

Clustering II

Lecture 15

Hierarchical Clustering

agglomerative (bottom-up) clustering

divisive (top-down) clustering

Agglomerative clustering components

Distance metric

How we measure distance/dissimilarity

Euclidean distance
(L_2 norm) $D(\mathbf{a}, \mathbf{b}) = \|\mathbf{a} - \mathbf{b}\|_2$

Squared Euclidean
distance $D(\mathbf{a}, \mathbf{b}) = \|\mathbf{a} - \mathbf{b}\|_2^2$

Manhattan distance
(L_1 norm) $D(\mathbf{a}, \mathbf{b}) = \|\mathbf{a} - \mathbf{b}\|_1$

Maximum distance $D(\mathbf{a}, \mathbf{b}) = \|\mathbf{a} - \mathbf{b}\|_\infty$
 $= \max_i |a_i - b_i|$

Linkage criterion

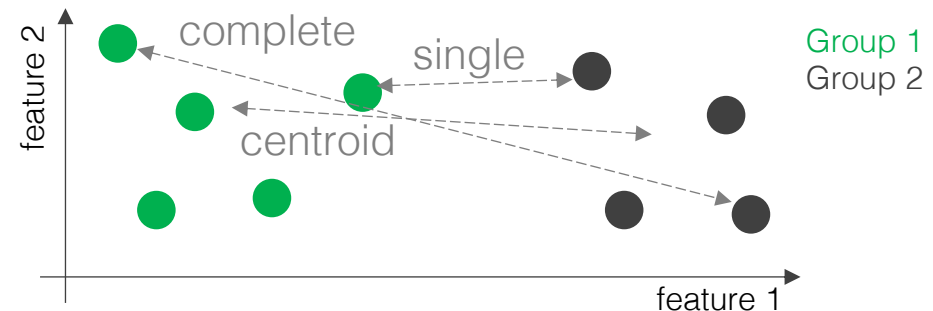
How to measure distance/dissimilarity
between groups or sets

Complete = maximum intercluster dissimilarity

Single = minimum intercluster dissimilarity

Average = average intercluster dissimilarity (calculate the dissimilarity between all pairs of points, take the average)

Centroid = dissimilarity between cluster centroids



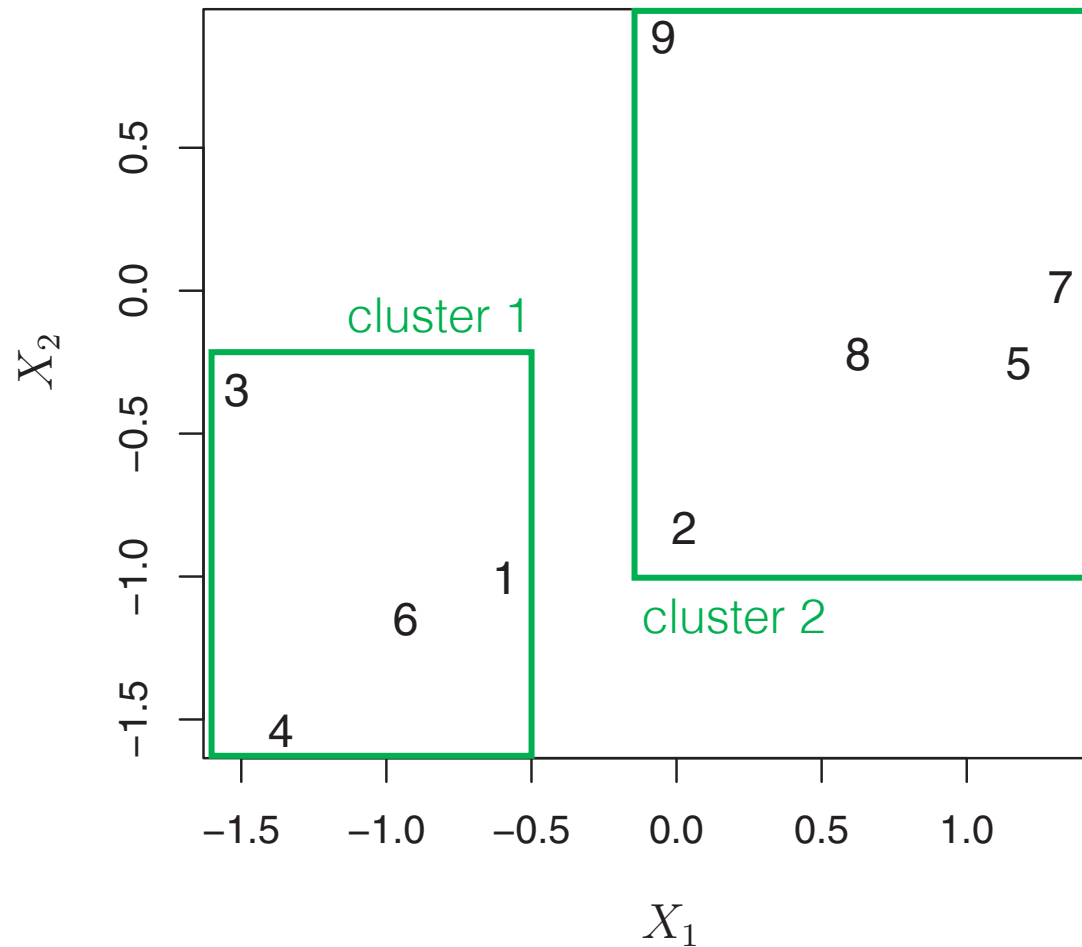
Agglomerative clustering

With complete linkage and Euclidean distance

Algorithm:

1. Select a measure of dissimilarity and linkage
2. Set each observation as a unique cluster
3. Group the two closest clusters together
4. Repeat until there is only one cluster

Data in 2-D feature space



Dendrogram

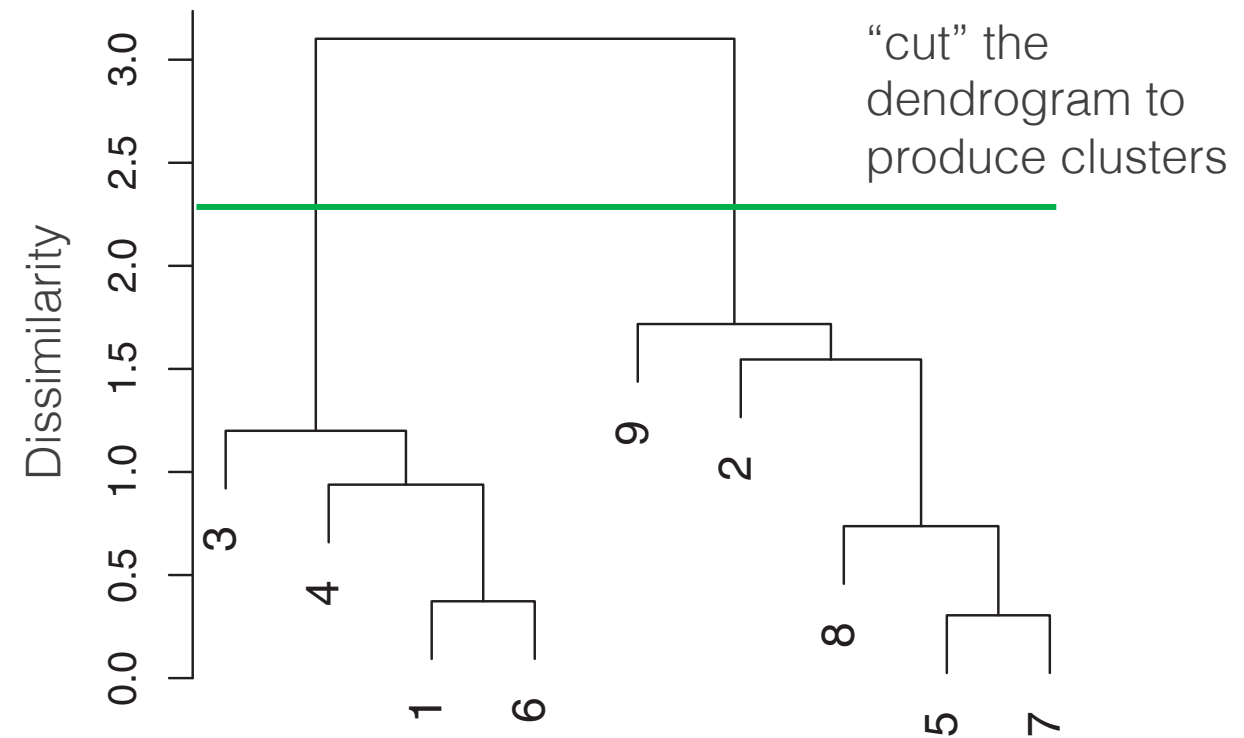


Image from James et al., Introduction to Statistical Learning, 2013

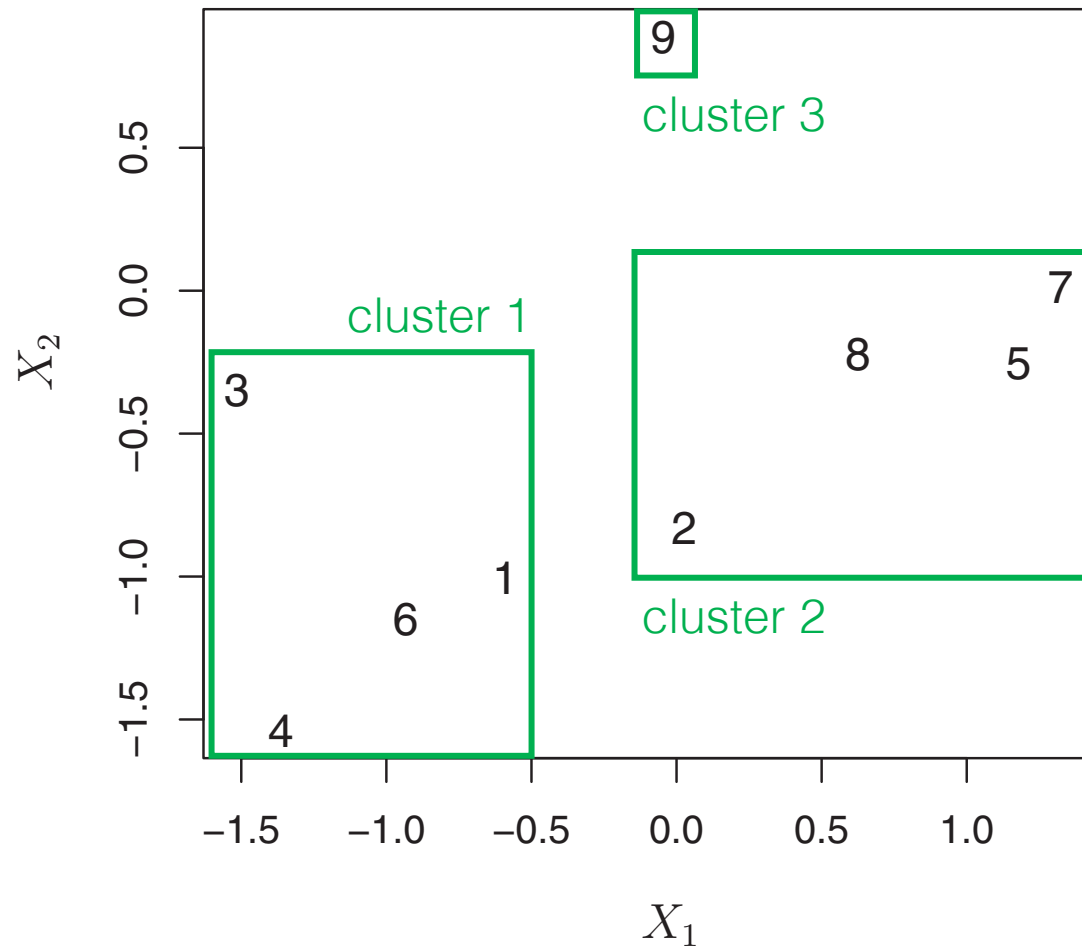
Agglomerative clustering

With complete linkage and Euclidean distance

Algorithm:

1. Select a measure of dissimilarity and linkage
2. Set each observation as a unique cluster
3. Group the two closest clusters together
4. Repeat until there is only one cluster

Data in 2-D feature space



Dendrogram

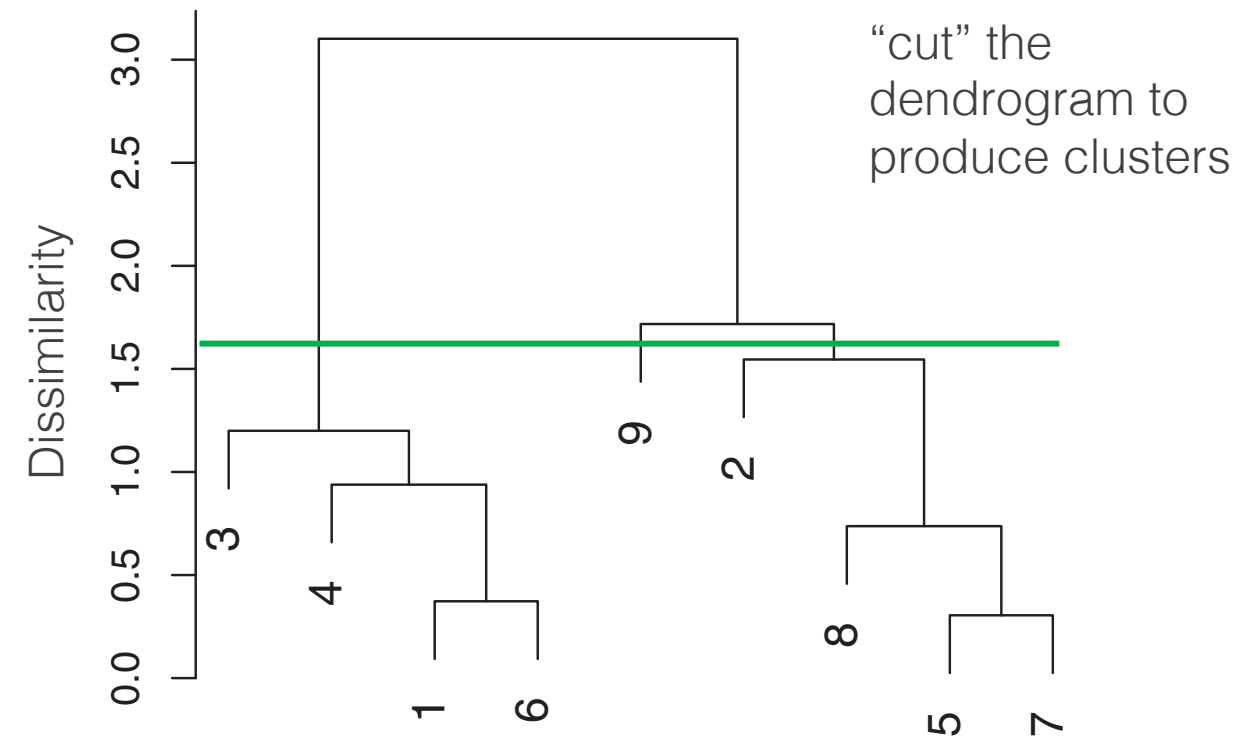


Image from James et al., Introduction to Statistical Learning, 2013

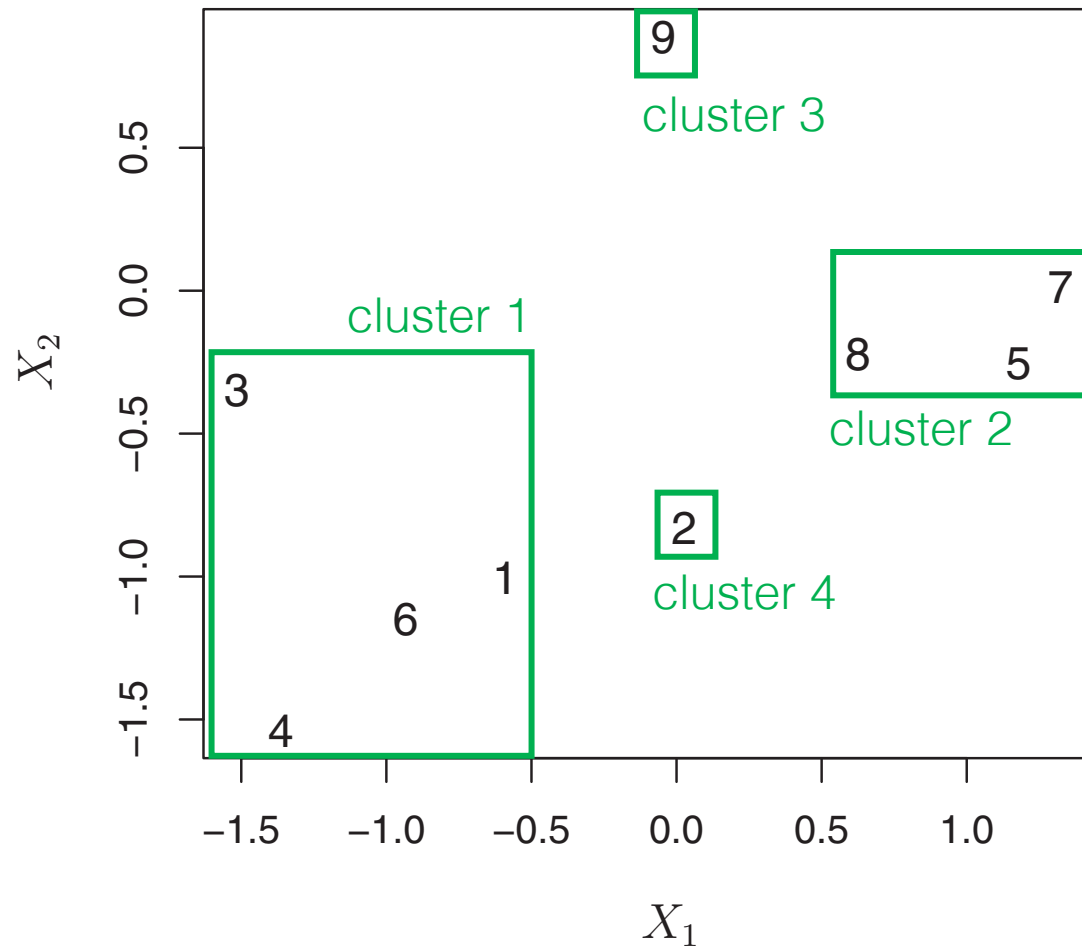
Agglomerative clustering

With complete linkage and Euclidean distance

Algorithm:

1. Select a measure of dissimilarity and linkage
2. Set each observation as a unique cluster
3. Group the two closest clusters together
4. Repeat until there is only one cluster

Data in 2-D feature space



Dendrogram

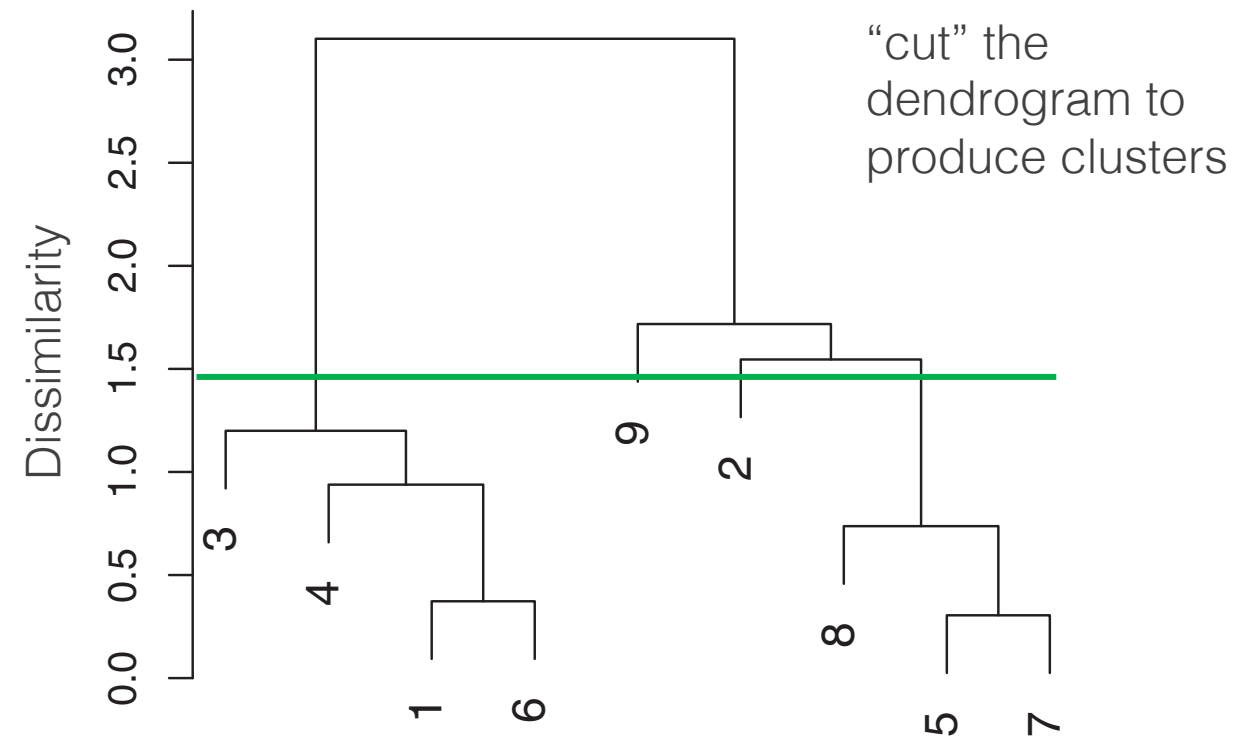
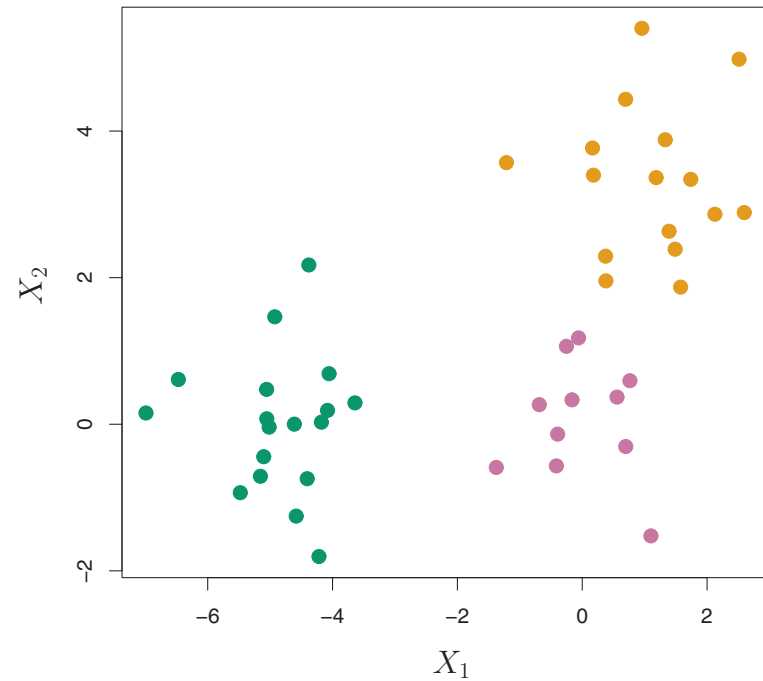


