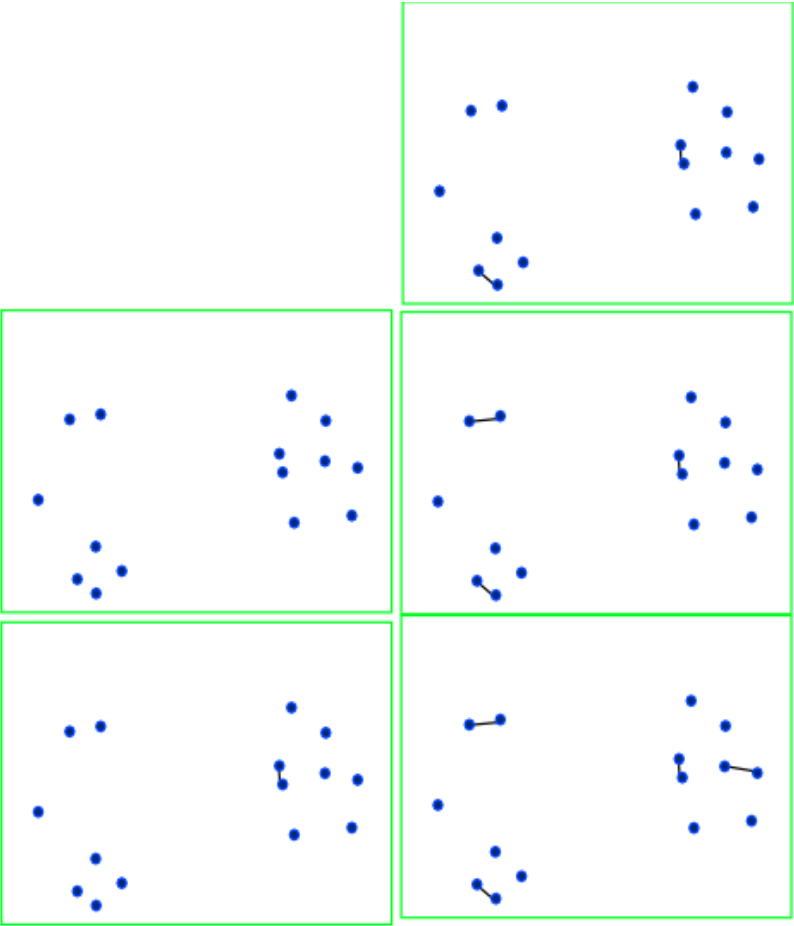
Single Linkage Hierarchical Clustering and the and Kruskal's algorithm



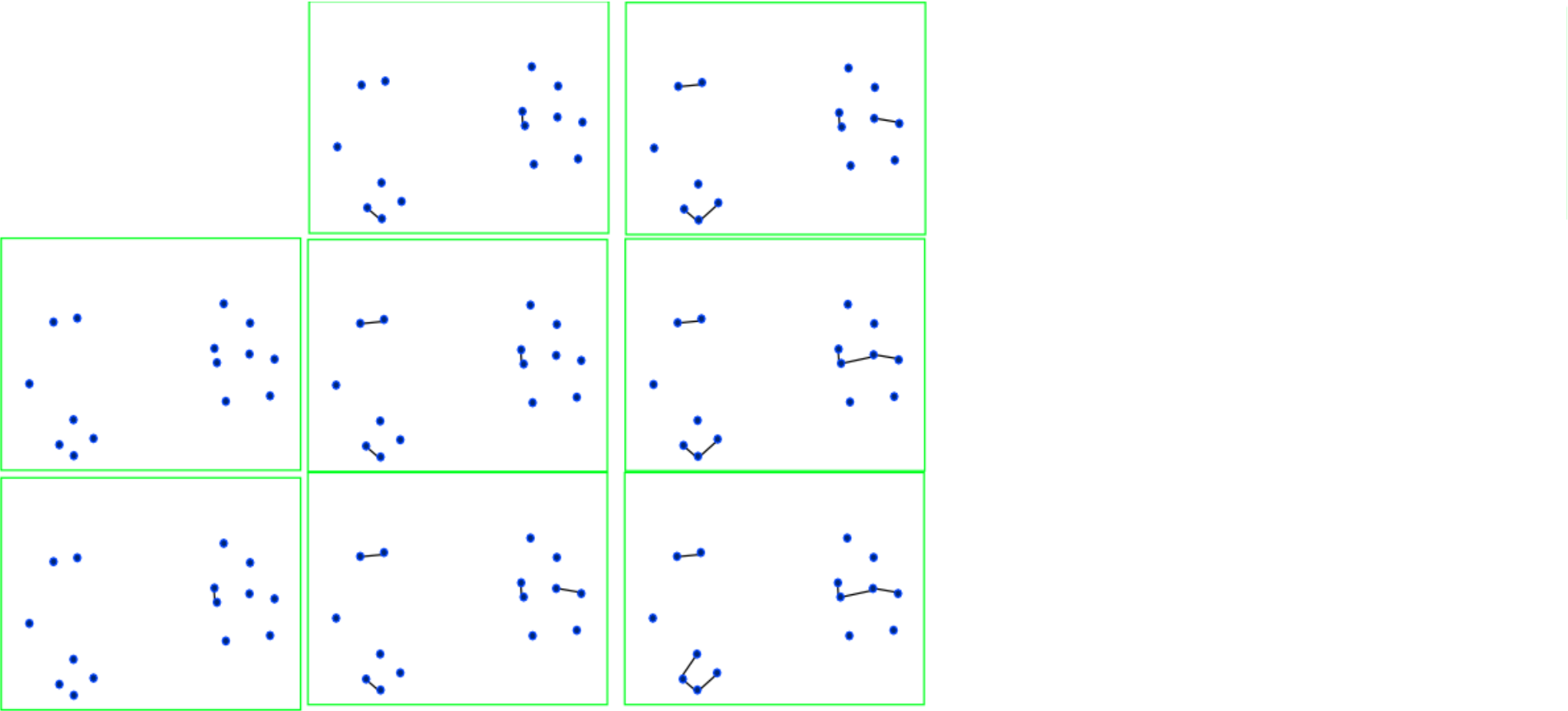
Single Linkage Hierarchical Clustering and the and Kruskal's algorithm



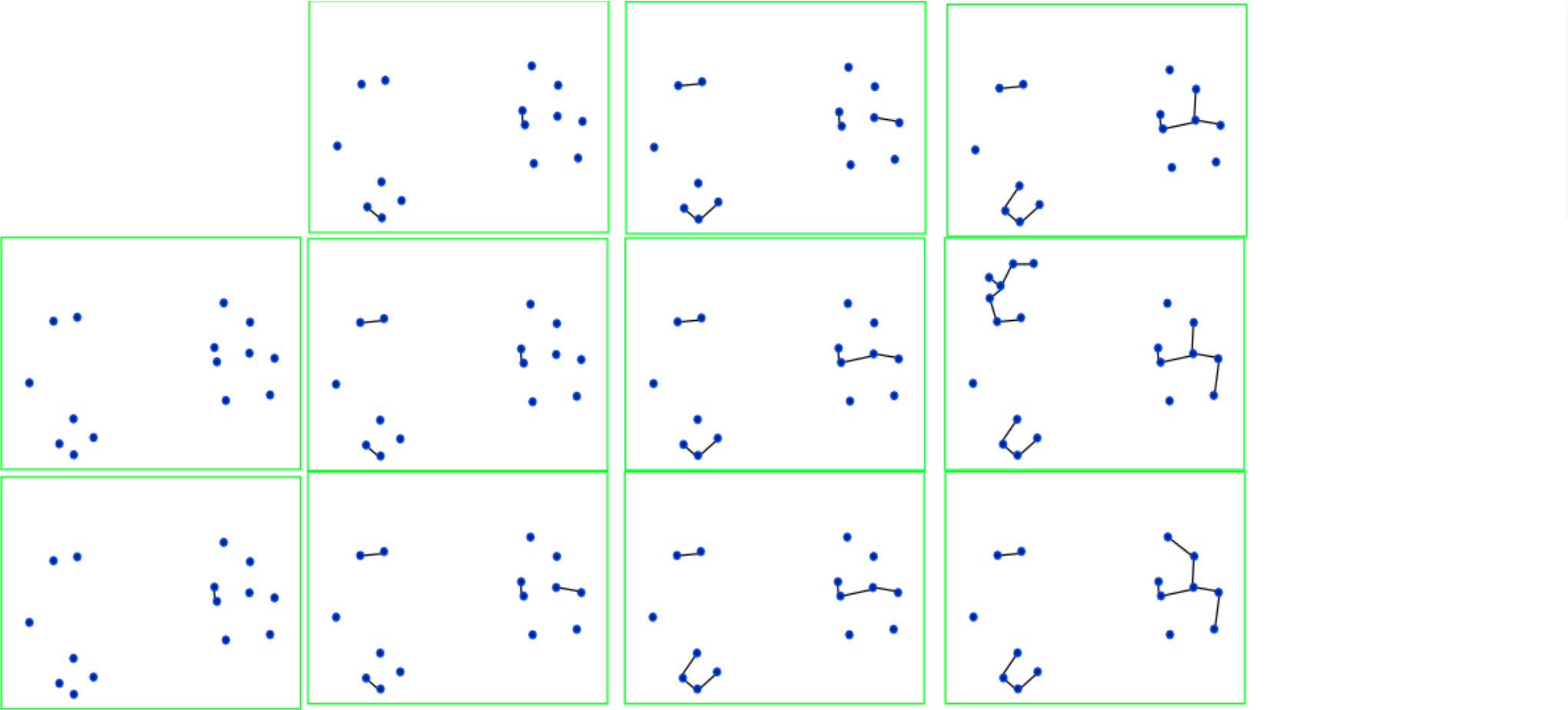
Single Linkage Hierarchical Clustering and the and Kruskal's algorithm



