

Clustering

Lecture 12

Termeh Shafie



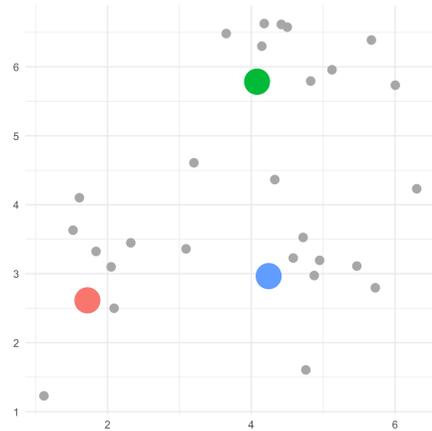
1

K-Means

2

K-Means

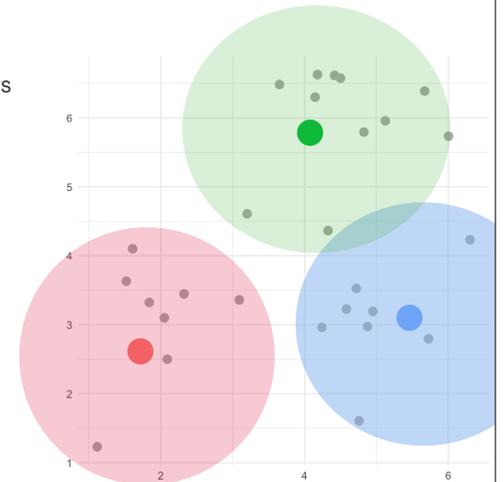
1. Choose k random points as cluster centers



3

K-Means

1. Choose k random points as cluster centers
2. For each data point, assign it the cluster whose centroid is the closest



4

K-Means

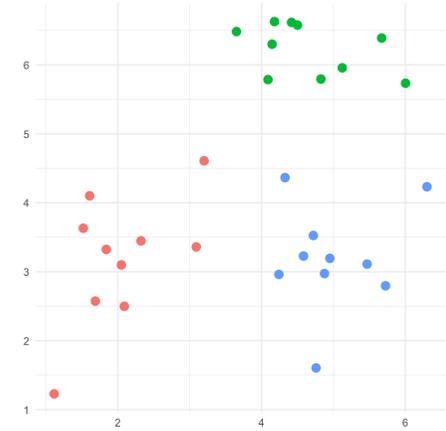
1. Choose k random points as cluster centers
2. For each data point, assign it the cluster whose centroid is the closest
3. Using these assignments, recalculate the centers



5

K-Means

1. Choose k random points as cluster centers
2. For each data point, assign it the cluster whose centroid is the closest
3. Using these assignments, recalculate the centers
4. Reiterate from step (2) until **convergence**:
 - cluster membership does not change
 - center only changes very very little



6

K-Means

$$C_1 \cup C_2 \cup \dots \cup C_K = \{1, \dots, n\}$$

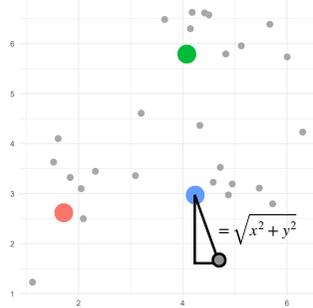
$$C_k \cap C_{k'} = \emptyset \text{ for all } k \neq k'$$

$$\min_{C_1, \dots, C_K} \left\{ \sum_{k=1}^K W(C_k) \right\}$$

within cluster variance

$$\text{where } W(C_k) = \frac{1}{|C_k|} \sum_{i, i' \in C_k} \sum_{j=1}^p (x_{ij} - x_{i'j})^2$$