Image from James et al., Introduction to Statistical Learning, 2013

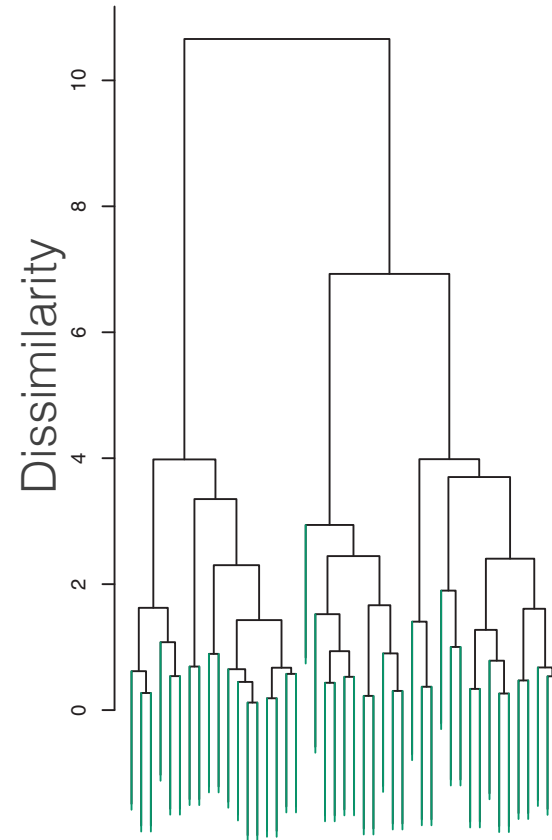
Example of agglomerative clustering

With complete linkage and Euclidean distance

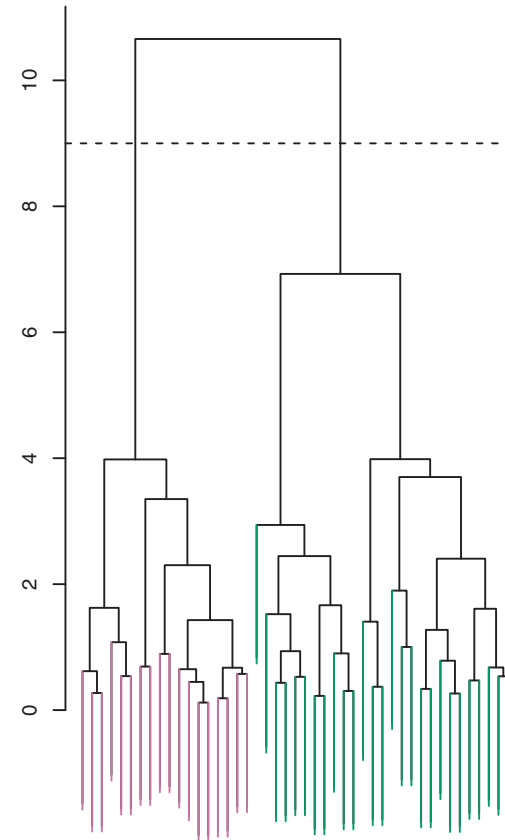
Data in 2-D feature space



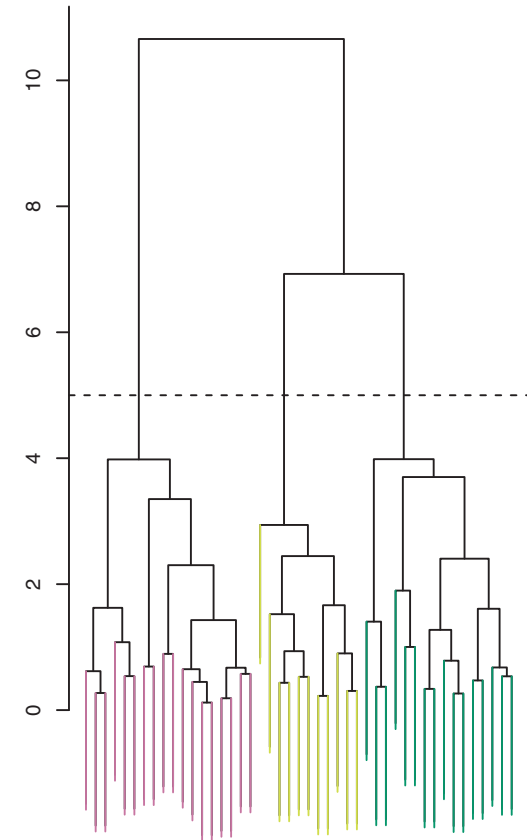
Original Dendrogram



Dendrogram cut for 2 clusters



Dendrogram cut for 3 clusters



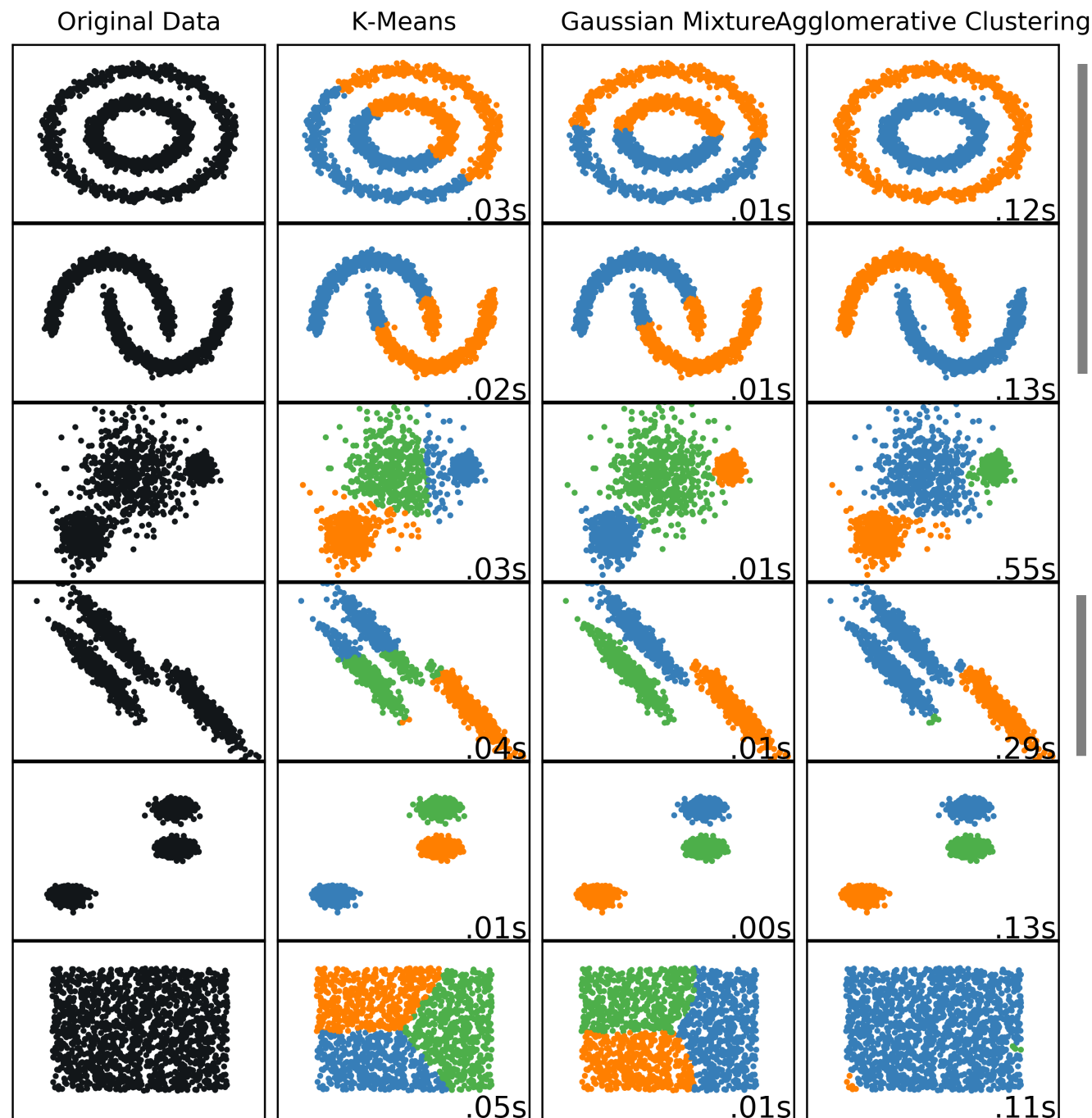
Note: colors do not directly map to plot on the left

Image from James et al., Introduction to Statistical Learning, 2013

Examples: Agglomerative clustering

Need to choose
where to cut the
dendrogram

Can be slow since
all pairwise
distances between
clusters need to
be evaluated



Performs well
when clusters are
well-separated

Struggles when
intercluster
distance is not
sufficient to
distinguish
between clusters

DBSCAN Clustering

Density-based spatial clustering of applications with noise

By Martin Ester, Hans-Peter Kriegel, Jörg Sander, and Xiaowei Xu, 1996

DBSCAN

Parameters:

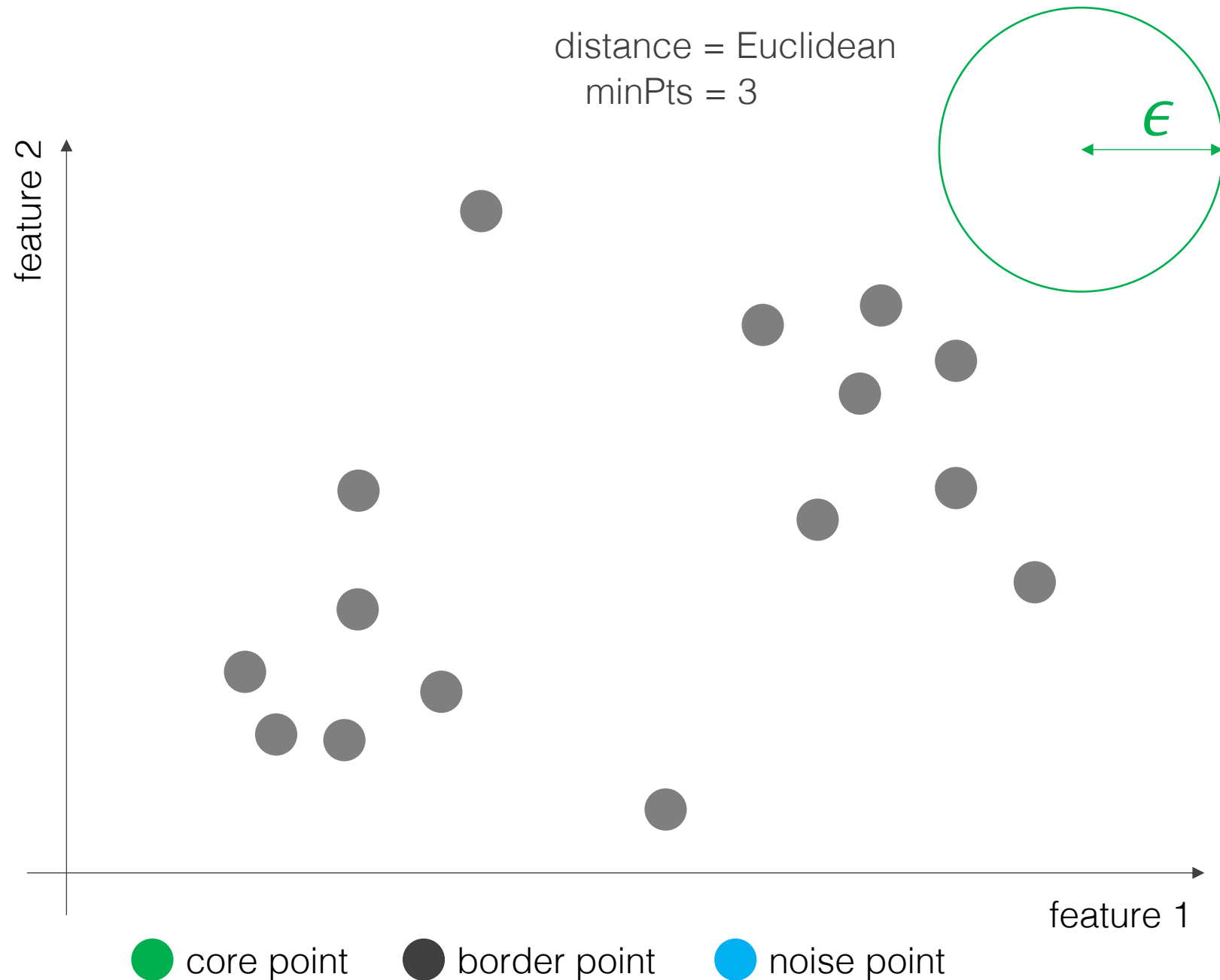
1. Distance measure
2. The radius of a neighbor, ϵ
3. 'minPts': The number of neighbors for a point to be considered a core point

Types of points:

- **Core**: a point with at least minPts neighbors
- **Border**: a non-core point that neighbors a core point
- **Noise**: Other points

Algorithm:

1. Label core and border points
2. Group neighboring core points
3. Add border points that are neighbors of core points



DBSCAN

Parameters:

1. Distance measure
2. The radius of a neighbor, ϵ
3. 'minPts': The number of neighbors for a point to be considered a core point

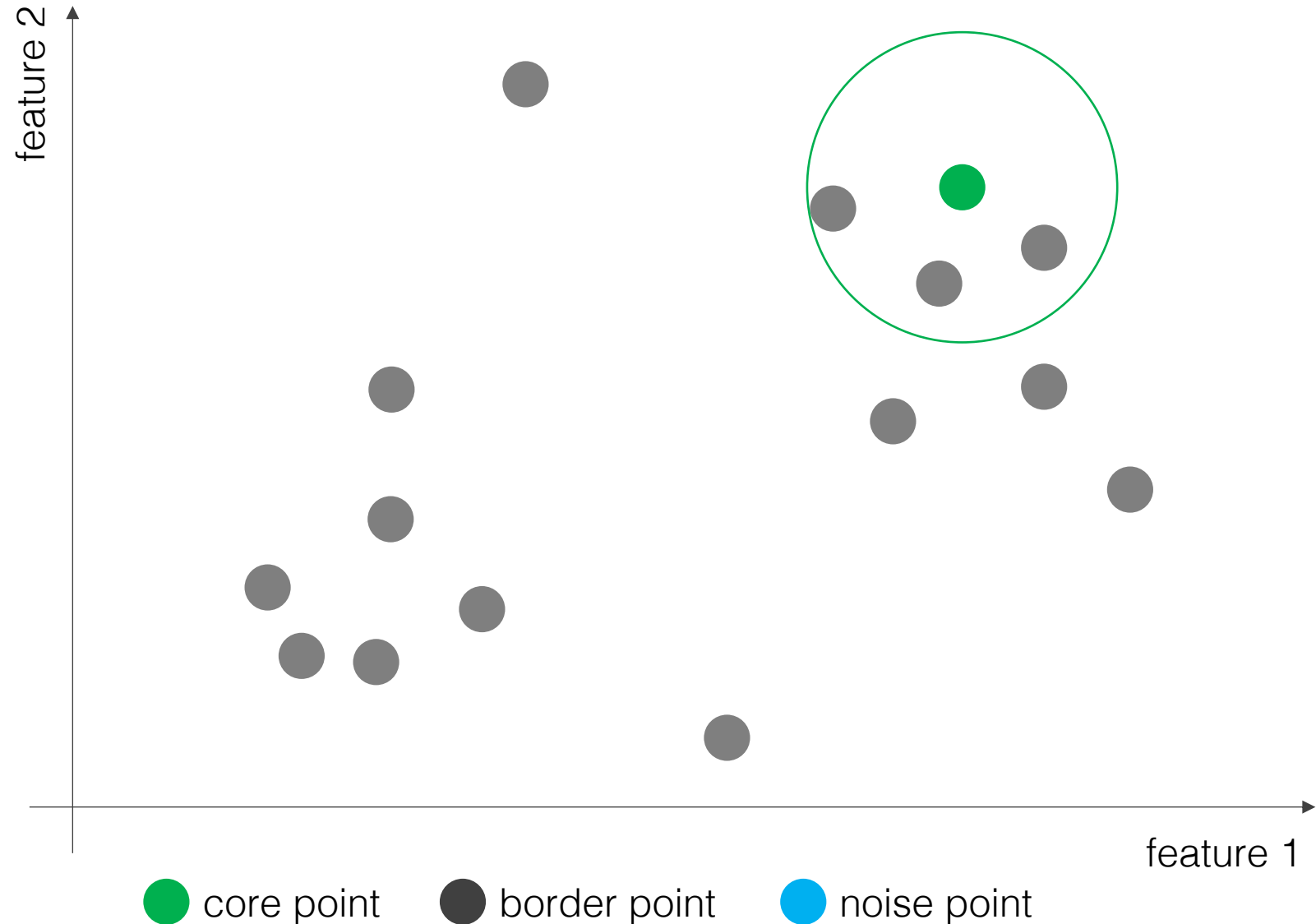
Types of points:

- **Core:** a point with at least minPts neighbors
- **Border:** a non-core point that neighbors a core point
- **Noise:** Other points

Algorithm:

1. Label core and border points
2. Group neighboring core points
3. Add border points that are neighbors of core points

distance = Euclidean
minPts = 3



DBSCAN

Parameters:

1. Distance measure
2. The radius of a neighbor, ϵ
3. 'minPts': The number of neighbors for a point to be considered a core point

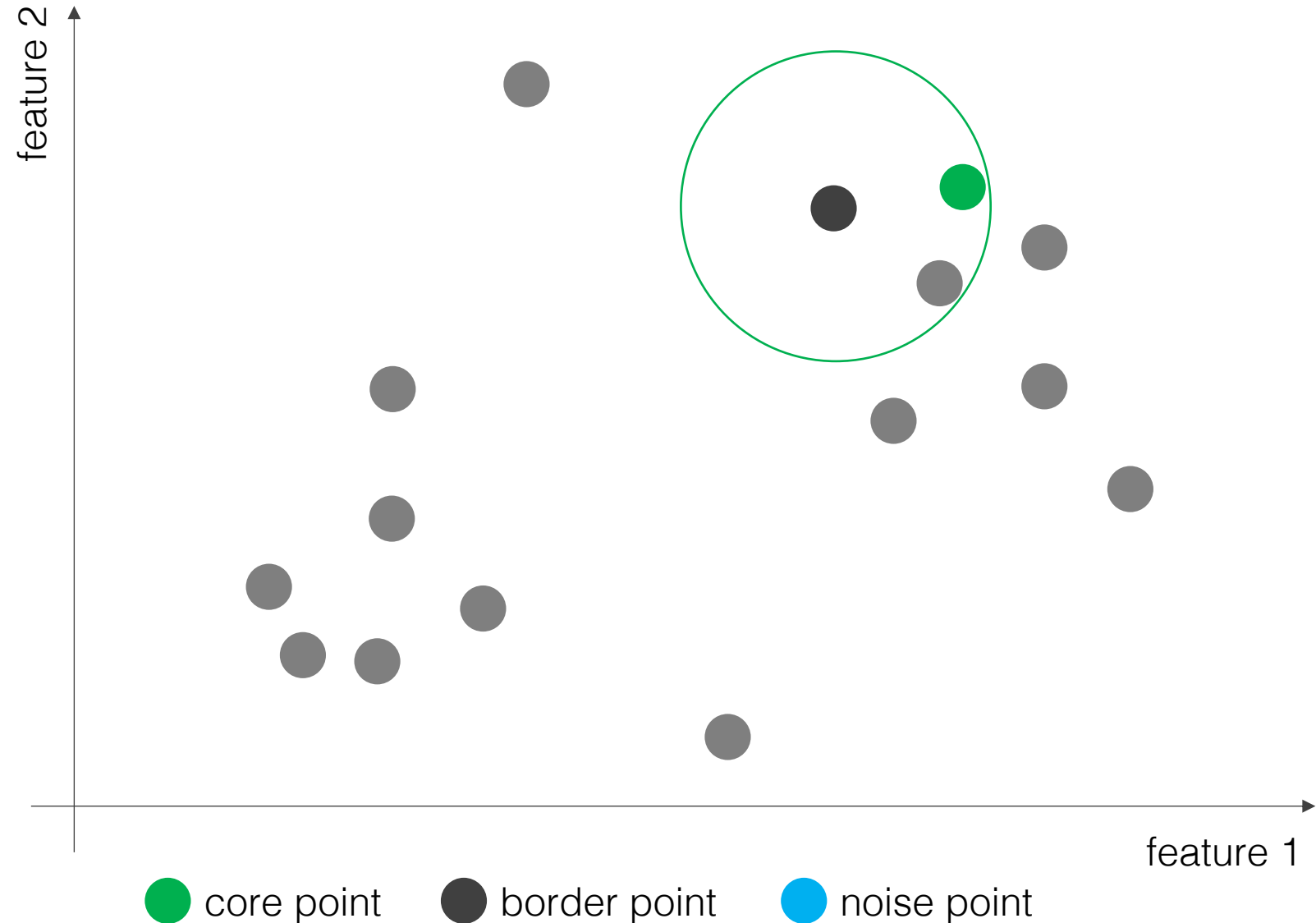
Types of points:

- **Core**: a point with at least minPts neighbors
- **Border**: a non-core point that neighbors a core point
- **Noise**: Other points

Algorithm:

1. Label core and border points
2. Group neighboring core points
3. Add border points that are neighbors of core points

distance = Euclidean
minPts = 3



DBSCAN

Parameters:

1. Distance measure
2. The radius of a neighbor, ϵ
3. 'minPts': The number of neighbors for a point to be considered a core point

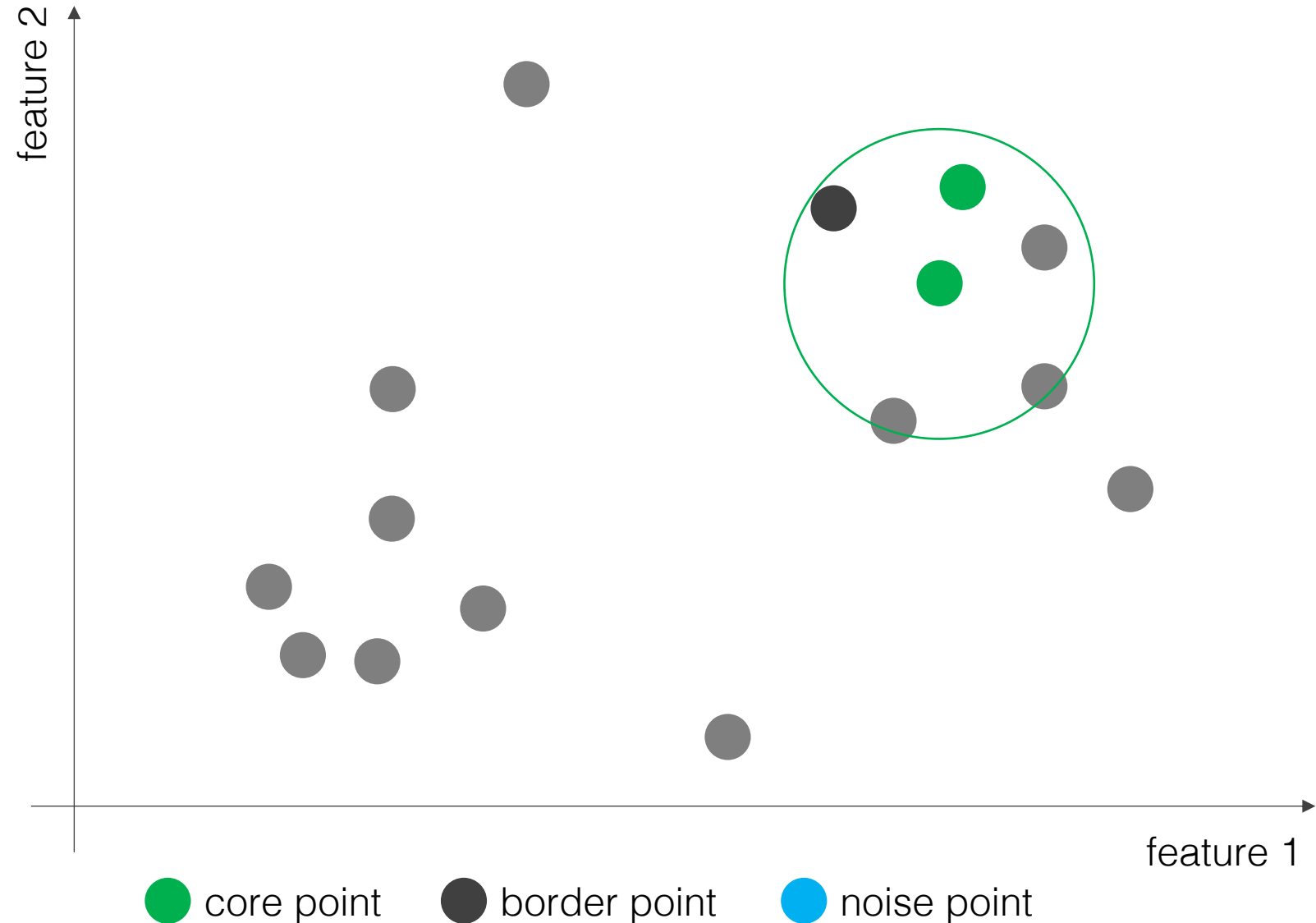
Types of points:

- **Core:** a point with at least minPts neighbors
- **Border:** a non-core point that neighbors a core point
- **Noise:** Other points

Algorithm:

1. Label core and border points
2. Group neighboring core points
3. Add border points that are neighbors of core points

distance = Euclidean
minPts = 3



DBSCAN

Parameters:

1. Distance measure
2. The radius of a neighbor, ϵ
3. 'minPts': The number of neighbors for a point to be considered a core point

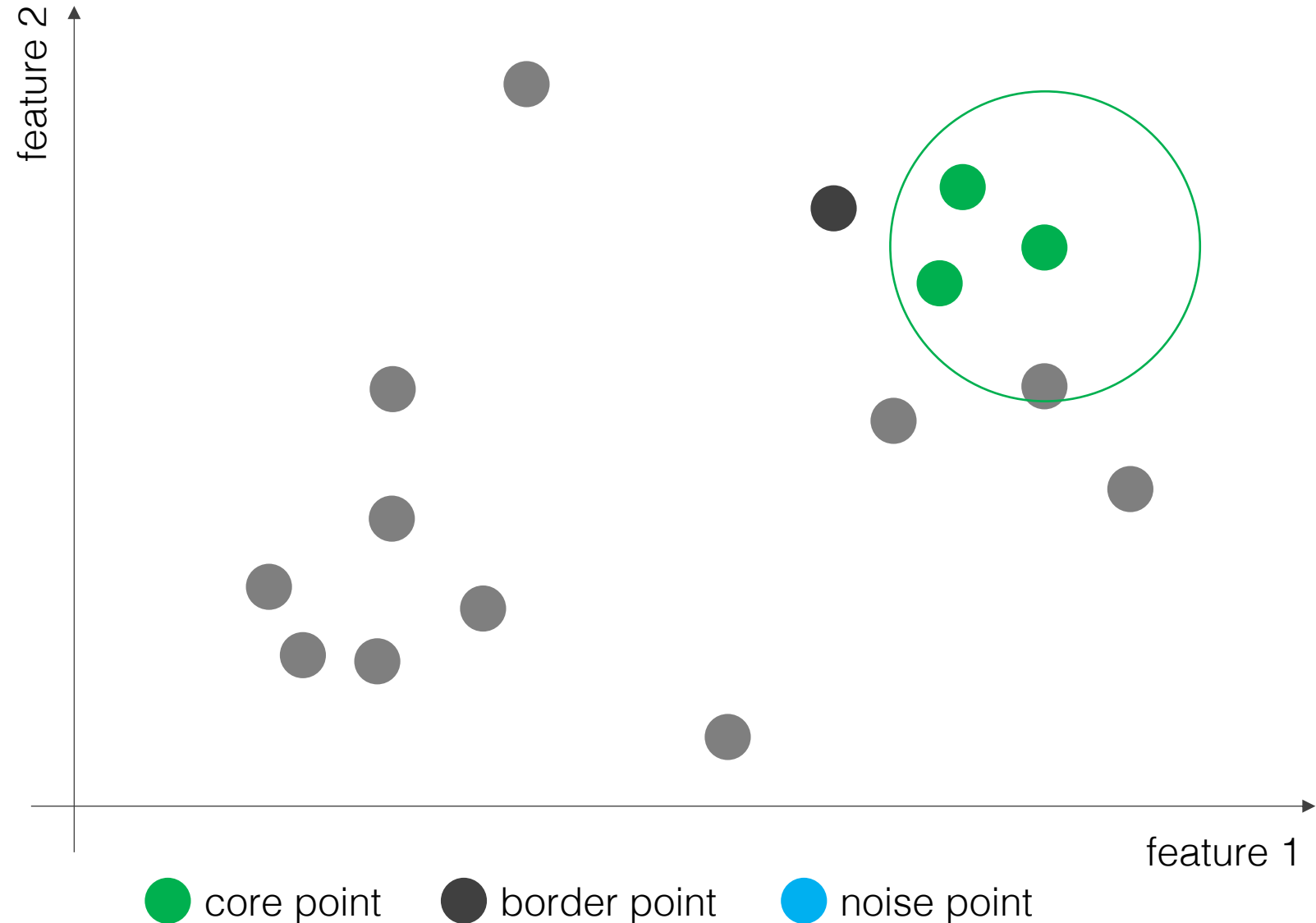
Types of points:

- **Core**: a point with at least minPts neighbors
- **Border**: a non-core point that neighbors a core point
- **Noise**: Other points

Algorithm:

1. Label core and border points
2. Group neighboring core points
3. Add border points that are neighbors of core points

distance = Euclidean
minPts = 3



DBSCAN

Parameters:

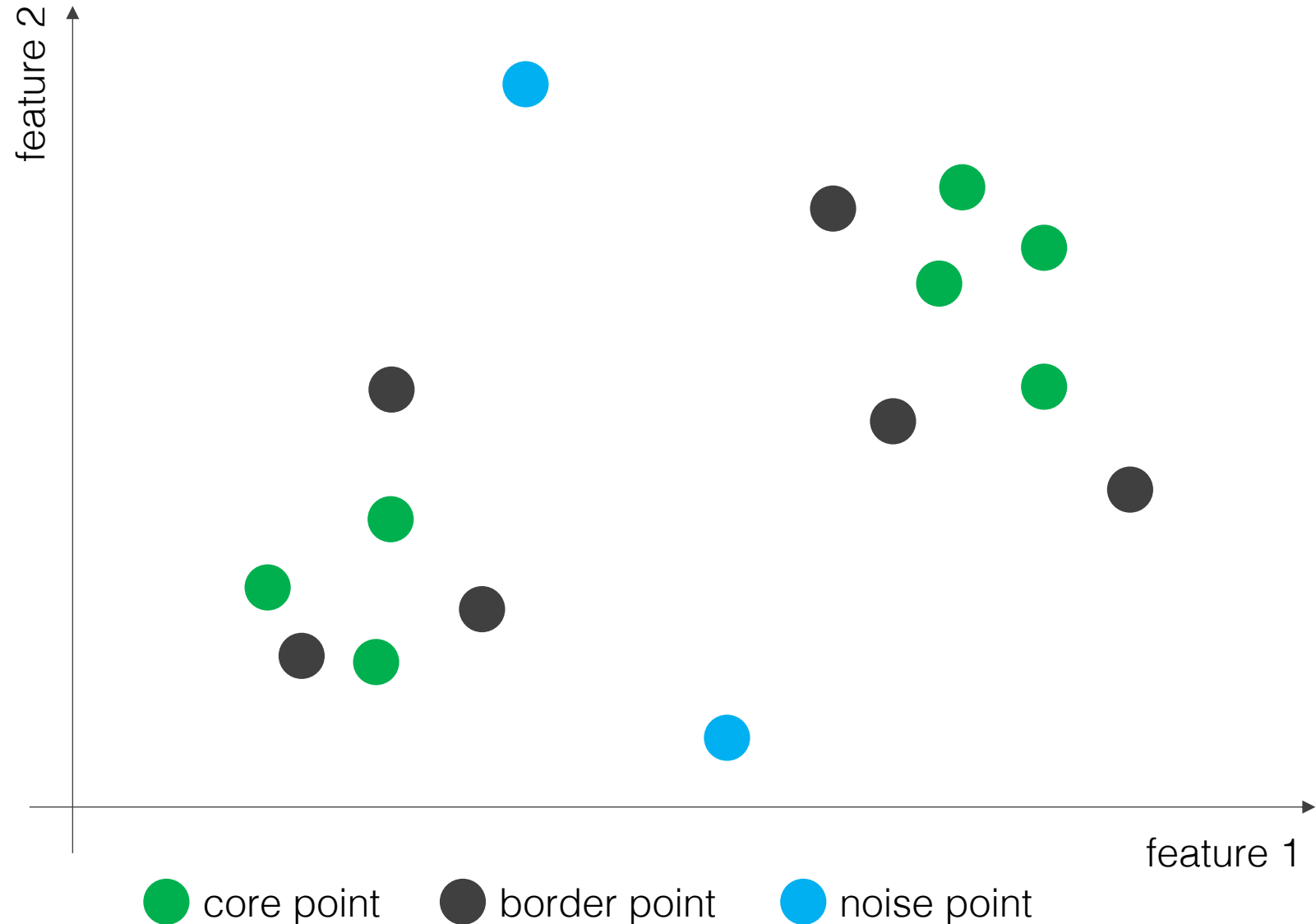
1. Distance measure
2. The radius of a neighbor, ϵ
3. 'minPts': The number of neighbors for a point to be considered a core point