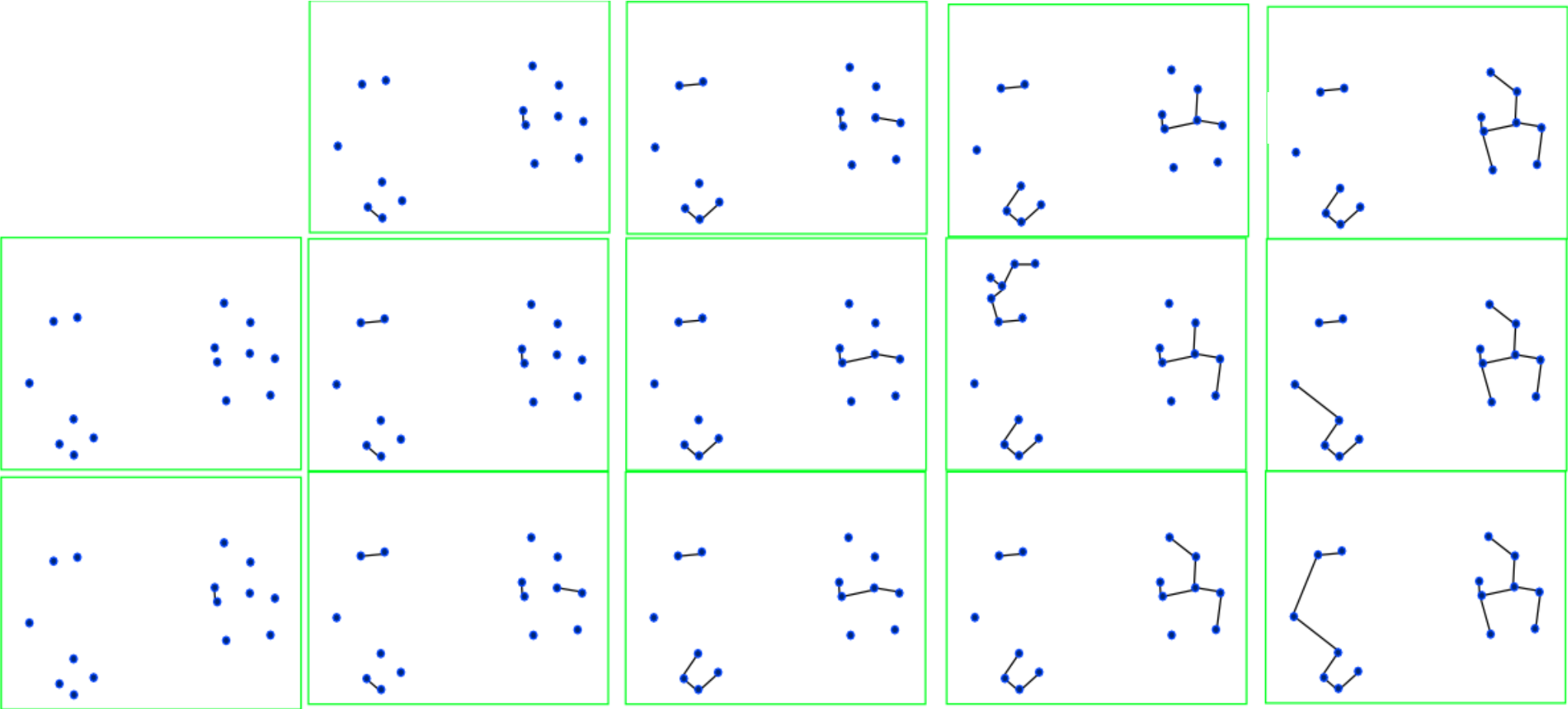
Single Linkage Hierarchical Clustering and the and Kruskal's algorithm



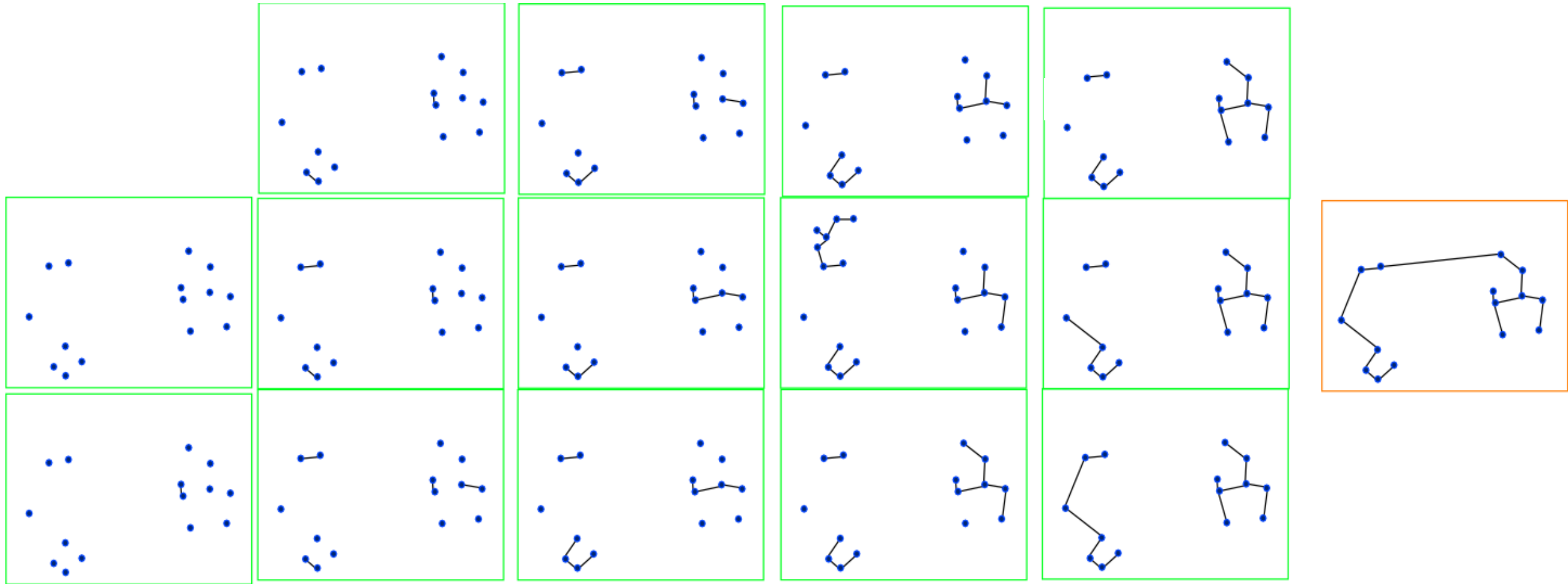
Single Linkage Hierarchical Clustering and the and Kruskal's algorithm



Single Linkage Hierarchical Clustering and the and Kruskal's algorithm



Single Linkage Hierarchical Clustering and the and Kruskal's algorithm



Lance–Williams algorithms

1. Compute the distance matrix : $D = \{D_{ij} : \text{distance between } i \text{ and } j \text{ for } i, j \text{ between } 1 \text{ and } n\}$
2. Iterate n times :
 1. Find i and j with $\min_{i,j} D(c_i, c_j)$
 2. Add the cluster $i + j$ and delete the clusters i and j
 3. For each remaining cluster k
 1. $D_{k,i+j} = \min\{D_{k,j}, D_{k,i}\}$

Single Linkage H-Clustering



Lance–Williams algorithms

1. Compute the distance matrix : $D = \{D_{ij} : \text{distance between } i \text{ and } j \text{ for } i, j \text{ between } 1 \text{ and } n\}$
2. Iterate n times :
 1. Find i and j with $\min_{i,j} D(c_i, c_j)$
 2. Add the cluster $i + j$ and delete the clusters i and j
 3. For each remaining cluster k
 1. $D_{k,i+j} = \min\{D_{k,j}, D_{k,i}\}$

Single Linkage H-Clustering

This can be replaced by a more general condition.