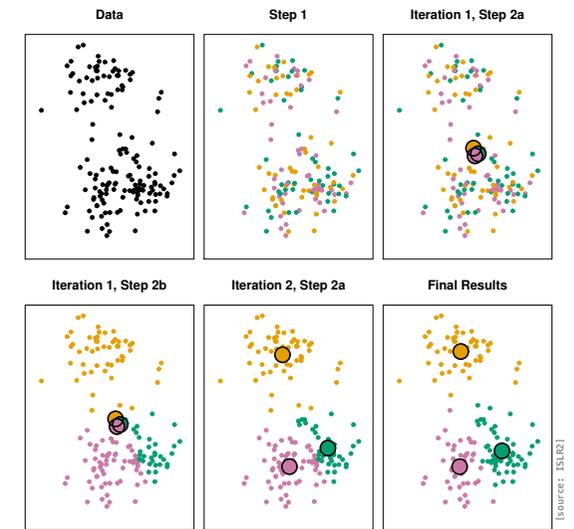
squared Euclidean distance



7

K-Means: Algorithm

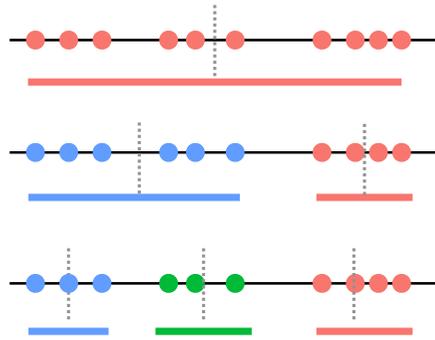
1. Choose k random points as cluster centers
2. For each data point, assign it the cluster whose centroid is the closest
3. Using these assignments, recalculate the centers
4. Reiterate from step (2) until **convergence**:
 - cluster membership does not change
 - center only changes very very little



8

How to decide on K

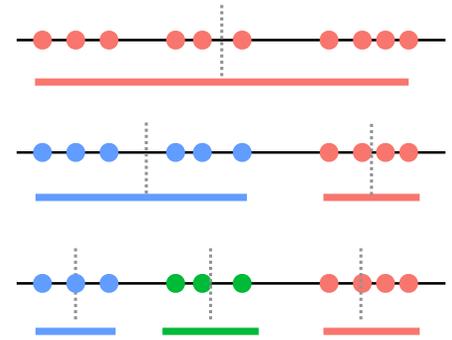
- Start with K=1: compute total variation
- Then K=2: compute total variation
- Then K=3: compute total variation
- ⋮



9

How to decide on K

- Start with K=1: compute total variation
- Then K=2: compute total variation
- Then K=3: compute total variation
- ⋮



10

How to decide on K

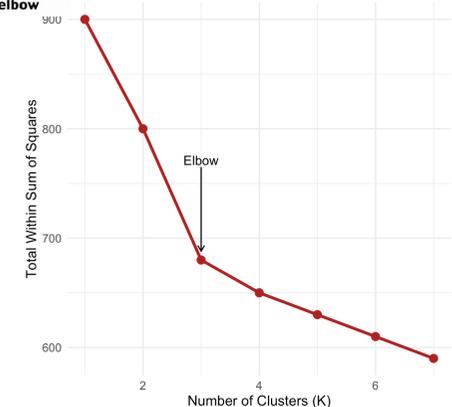
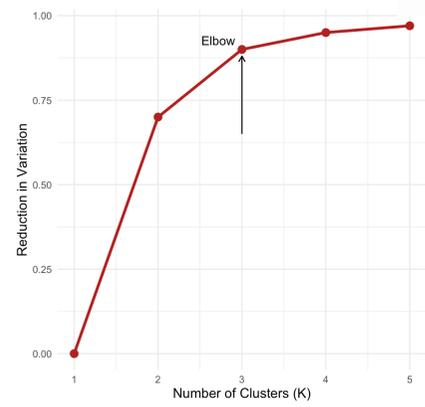
- Start with K=1: compute total variation
- Then K=2: compute total variation
- Then K=3: compute total variation
- ⋮



the reduction in total variation will get less and less as you increase K

11

How to decide on K



12

Cluster Assignments: Distortion

metric that assesses the performance of K-means (smaller values better)