Types of points:

- **Core:** a point with at least minPts neighbors
- **Border:** a non-core point that neighbors a core point
- **Noise:** Other points

Algorithm:

1. Label core and border points
2. Group neighboring core points
3. Add border points that are neighbors of core points



DBSCAN

Parameters:

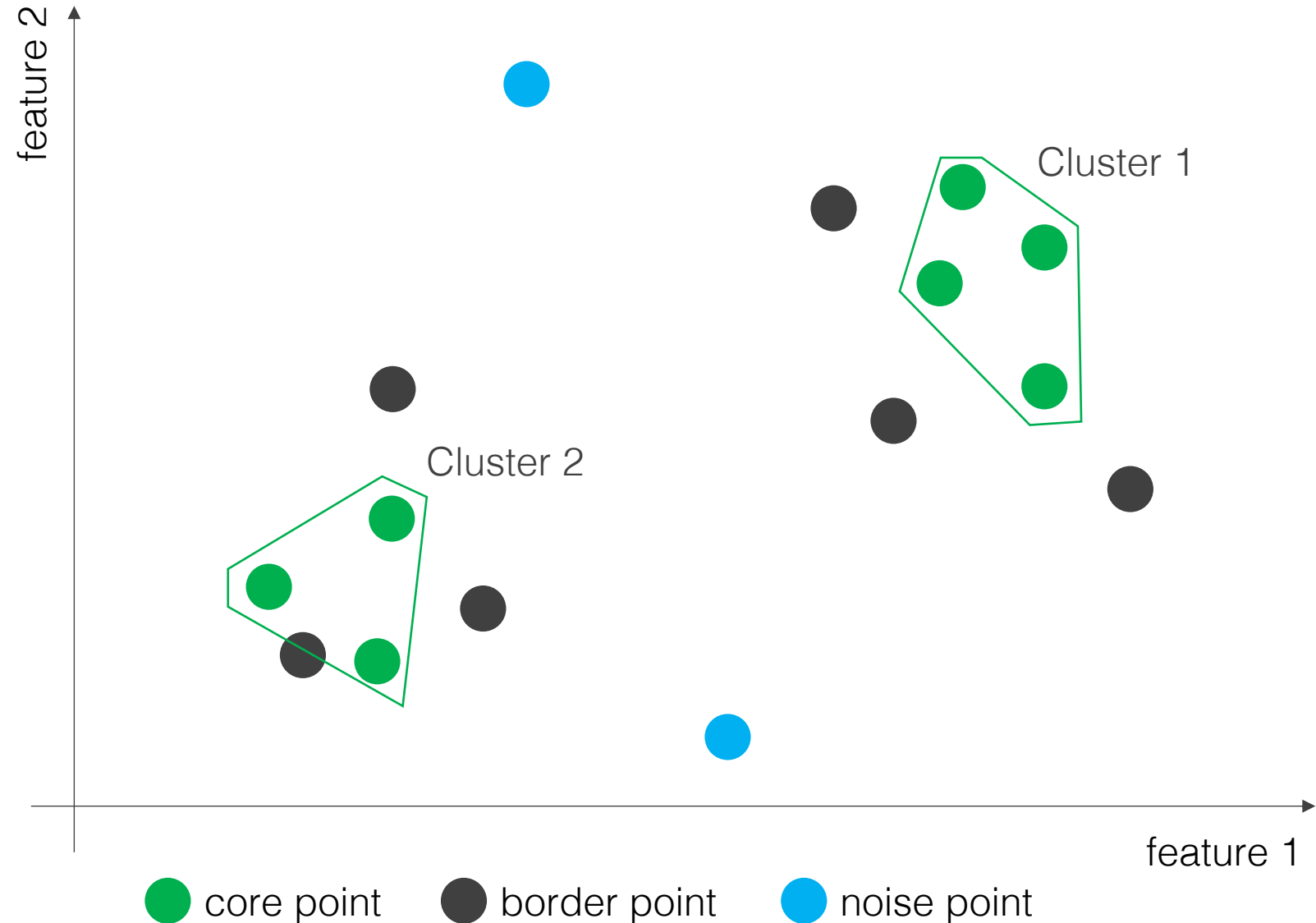
1. Distance measure
2. The radius of a neighbor, ϵ
3. 'minPts': The number of neighbors for a point to be considered a core point

Types of points:

- **Core:** a point with at least minPts neighbors
- **Border:** a non-core point that neighbors a core point
- **Noise:** Other points

Algorithm:

1. Label core and border points
2. Group neighboring core points
3. Add border points that are neighbors of core points



DBSCAN

Parameters:

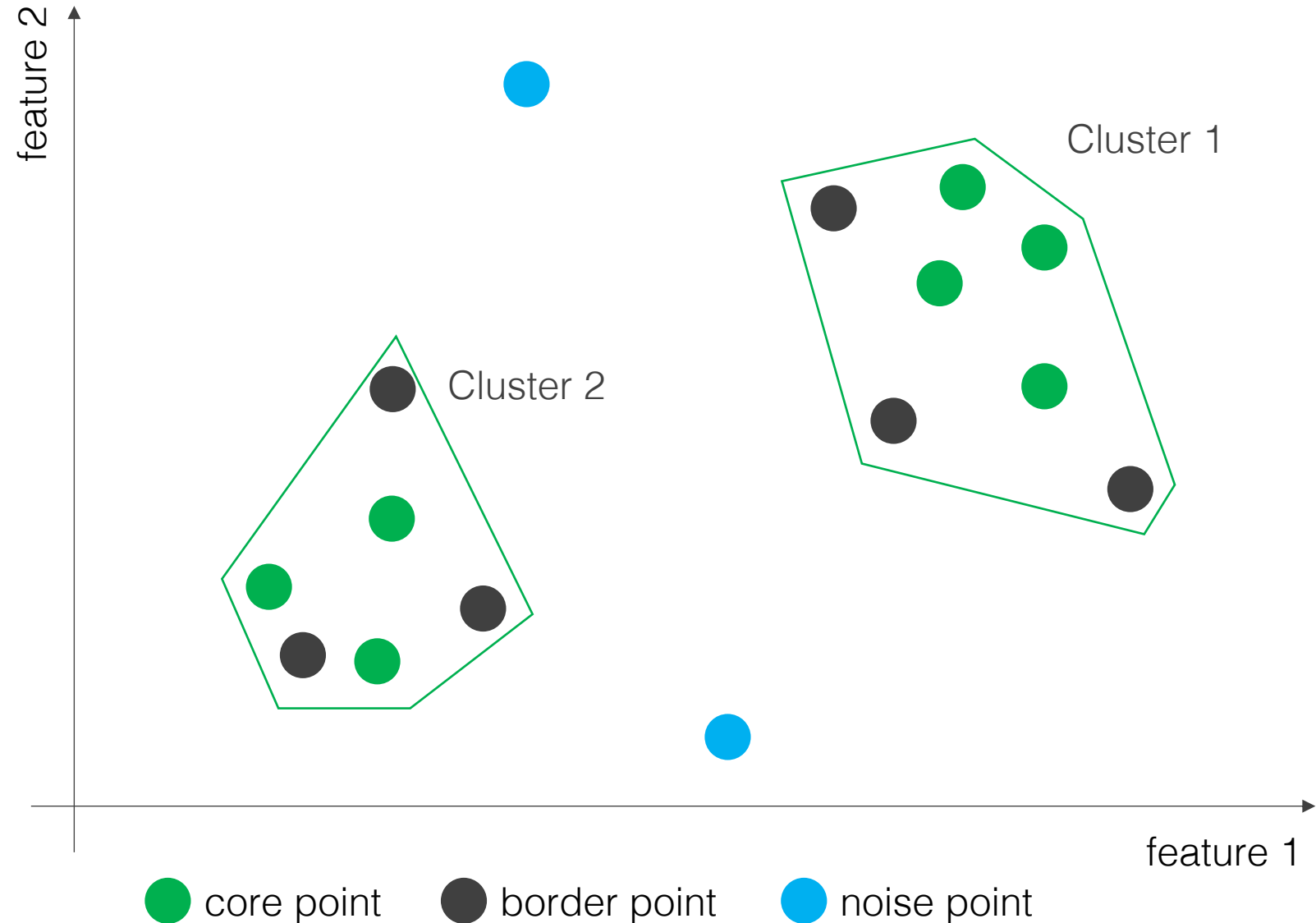
1. Distance measure
2. The radius of a neighbor, ϵ
3. 'minPts': The number of neighbors for a point to be considered a core point

Types of points:

- **Core:** a point with at least minPts neighbors
- **Border:** a non-core point that neighbors a core point
- **Noise:** Other points

Algorithm:

1. Label core and border points
2. Group neighboring core points
3. Add border points that are neighbors of core points



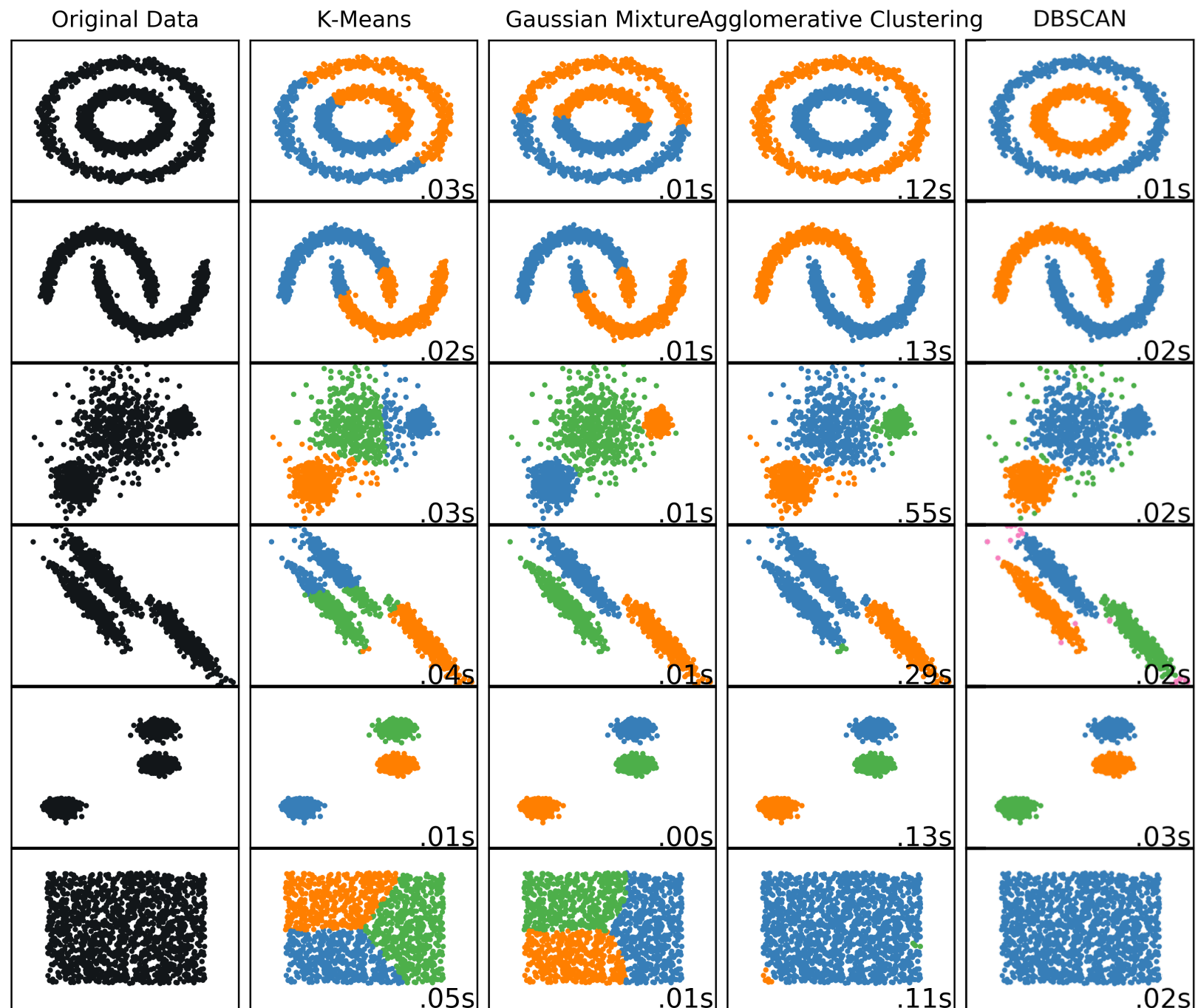
DBSCAN

- The number of clusters is chosen as part of the algorithm
- Can find arbitrarily shaped clusters
- Robust to outliers
- Cannot handle significant variation in cluster density
- Not entirely deterministic (border points reachable from more than one cluster may be assigned to either)

Examples: DBSCAN

Need to choose the
density parameters

Does not require
selecting the
number of clusters
beforehand



Spectral Clustering

Clustering in a low dimensional space based on data similarity

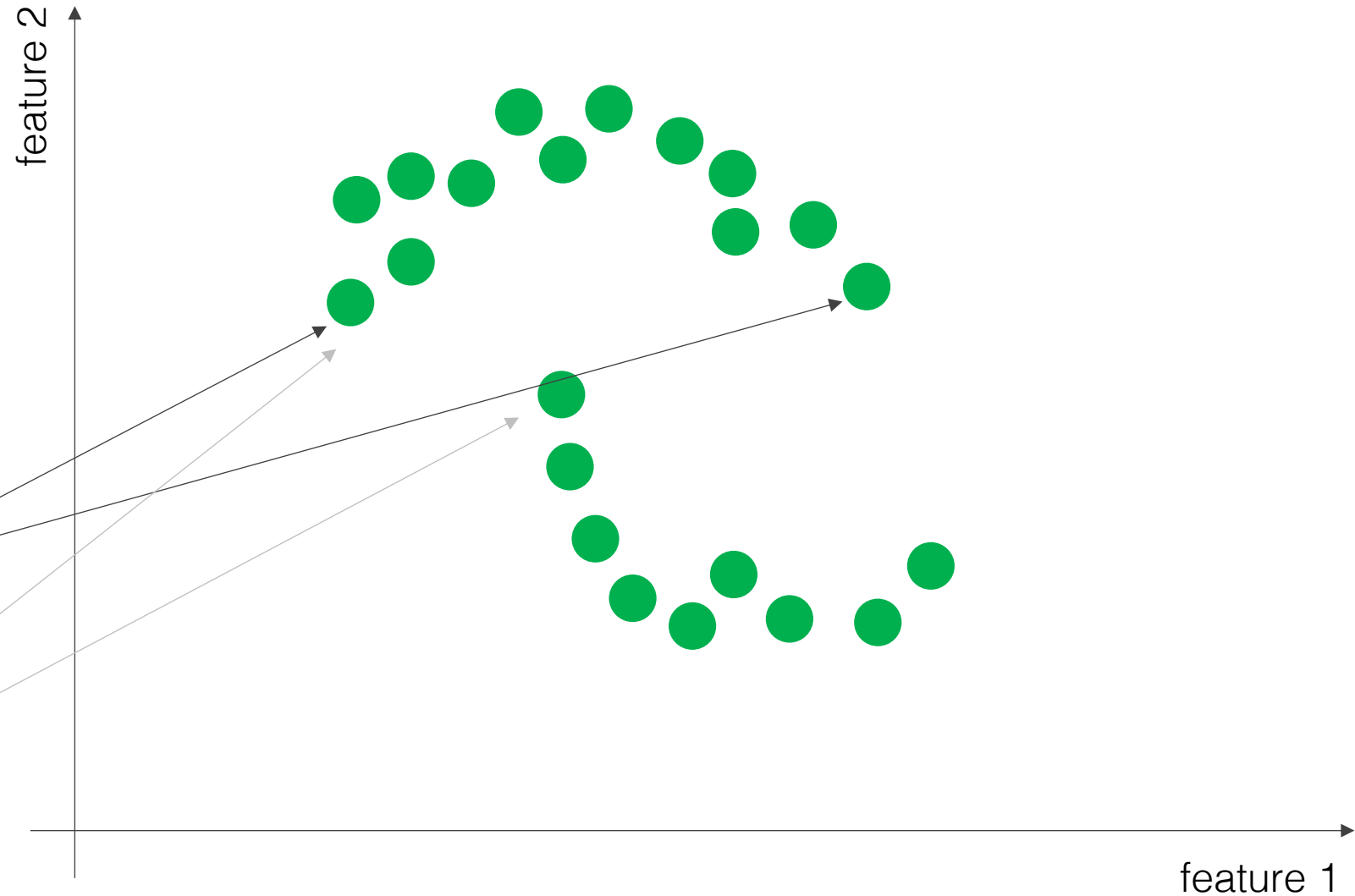
Spectral Clustering

Focuses on **connectedness** instead of compactness

The location alone does not determine **similarity** or “**affinity**”

These two points are likely connected by a cluster

These two points are NOT likely connected by a cluster



Concept from Sebastian Thrun and Peter Norvig

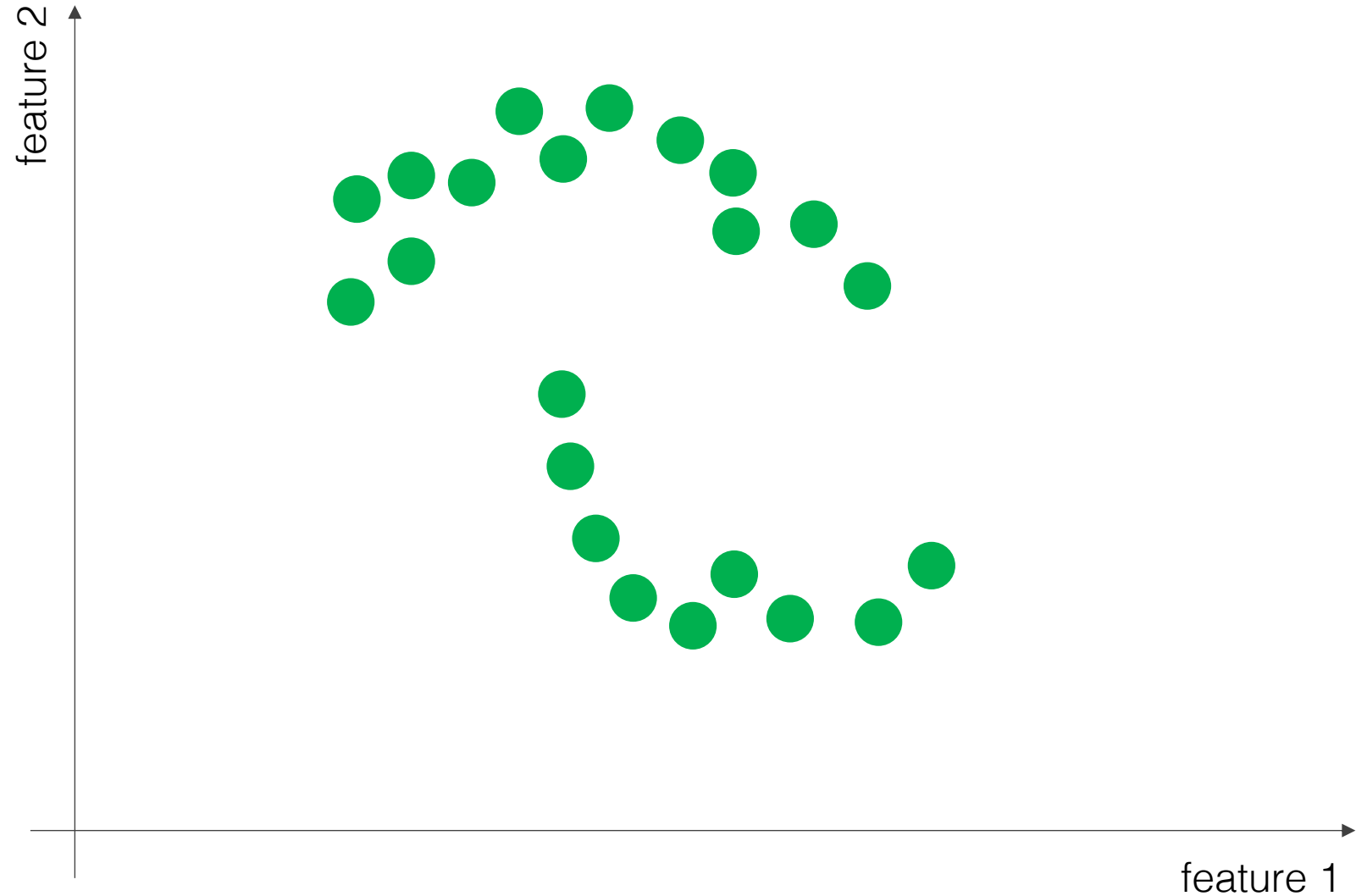
Spectral Clustering

Define **similarity** or **affinity** as the opposite of distance:

$$A(\mathbf{a}, \mathbf{b}) = -D(\mathbf{a}, \mathbf{b})$$

For example, using Euclidean distance, we could define affinity as:

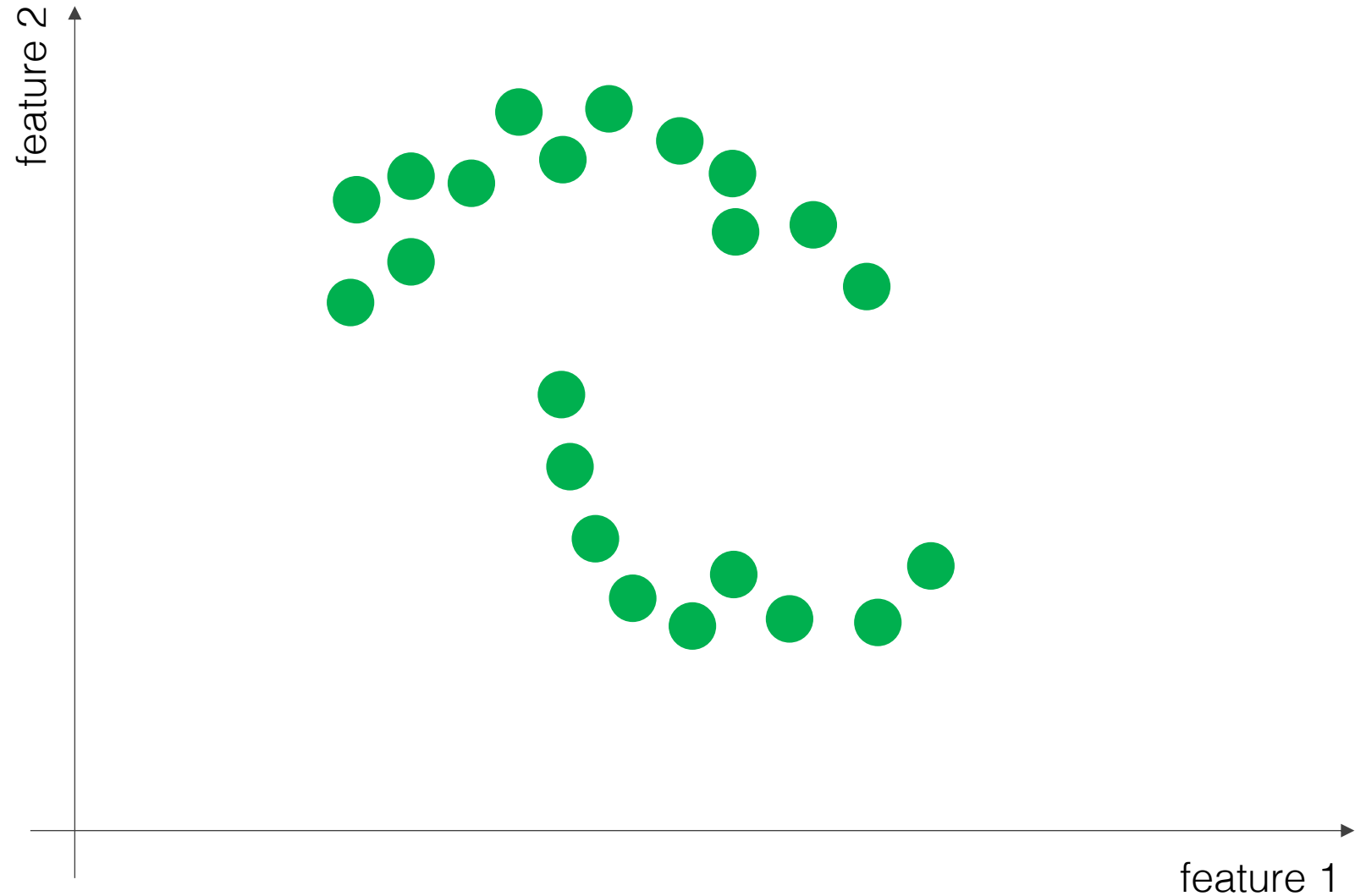
$$A(\mathbf{a}, \mathbf{b}) = -\|\mathbf{a} - \mathbf{b}\|_2$$



Spectral Clustering

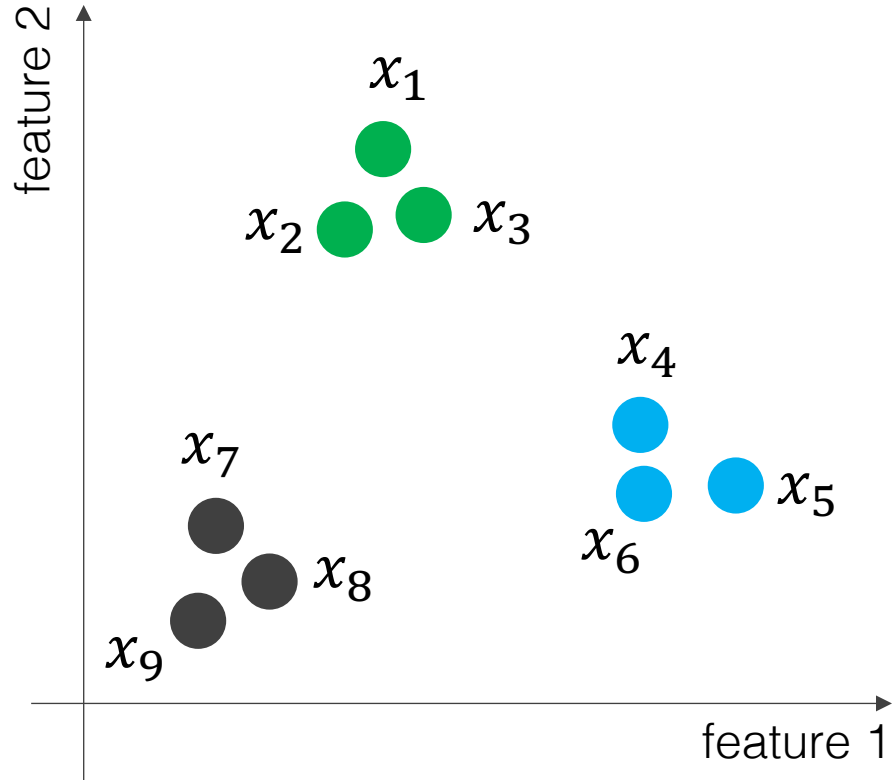
Algorithm

1. Construct an affinity matrix based on the data (works best when this matrix is sparse)
2. Reduce dimensions of the data using the affinity matrix
3. Perform clustering in this new lower dimensional space



Concept from Sebastian Thrun and Peter Norvig

Spectral Clustering



Affinity Matrix

	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9
x_1									
x_2									
x_3									
x_4									
x_5									
x_6									
x_7									
x_8									
x_9									

Affinity matrix

	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9
x_1									
x_2									
x_3									
x_4									
x_5									
x_6									
x_7									
x_8									
x_9									

Spectral Clustering

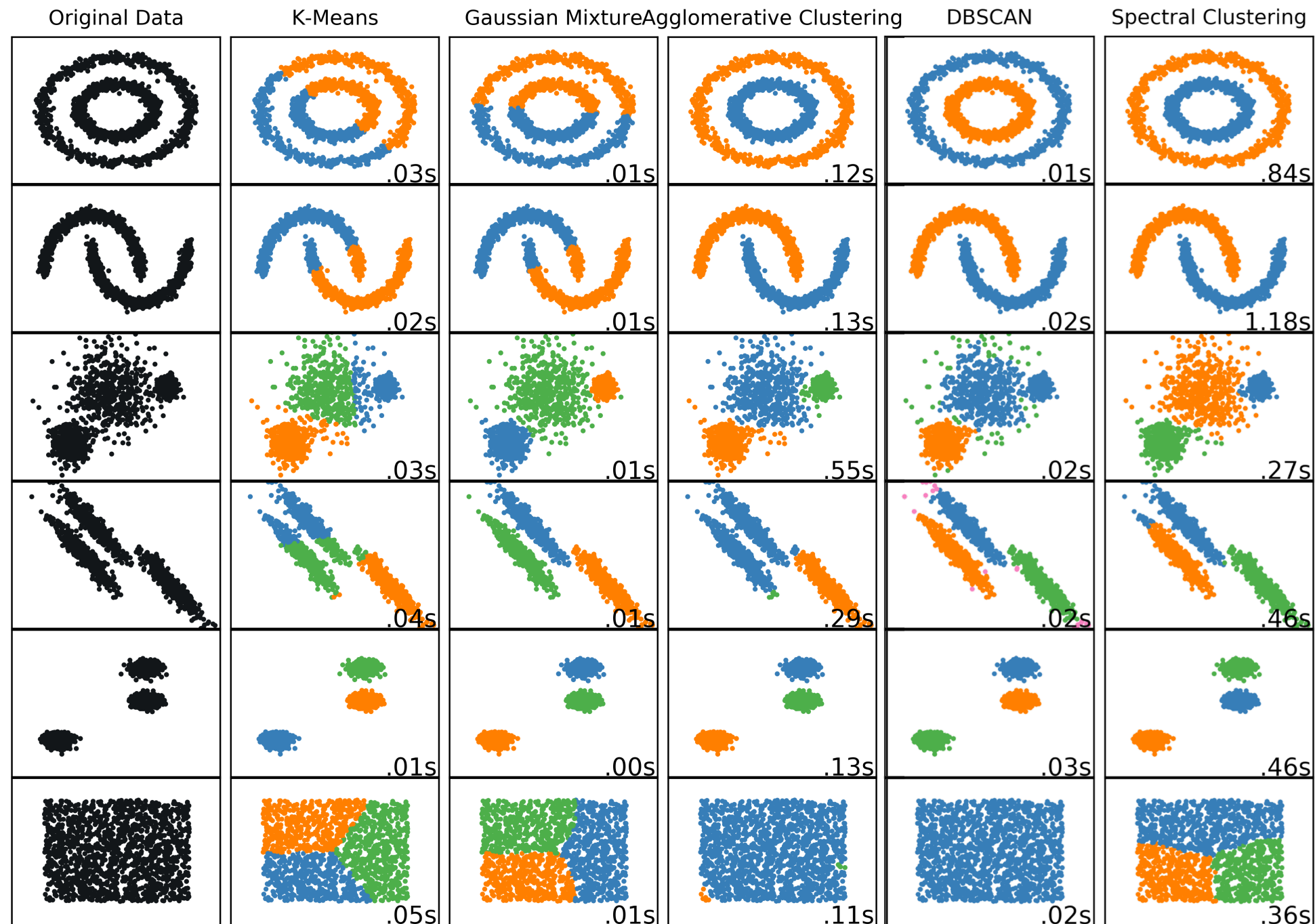
1. Project into lower dimension using the affinity matrix
2. Perform clustering in the lower dimension (often K-means)

Examples: Spectral Clustering

Makes few assumptions about data, so often produces good clustering results

Slow for large datasets

Requires specifying number of clusters



Types of clustering algorithms

Methods

Centroid-based clustering (e.g. **K-Means**)

Distribution-based clustering (e.g. **Gaussian mixture model**)

Density-based clustering (e.g. **DBSCAN**, mean-shift)

Hierarchical clustering (e.g. **agglomerative clustering**)
a.k.a. connectivity-based clustering

Graph-based clustering (e.g. **spectral clustering**, affinity propagation)

Cluster assignment

Hard clustering

Soft clustering (a.k.a. fuzzy clustering)

Clustering choices:

1. How should the data be scaled?
2. For K-means and GMMs: how many clusters to estimate?
3. For hierarchical clustering: dissimilarity measure, linkage, where to cut dendrogram

Approach: try multiple options, and select the one with the most useful or interpretable solution