Lance–Williams algorithms

1. Compute the distance matrix : $D = \{D_{ij} : \text{distance between } i \text{ and } j \text{ for } i, j \text{ between } 1 \text{ and } n\}$
2. Iterate n times :
 1. Find i and j with $\min_{i,j} D(c_i, c_j)$
 2. Add the cluster $i + j$ and delete the clusters i and j
 3. For each remaining cluster k
 1. $D_{k,i+j} = a_i D_{k,i} + a_j D_{k,j} + \beta D_{i,j} + \alpha |D_{k,j} - D_{k,i}|$

Every algorithm has a special
 a_i, a_j, β and α

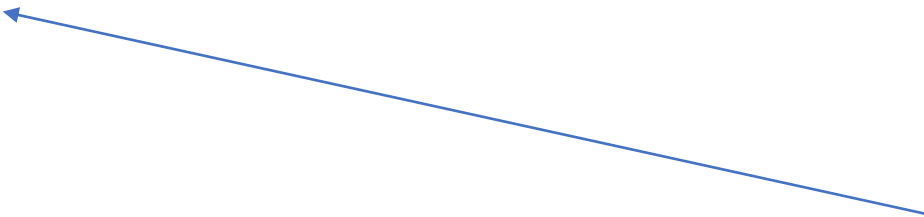
The neat thing about the algorithm : all needs to be changed
is how to *update* the new distance



Lance–Williams algorithms

1. Compute the distance matrix : $D = \{D_{ij} : \text{distance between } i \text{ and } j \text{ for } i, j \text{ between } 1 \text{ and } n\}$
2. Iterate n times :
 1. Find i and j with $\min_{i,j} D(c_i, c_j)$
 2. Add the cluster $i + j$ and delete the clusters i and j
 3. For each remaining cluster k
 1. $D_{k,i+j} = a_i D_{k,i} + a_j D_{k,j} + \beta D_{i,j} + \alpha |D_{k,j} - D_{k,i}|$

Every algorithm has a special
 a_i, a_j, β and α




For complete linkage, choose $a_i = ?$ $a_j = ?$ $\alpha = ?$ and $\beta = ?$

Lance–Williams algorithms

1. Compute the distance matrix : $D = \{D_{ij} : \text{distance between } i \text{ and } j \text{ for } i, j \text{ between } 1 \text{ and } n\}$
2. Iterate n times :
 1. Find i and j with $\min_{i,j} D(c_i, c_j)$
 2. Add the cluster $i + j$ and delete the clusters i and j
 3. For each remaining cluster k
 1. $D_{k,i+j} = a_i D_{k,i} + a_j D_{k,j} + \beta D_{i,j} + \alpha |D_{k,j} - D_{k,i}|$

Ward's method : $a_i = \frac{n_i + n_k}{n_i + n_j + n_k}, \beta = \frac{-n_k}{n_i + n_j + n_k}, \alpha = 0$



Lance–Williams algorithms

What is the complexity here ?

1. Compute the distance matrix : $D = \{D_{ij} : \text{distance between } i \text{ and } j \text{ for } i, j \text{ between } 1 \text{ and } n\}$ $\longleftarrow n^2$ steps
2. Iterate n times : $\longleftarrow n$ steps
 1. Find i and j with $\min_{i,j} D(c_i, c_j)$ $\longleftarrow n^2$ steps
 2. Add the cluster $i + j$ and delete the clusters i and j
 3. For each remaining cluster k $\longleftarrow n$ steps at most
 1. $D_{k,i+j} = a_i D_{k,i} + a_j D_{k,j} + \beta D_{i,j} + \alpha |D_{k,j} - D_{k,i}|$ \longleftarrow Constant time!

This is just the naïve implementation, there are better algorithms that perform with $O(n^2)$

In Sklearn

Scikit learn supports [Agglomerative Clustering](#). Many features discussed in this lecture are also supported.

An Introduction to Multidimensional Scaling and ISOMAP

MDS

Let $D=[d_{ij}]$ be an $N \times N$ dissimilarity matrix.

In MDS, we want to find n vectors x_1, \dots, x_N in R^d such that $\|x_i - x_j\| \approx d_{ij}$

MDS

Let $D=[d_{ij}]$ be an $N \times N$ dissimilarity matrix.

In MDS, we want to find n vectors x_1, \dots, x_N in R^d such that $\|x_i - x_j\| \approx d_{ij}$

- Usually if we choose d to be large enough, we can construct the vectors x_1, \dots, x_N with exact solutions : $\|x_i - x_j\| = d_{ij}$.

MDS

Let $D=[d_{ij}]$ be an $N \times N$ dissimilarity matrix.

In MDS, we want to find n vectors x_1, \dots, x_N in R^d such that $\|x_i - x_j\| \approx d_{ij}$

- Usually if we choose d to be large enough, we can construct the vectors x_1, \dots, x_N with exact solutions : $\|x_i - x_j\| = d_{ij}$.
- In this case the distance d above is the usual Euclidean distance.

MDS

Let $D=[d_{ij}]$ be an $N \times N$ dissimilarity matrix.

In MDS, we want to find n vectors x_1, \dots, x_N in R^d such that $||x_i - x_j|| \approx d_{ij}$

- Usually if we choose d to be large enough, we can construct the vectors x_1, \dots, x_N with exact solutions : $||x_i - x_j|| = d_{ij}$.
- In this case the distance d above is the usual Euclidean distance.
- There are cases where the matrix D is valid distance matrix, but still there exists no set of vectors x_1, \dots, x_N in any R^d with perfect $||x_i - x_j|| = d_{ij}$. Such a distance is called non-Euclidean distance.

Classical MDS Algorithm

1-Construct the matrix of squares of the distances $P^{(2)} = [d_{ij}^2]$.

Classical MDS Algorithm

1-Construct the matrix of squares of the distances $P^{(2)} = [d_{ij}^2]$.

2. Apply the double centering: $B = -\frac{1}{2}JP^{(2)}J$ where $J = I - \frac{1}{n}11'$, where n is the number of elements.

Classical MDS Algorithm

- 1-Construct the matrix of squares of the distances $P^{(2)} = [d_{ij}^2]$.
2. Apply the double centering: $B = -\frac{1}{2}JP^{(2)}J$ where $J = I - \frac{1}{n}11'$, where n is the number of elements.
3. Extract the largest d positive eigenvalues $\lambda_1 \dots \lambda_d$ of B and the corresponding m eigenvectors $e_1 \dots e_d$.

Classical MDS Algorithm

- 1-Construct the matrix of squares of the distances $P^{(2)} = [d_{ij}^2]$.
2. Apply the double centering: $B = -\frac{1}{2}JP^{(2)}J$ where $J = I - \frac{1}{n}11'$, where n is the number of elements.
3. Extract the largest d positive eigenvalues $\lambda_1 \dots \lambda_d$ of B and the corresponding d eigenvectors $e_1 \dots e_d$.
4. A d -dimensional MDS coordinates of n objects is derived from the coordinate matrix $X = E_d\Lambda_d^{\frac{1}{2}}$, where E_d is the matrix of d eigenvectors and Λ_d is the diagonal matrix of d eigenvalues of B , respectively

Example

Start with a distance matrix D

$$D = \begin{matrix} & \begin{matrix} 0 & 93 & 82 & 133 \end{matrix} \\ \begin{matrix} 93 \\ 82 \\ 133 \end{matrix} & \begin{matrix} 0 & 52 & 111 \\ 52 & 0 & 60 \\ 60 & 111 & 0 \end{matrix} \end{matrix}$$

See [this](#) for more details

Example

Take the square the elements of D

$$D = \begin{bmatrix} 0 & 93 & 82 & 133 \\ 93 & 0 & 52 & 60 \\ 82 & 52 & 0 & 111 \\ 133 & 60 & 111 & 0 \end{bmatrix} \quad \longrightarrow \quad P^{(2)} = \begin{bmatrix} 0 & 8649 & 6724 & 17689 \\ 8649 & 0 & 2704 & 3600 \\ 6724 & 2704 & 0 & 12321 \\ 17689 & 3600 & 12321 & 0 \end{bmatrix}$$

Example

Construct the J matrix

$$D = \begin{bmatrix} 0 & 93 & 82 & 133 \\ 93 & 0 & 52 & 60 \\ 82 & 52 & 0 & 111 \\ 133 & 60 & 111 & 0 \end{bmatrix} \quad \longrightarrow \quad \mathbf{P}^{(2)} = \begin{bmatrix} 0 & 8649 & 6724 & 17689 \\ 8649 & 0 & 2704 & 3600 \\ 6724 & 2704 & 0 & 12321 \\ 17689 & 3600 & 12321 & 0 \end{bmatrix}$$

$$\mathbf{J} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} - 0.25 \times \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 0.75 & -0.25 & -0.25 & -0.25 \\ -0.25 & 0.75 & -0.25 & -0.25 \\ -0.25 & -0.25 & 0.75 & -0.25 \\ -0.25 & -0.25 & -0.25 & 0.75 \end{bmatrix}$$

Example

Construct double centering matrix

$$D = \begin{bmatrix} 0 & 93 & 82 & 133 \\ 93 & 0 & 52 & 60 \\ 82 & 52 & 0 & 111 \\ 133 & 60 & 111 & 0 \end{bmatrix} \quad \longrightarrow \quad P^{(2)} = \begin{bmatrix} 0 & 8649 & 6724 & 17689 \\ 8649 & 0 & 2704 & 3600 \\ 6724 & 2704 & 0 & 12321 \\ 17689 & 3600 & 12321 & 0 \end{bmatrix}$$

$$J = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} - 0.25 \times \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 0.75 & -0.25 & -0.25 & -0.25 \\ -0.25 & 0.75 & -0.25 & -0.25 \\ -0.25 & -0.25 & 0.75 & -0.25 \\ -0.25 & -0.25 & -0.25 & 0.75 \end{bmatrix}$$

$$B = -\frac{1}{2}JP^{(2)}J = \begin{bmatrix} 5035.0625 & -1553.0625 & 258.9375 & -3740.938 \\ -1553.0625 & 507.8125 & 5.3125 & 1039.938 \\ 258.9375 & 5.3125 & 2206.8125 & -2471.062 \\ -3740.9375 & 1039.9375 & -2471.0625 & 5172.062 \end{bmatrix}$$