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$$

Goal: choose r_{nk} and μ_k that minimizes J

hard assignments!

$$r_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_j \|x_n - \mu_j\|^2 \\ 0 & \text{otherwise} \end{cases}$$

$$\frac{dJ}{d\mu_k} = 2 \sum_{n=1}^N r_{nk} (x_n - \mu_k) = 0 \implies \mu_k = \frac{\sum_n r_{nk} x_n}{\sum_n r_{nk}} = \frac{1}{N_k} \sum_n r_{nk} x_n$$

optimal value for μ_k minimizing our loss is the mean of all data points in that cluster



13

Cluster Assignments: Distortion

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$$

- choose k random points as cluster centers
- for each data point, assign it the cluster whose centroid is the closest
- using these assignments, recalculate the centers
- reiterate from step (2) until convergence:
 - cluster membership does not change
 - center only changes very little

$$\frac{dJ}{d\mu_k} = 2 \sum_{n=1}^N r_{nk} (x_n - \mu_k) = 0 \implies \mu_k = \frac{\sum_n r_{nk} x_n}{\sum_n r_{nk}} = \frac{1}{N_k} \sum_n r_{nk} x_n$$



14

Evaluate Clustering: Silhouette Score

Cohesion: data points are similar within cluster

Separation: how far apart clusters are from other clusters

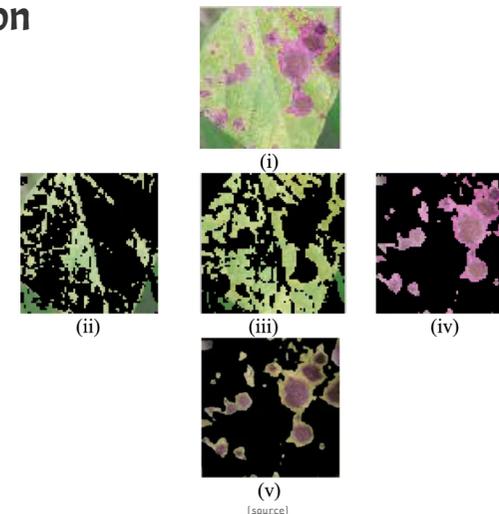
$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}$$

$a(i)$: average distance between data points and other members of its own cluster

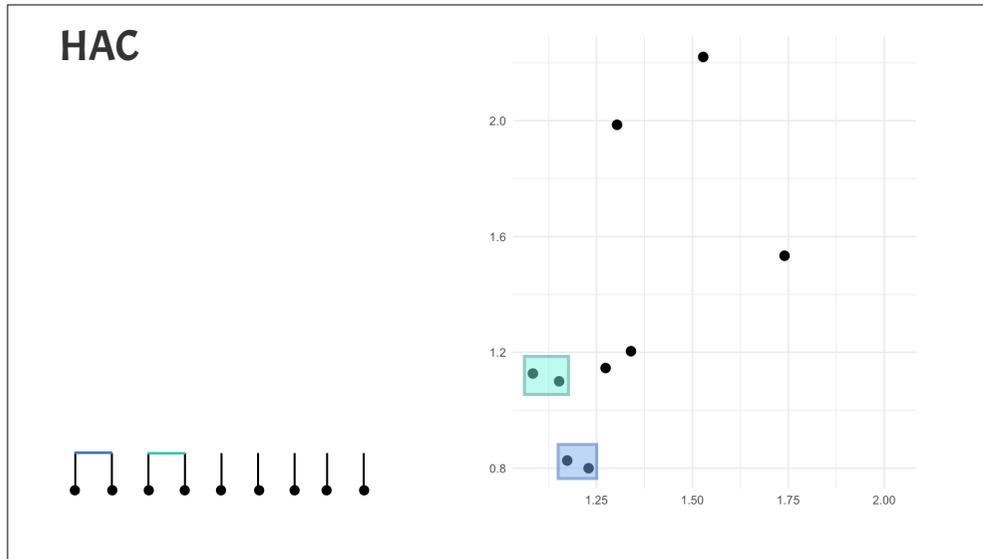
$b(i)$: average distance between a data point and the members of the next closest cluster