Example

Solve the largest 2 eigenvalues and eigenvector of B

$$D = \begin{bmatrix} 0 & 93 & 82 & 133 \\ 93 & 0 & 52 & 60 \\ 82 & 52 & 0 & 111 \\ 133 & 60 & 111 & 0 \end{bmatrix} \quad \longrightarrow \quad P^{(2)} = \begin{bmatrix} 0 & 8649 & 6724 & 17689 \\ 8649 & 0 & 2704 & 3600 \\ 6724 & 2704 & 0 & 12321 \\ 17689 & 3600 & 12321 & 0 \end{bmatrix}$$

$$J = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} - 0.25 \times \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 0.75 & -0.25 & -0.25 & -0.25 \\ -0.25 & 0.75 & -0.25 & -0.25 \\ -0.25 & -0.25 & 0.75 & -0.25 \\ -0.25 & -0.25 & -0.25 & 0.75 \end{bmatrix}$$

$$B = -\frac{1}{2}JP^{(2)}J = \begin{bmatrix} 5035.0625 & -1553.0625 & 258.9375 & -3740.938 \\ -1553.0625 & 507.8125 & 5.3125 & 1039.938 \\ 258.9375 & 5.3125 & 2206.8125 & -2471.062 \\ -3740.9375 & 1039.9375 & -2471.0625 & 5172.062 \end{bmatrix}$$

$$\lambda_1 = 9724.168, \lambda_2 = 3160.986, \quad e_1 = \begin{pmatrix} -0.637 \\ 0.187 \\ -0.253 \\ 0.704 \end{pmatrix}, \quad e_2 = \begin{pmatrix} -0.586 \\ 0.214 \\ 0.706 \\ -0.334 \end{pmatrix}$$

Example

Use that to construct the final MDS coordinates.

$$D = \begin{bmatrix} 0 & 93 & 82 & 133 \\ 93 & 0 & 52 & 60 \\ 82 & 52 & 0 & 111 \\ 133 & 60 & 111 & 0 \end{bmatrix} \longrightarrow P^{(2)} = \begin{bmatrix} 0 & 8649 & 6724 & 17689 \\ 8649 & 0 & 2704 & 3600 \\ 6724 & 2704 & 0 & 12321 \\ 17689 & 3600 & 12321 & 0 \end{bmatrix}$$

$$J = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} - 0.25 \times \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 0.75 & -0.25 & -0.25 & -0.25 \\ -0.25 & 0.75 & -0.25 & -0.25 \\ -0.25 & -0.25 & 0.75 & -0.25 \\ -0.25 & -0.25 & -0.25 & 0.75 \end{bmatrix}$$

$$B = -\frac{1}{2}JP^{(2)}J = \begin{bmatrix} 5035.0625 & -1553.0625 & 258.9375 & -3740.938 \\ -1553.0625 & 507.8125 & 5.3125 & 1039.938 \\ 258.9375 & 5.3125 & 2206.8125 & -2471.062 \\ -3740.9375 & 1039.9375 & -2471.0625 & 5172.062 \end{bmatrix}$$

$$\lambda_1 = 9724.168, \lambda_2 = 3160.986, \quad e_1 = \begin{pmatrix} -0.637 \\ 0.187 \\ -0.253 \\ 0.704 \end{pmatrix}, \quad e_2 = \begin{pmatrix} -0.586 \\ 0.214 \\ 0.706 \\ -0.334 \end{pmatrix} \longrightarrow X = \begin{bmatrix} -0.637 & -0.586 \\ 0.187 & 0.214 \\ -0.253 & 0.706 \\ 0.704 & -0.334 \end{bmatrix} \begin{bmatrix} \sqrt{9724.168} & 0 \\ 0 & \sqrt{3160.986} \end{bmatrix} = \begin{bmatrix} -62.831 & -32.97448 \\ 18.403 & 12.02697 \\ -24.960 & 39.71091 \\ 69.388 & -18.76340 \end{bmatrix}$$

Example

Use that to construct the final MDS coordinates.

$$D = \begin{bmatrix} 0 & 93 & 82 & 133 \\ 93 & 0 & 52 & 60 \\ 82 & 52 & 0 & 111 \\ 133 & 60 & 111 & 0 \end{bmatrix} \longrightarrow P^{(2)} = \begin{bmatrix} 0 & 8649 & 6724 & 17689 \\ 8649 & 0 & 2704 & 3600 \\ 6724 & 2704 & 0 & 12321 \\ 17689 & 3600 & 12321 & 0 \end{bmatrix}$$

$$J = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} - 0.25 \times \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 0.75 & -0.25 & -0.25 & -0.25 \\ -0.25 & 0.75 & -0.25 & -0.25 \\ -0.25 & -0.25 & 0.75 & -0.25 \\ -0.25 & -0.25 & -0.25 & 0.75 \end{bmatrix}$$

$$B = -\frac{1}{2}JP^{(2)}J = \begin{bmatrix} 5035.0625 & -1553.0625 & 258.9375 & -3740.938 \\ -1553.0625 & 507.8125 & 5.3125 & 1039.938 \\ 258.9375 & 5.3125 & 2206.8125 & -2471.062 \\ -3740.9375 & 1039.9375 & -2471.0625 & 5172.062 \end{bmatrix}$$

Coordinates of first point

$$\lambda_1 = 9724.168, \lambda_2 = 3160.986, \mathbf{e}_1 = \begin{pmatrix} -0.637 \\ 0.187 \\ -0.253 \\ 0.704 \end{pmatrix}, \mathbf{e}_2 = \begin{pmatrix} -0.586 \\ 0.214 \\ 0.706 \\ -0.334 \end{pmatrix} \longrightarrow \mathbf{X} = \begin{bmatrix} -0.637 & -0.586 \\ 0.187 & 0.214 \\ -0.253 & 0.706 \\ 0.704 & -0.334 \end{bmatrix} \begin{bmatrix} \sqrt{9724.168} & 0 \\ 0 & \sqrt{3160.986} \end{bmatrix} = \begin{bmatrix} -62.831 & -32.97448 \\ 18.403 & 12.02697 \\ -24.960 & 39.71091 \\ 69.388 & -18.76340 \end{bmatrix}$$

Stress Majorization

In the classical MDS algorithm the cost function that we are trying to optimize is called the stress function and it is given by :

$$Stress_D(x_1, x_2, \dots, x_N) = \left(\sum_{i \neq j=1, \dots, N} (d_{ij} - \|x_i - x_j\|)^2 \right)^{1/2}$$

In this function we try to find x_1, \dots, x_N in a certain dimension d such that $Stress(x_1, \dots, x_N)$ is as small as possible

Stress Majorization

The stress function has a more general form as :

$$\sigma(X) = \sum_{i < j \leq n} w_{ij} (d_{ij}(X) - \delta_{ij})^2$$

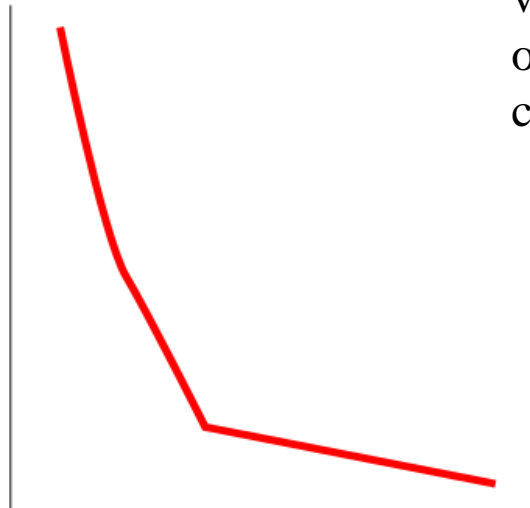
Here w_{ij} weight between a pair of points (i,j) that represents the confidence in the similarity between points (i,j) .

δ_{ij} the given distance between the points i,j

Stress Majorization

“Pressing” the data into 2 dimensions enables us to visualize the data. However, that comes with a price : high stress function value (which correlates with distorted representation)

stress

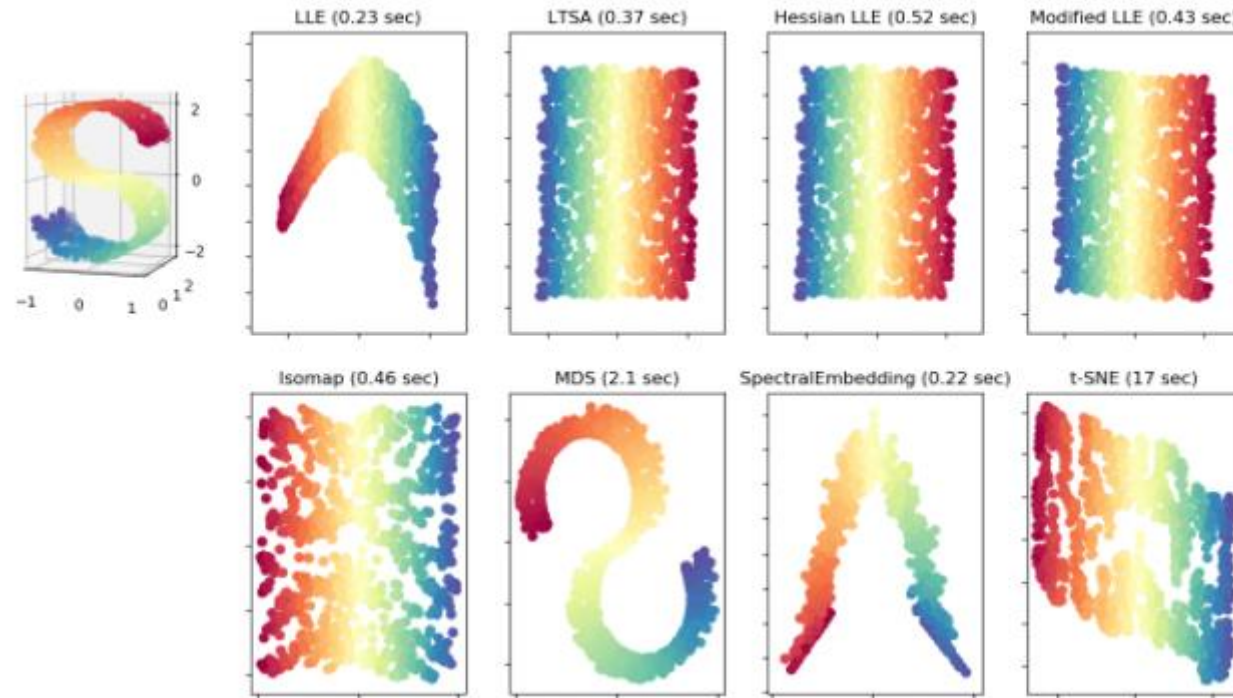


number of
dimensions

Mathematically non-zero stress values may occur only for only one reason: dimensionality of the chosen MDS projection is too low.

MDS in Sklearn

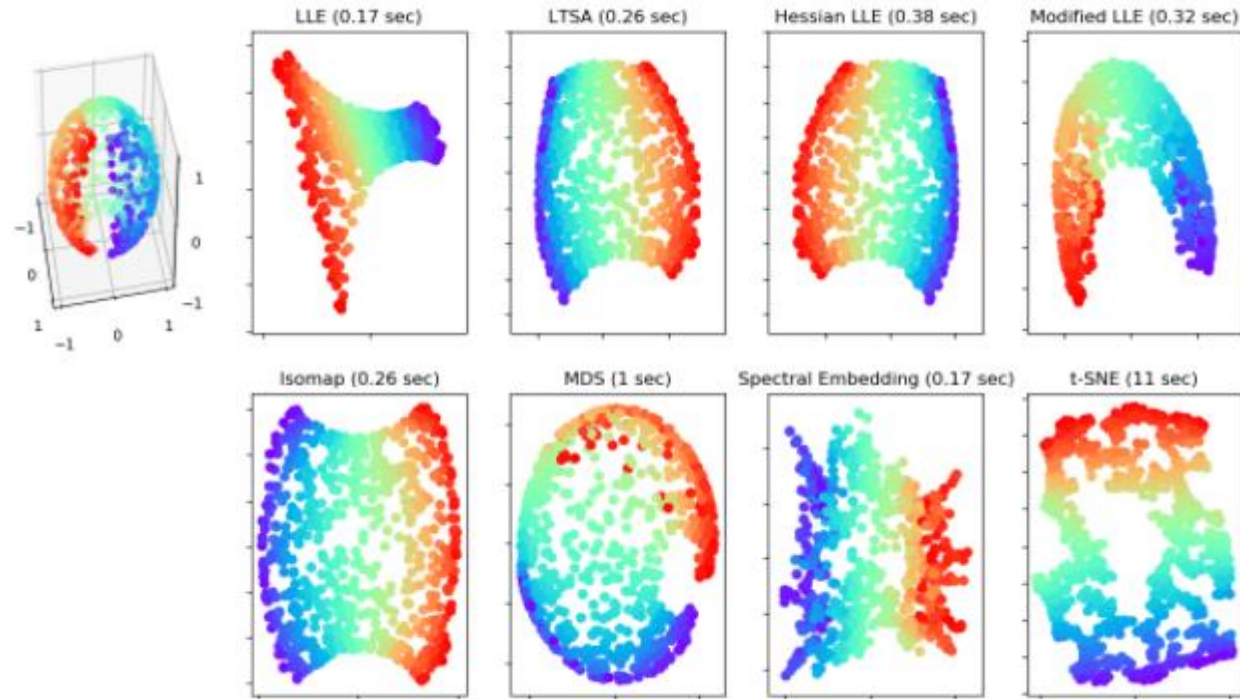
MDS is implemented in [Sklearn](#)



[Example](#)

MDS in Sklearn

MDS is implemented in [Sklearn](#)



[Example](#)

MDS in Sklearn

MDS is implemented in [Sklearn](#)

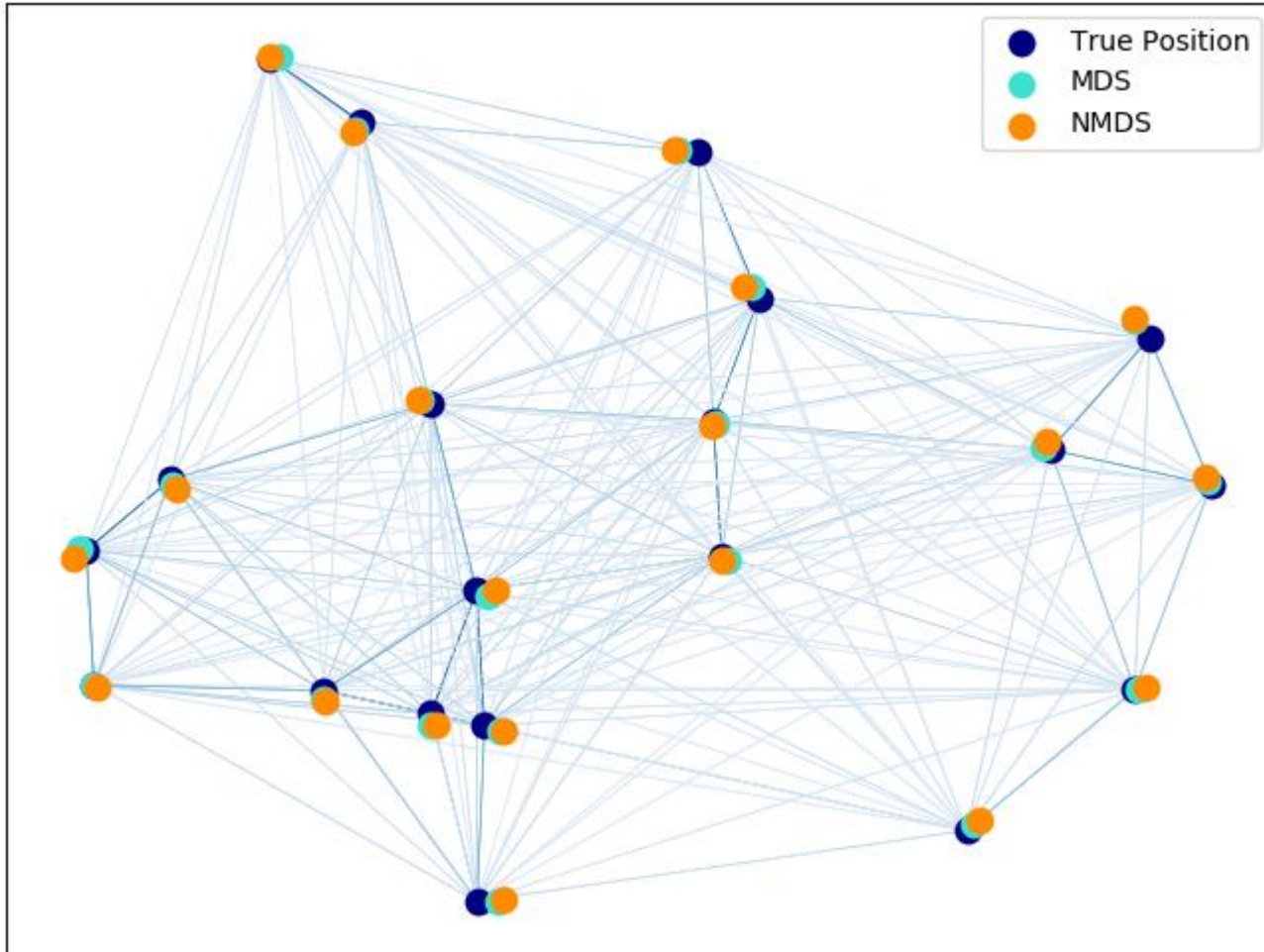
A selection from the 64-dimensional digits dataset

0	1	2	3	4	5	0	1	2	3	4	5	0	1	2	3	4	5	0	5
5	5	0	4	1	3	5	1	0	0	2	2	2	0	1	2	3	3	3	3
4	4	1	5	0	5	2	2	0	0	1	3	2	1	4	3	1	3	1	4
3	1	4	0	5	3	1	5	4	4	2	2	2	5	5	4	4	0	0	1
2	3	4	5	0	1	2	3	4	5	0	1	2	3	4	5	0	5	5	5
0	4	1	3	5	1	0	0	2	2	2	0	1	2	3	3	3	3	4	4
1	5	0	5	2	2	0	0	1	3	2	1	3	1	3	1	4	3	1	4
0	5	3	4	5	4	4	1	2	2	5	5	4	4	0	0	1	2	3	4
5	0	1	2	3	4	5	0	1	2	3	4	5	0	5	5	5	0	4	1
3	5	1	0	0	2	2	2	0	1	2	3	3	3	3	4	4	1	5	0
5	2	2	0	0	1	3	2	1	4	3	1	3	1	4	3	1	4	0	5
3	1	5	4	4	2	2	2	5	5	4	4	0	3	0	1	2	3	4	5
0	1	2	3	4	5	0	1	2	3	4	5	0	5	5	5	0	4	1	3
5	1	0	0	2	2	2	0	1	2	3	3	3	3	4	4	1	5	0	5
2	2	0	0	1	3	2	1	4	3	1	3	1	4	3	1	4	0	5	3
1	5	4	4	2	2	2	5	5	4	4	0	0	1	2	3	4	5	0	1
2	3	4	5	0	1	2	3	4	5	0	5	5	5	0	4	1	3	5	1
0	0	2	2	2	0	1	1	3	3	3	3	4	4	1	5	0	5	2	2
0	0	1	3	2	1	4	3	1	3	1	4	3	1	4	0	5	3	1	5
4	4	2	2	1	5	5	4	4	0	0	1	2	3	4	5	0	1	2	3

[Example](#)

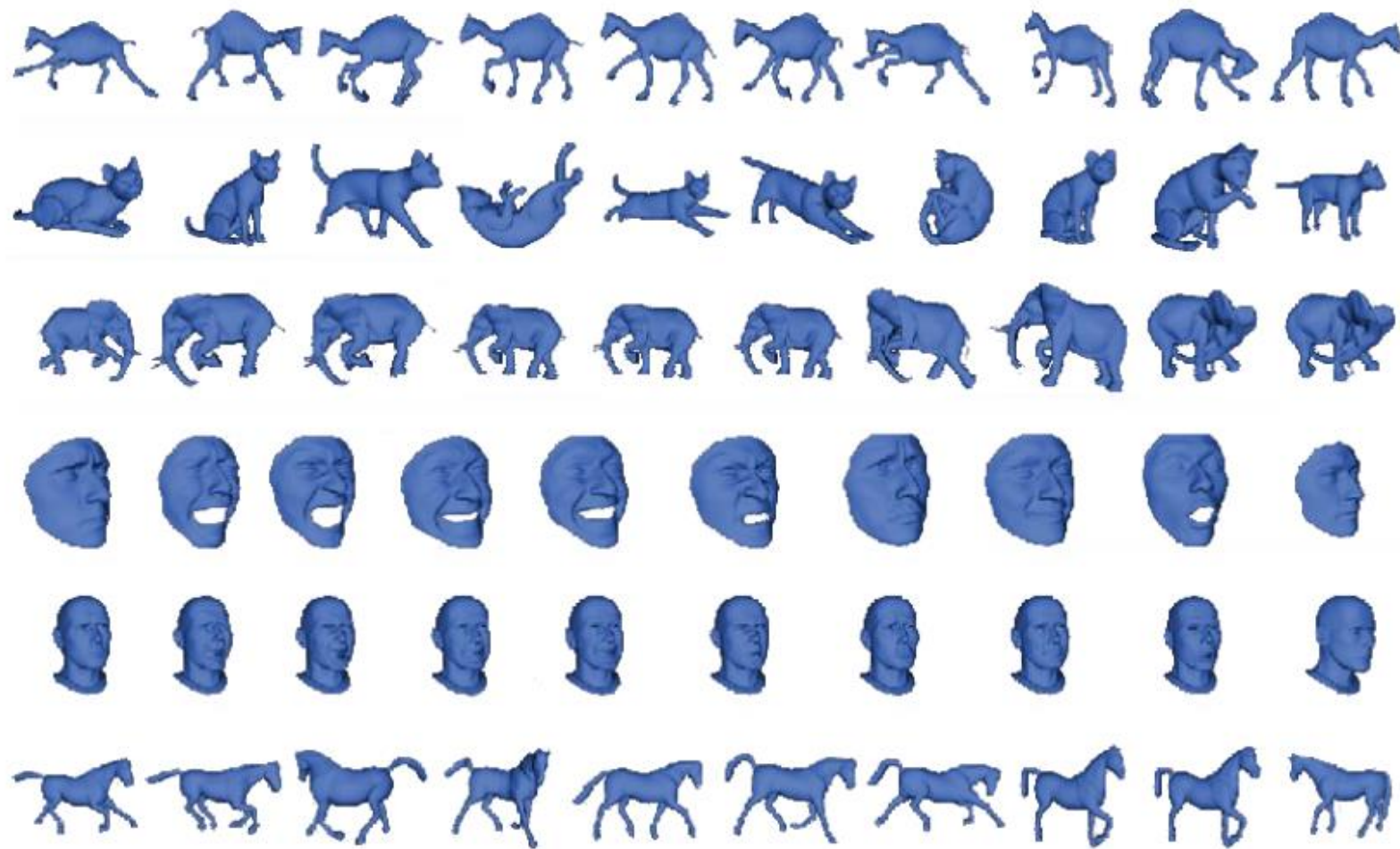
MDS in Sklearn

MDS is implemented in [Sklearn](#)

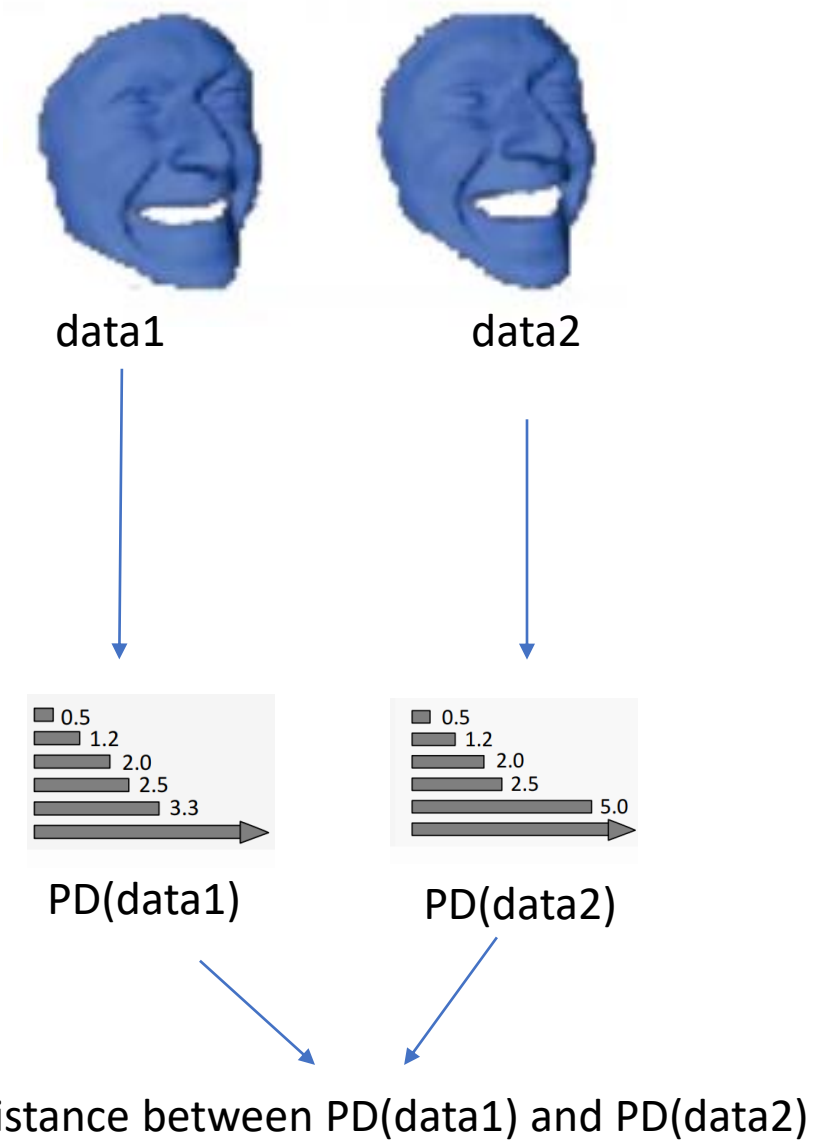
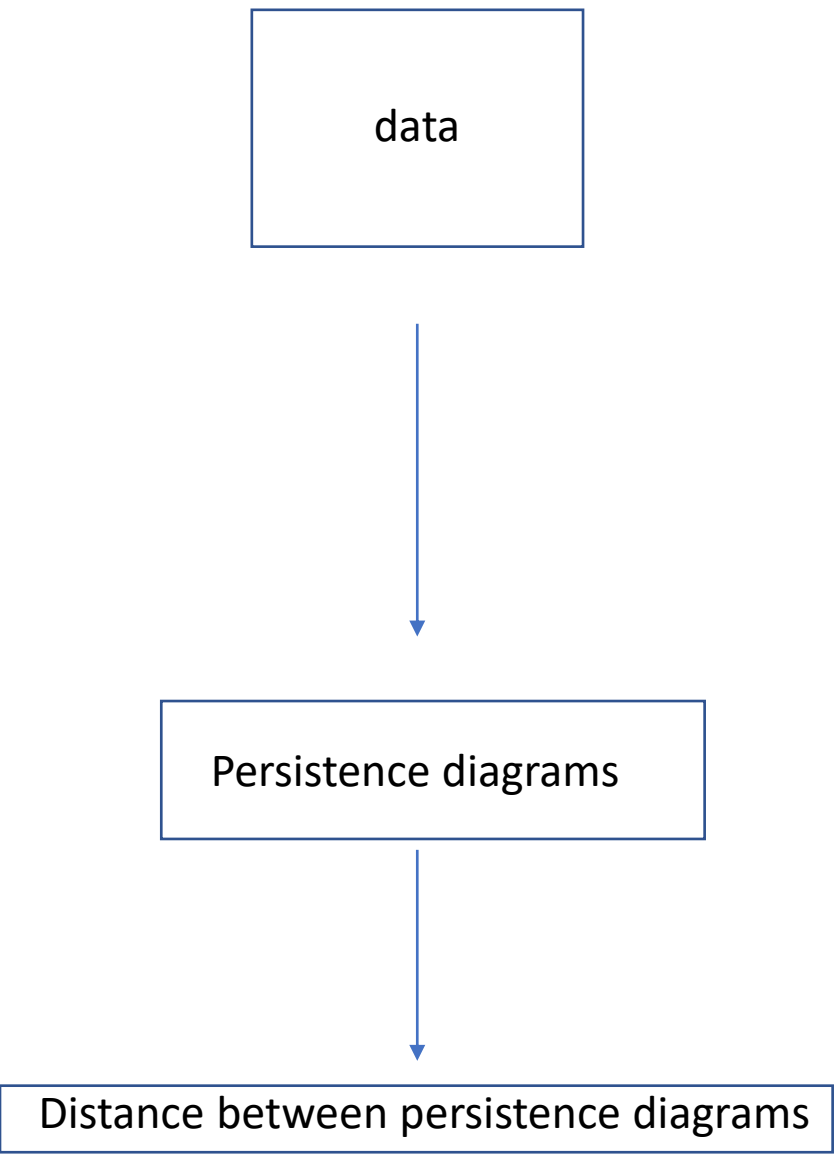


[Example](#)

Application



Measuring distance between two persistence diagrams



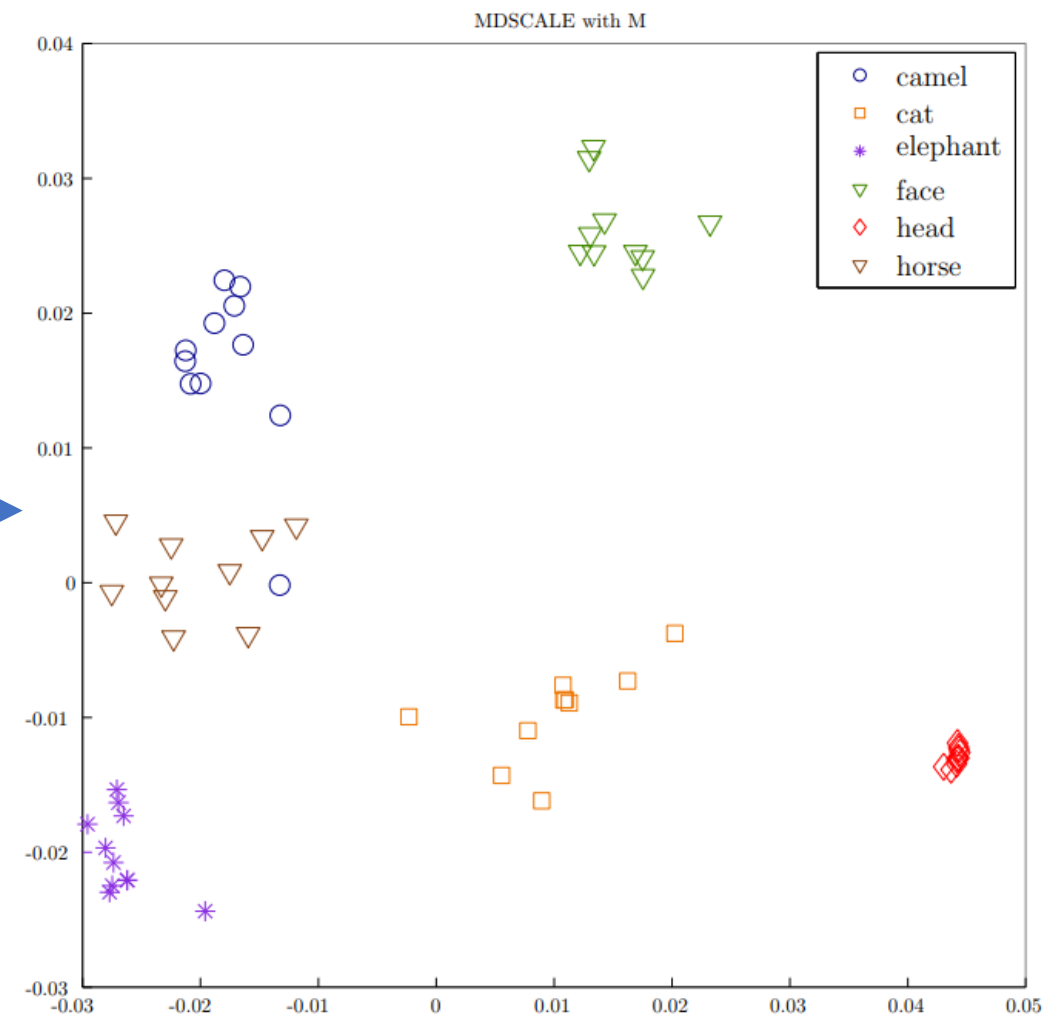
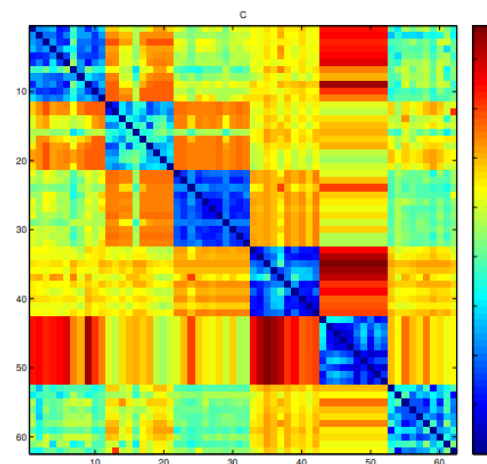
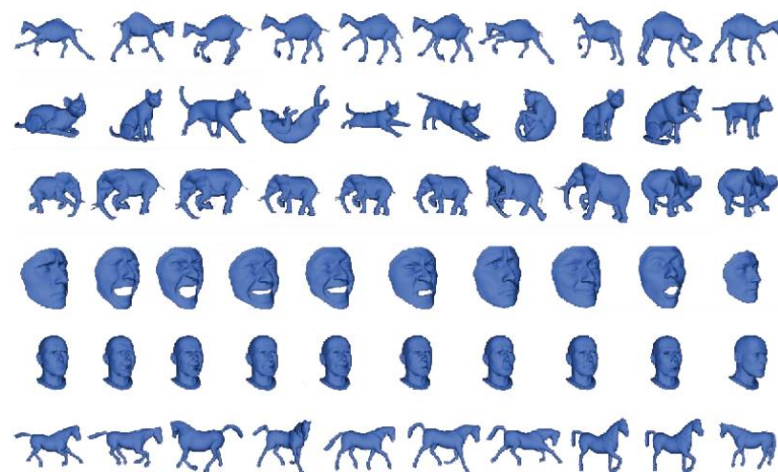
Bottleneck distance between two persistent diagrams

- Given two persistence diagrams X and Y , let η be a bijection between points in the diagram. The following two distances are commonly used in the context of PH to measure the distance between two persistence diagrams:

$$W_{\infty}(X, Y) = \inf_{\eta: X \rightarrow Y} \sup_{x \in X} \|x - \eta(x)\|_{\infty}$$

$$W_q(X, Y) = \left[\inf_{\eta: X \rightarrow Y} \sum_{x \in X} \|x - \eta(x)\|_{\infty}^q \right]^{1/q}$$

Bottleneck distance between two persistent diagrams



Input data



Matrix M describes the
pair-wise distance
between the persistence
diagrams of each data
element



MDS plot of the matrix M with labels corresponding
to each class.

Remarks

- The axes obtained when drawing the MDS coordinates are , in themselves, meaningless.

Remarks

- The axes obtained when drawing the MDS coordinates are , in themselves, meaningless.
- The orientation of the result MDS “picture” is arbitrary.

Remarks

- The axes obtained when drawing the MDS coordinates are , in themselves, meaningless.
- The orientation of the result MDS “picture” is arbitrary.
- When we obtain MDS coordinates that have non-zero stress, we should remember that the distances among the resulting items are distorted representations of the relationships given by the input data. This distortion is greater when the stress is greater.
- That being said, we, in general, can rely on the larger distances as being more accurate than smaller distances.

Non-Matric MDS

- Sometimes, there is no defined metric on points and all we are given is a similarity measure between the points.

Non-Matric MDS

- Sometimes, there is no defined metric on points and all we are given is a similarity measure between the points.

The main idea in non-metric MDS :

- The actual values given to us are not that meaningful
- *Ranking among different points* is important
- Non-metric MDS finds a low-dimensional representation, which respects the ranking of distances as much as possible

Non-Matric MDS

- Recall that in MDS we seek to find an optimal configuration x_i that *gives* $d_{ij} \approx d'_{ij} = ||x_i - x_j||$ as close as possible.

Non-Metric MDS

- Recall that in MDS we seek to find an optimal configuration x_i that *gives* $d_{ij} \approx d'_{ij} = ||x_i - x_j||$ as close as possible.
- Relaxing $d_{ij} \approx d'_{ij}$ from MDS by allowing $d'_{ij} \approx f(d_{ij})$, for some monotone function f

Monotonic means : $d_{ij} < d_{kl} \Leftrightarrow f(d_{ij}) \leq f(d_{kl})$

Non-Matric MDS

- Recall that in MDS we seek to find an optimal configuration x_i that gives $d_{ij} \approx d'_{ij} = ||x_i - x_j||$ as close as possible.
- Relaxing $d_{ij} \approx d'_{ij}$ from MDS by allowing $d'_{ij} \approx f(d_{ij})$, for some monotone function f

Monotonic means : $d_{ij} < d_{kl} \Leftrightarrow f(d_{ij}) \leq f(d_{kl})$

Given a dimension d , non-metric MDS seeks to find an optimal configuration $X \subset R^d$ that gives $f(d_{ij}) \approx \hat{d}_{ij} = ||x_i - x_j||$ as close as possible.

- $f(d_{ij}) = d^*_{ij}$ is only required to preserve the order of d_{ij} ,

i.e., $d_{ij} < d_{kl} \Leftrightarrow f(d_{ij}) \leq f(d_{kl}) \Leftrightarrow d^*_{ij} \leq d^*_{kl}$

Non-Metric MDS

The stress function for non-metric MDS is given by :

$$Stress = \left(\sum_{i < j} (\hat{d}_{ij} - f(d_{ij}))^2 / \sum d_{ij}^2 \right)^{\frac{1}{2}}$$

Non-metric MDS optimizes over both position of the points of points and f

Non-Metric MDS

The stress function for non-metric MDS is given by :

$$Stress = \left(\sum_{i < j} (\hat{d}_{ij} - f(d_{ij}))^2 / \sum d_{ij}^2 \right)^{\frac{1}{2}}$$

Non-metric MDS optimizes over both position of the points and f

Solved numerically using ([isotonic regression](#)); we usually use classical MDS as starting initial position.

Non-Metric MDS

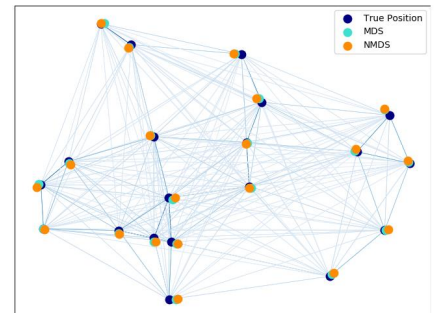
The stress function for non-metric MDS is given by :

$$Stress = \left(\sum_{i < j} (\hat{d}_{ij} - f(d_{ij}))^2 / \sum d_{ij}^2 \right)^{\frac{1}{2}}$$

Non-metric MDS optimizes over both position of the points and f

Solved numerically using ([isotonic regression](#)); we usually use classical MDS as starting initial position.

[Example](#)



ISOMAP

- Isomap extends MDS by utilizing geodesic distances induced by some neighborhood graph.
- There is no essential difference between MDS and ISOMAP algorithm once we find the distance matrix induced by the neighborhood graph.

ISOMAP

- Isomap extends MDS by utilizing geodesic distances induced by some neighborhood graph.
- There is no essential difference between MDS and ISOMAP algorithm once we find the distance matrix induced by the neighborhood graph.

So the steps for ISOMAP on a given data :

1- Construct the neighborhood graph of the data X using one of the neighborhood graphs we studied earlier in the course

ISOMAP

- Isomap extends MDS by utilizing geodesic distances induced by some neighborhood graph.
- There is no essential difference between MDS and ISOMAP algorithm once we find the distance matrix induced by the neighborhood graph.

So the steps for ISOMAP on a given data :

- 1- Construct the neighborhood graph of the data X using one of the neighborhood graphs we studied earlier in the course
- 2- Use the Dijkstra algorithm or the Floyd–Warshall algorithm to find the distance between nodes on the graph

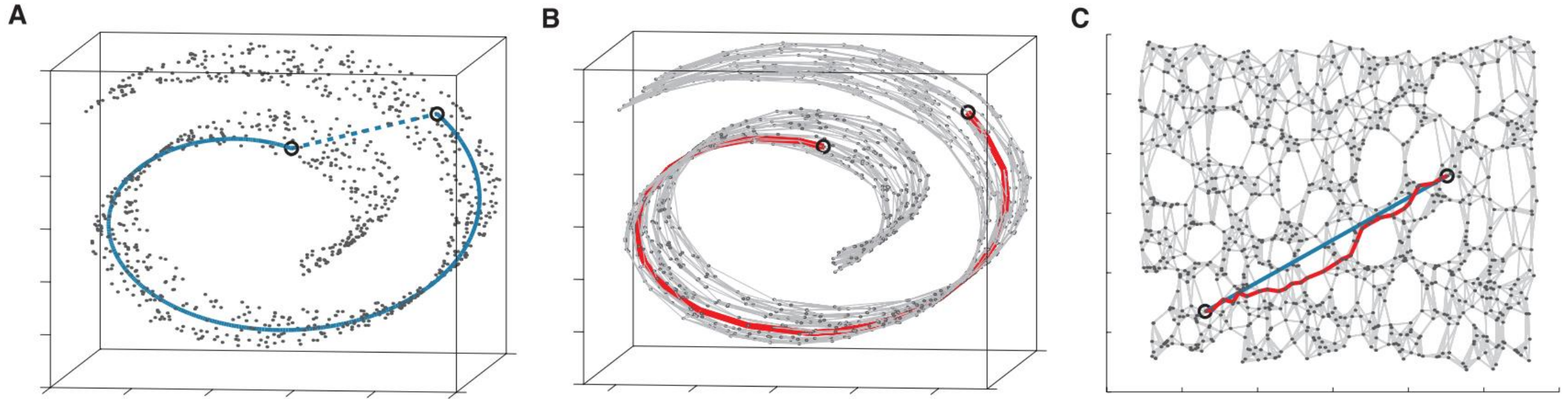
ISOMAP

- Isomap extends MDS by utilizing geodesic distances induced by some neighborhood graph.
- There is no essential difference between MDS and ISOMAP algorithm once we find the distance matrix induced by the neighborhood graph.

So the steps for ISOMAP on a given data :

- 1- Construct the neighborhood graph of the data X using one of the neighborhood graphs we studied earlier in the course
- 2- Use the Dijkstra algorithm or the Floyd–Warshall algorithm to find the distance between nodes on the graph
- 3- Apply MDS on the distance matrix above and extract the coordinates with the desired dimension

ISOMAP



A- Euclidian distance might not represent the actual distance between the points in the data.

B- We can construct the neighborhood graph of the data and then compute the geodesic distance between the points of the graph

C- Embedding the space we obtained in B into the plane.

Appendix : Floyd–Warshall algorithm

```
let dist be a  $|V| \times |V|$  array of minimum distances  
initialized to infinity  
  
for each edge  $(u, v)$   
     $\text{dist}[u][v] \leftarrow w(u, v)$  // the weight of the edge  $(u, v)$   
for each vertex  $v$   
     $\text{dist}[v][v] \leftarrow 0$   
for  $k$  from 1 to  $|V|$   
    for  $i$  from 1 to  $|V|$   
        for  $j$  from 1 to  $|V|$   
            if  $\text{dist}[i][j] > \text{dist}[i][k] + \text{dist}[k][j]$   
                 $\text{dist}[i][j] \leftarrow \text{dist}[i][k] + \text{dist}[k][j]$   
            end if
```

Floyd algorithm is good to use when we want to compute the distance matrix on a dense graph.
When the graph G is sparse, Dijkstra algorithm is a better choice.