15

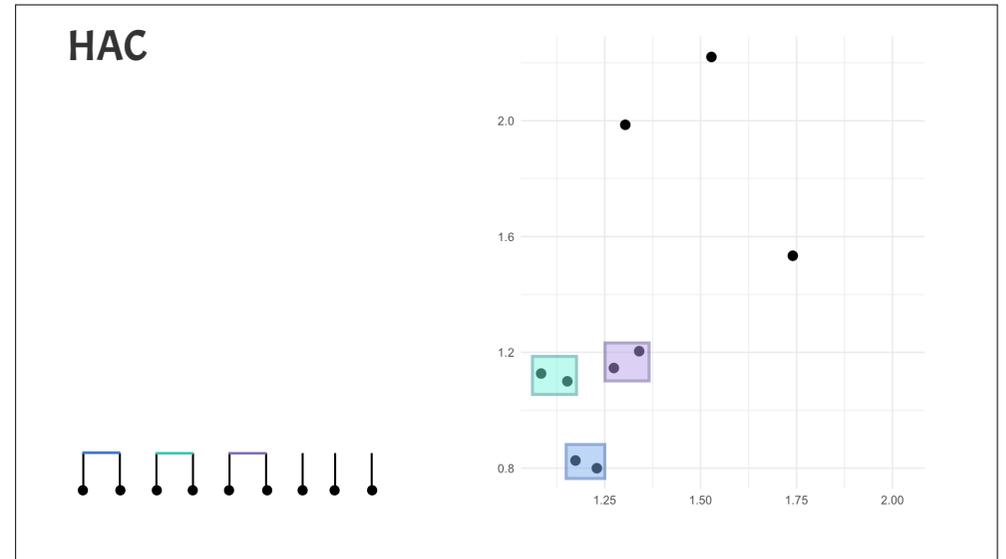
Application



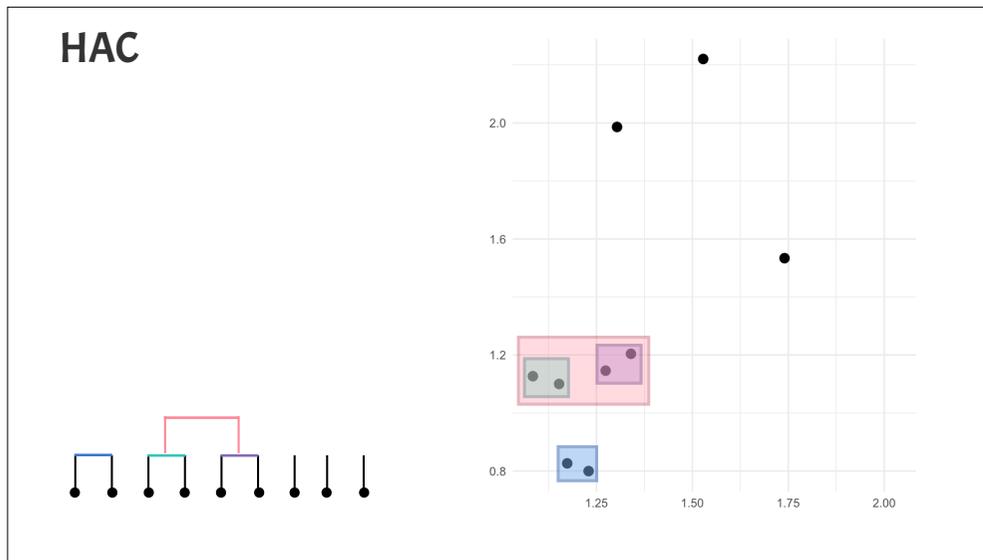
16



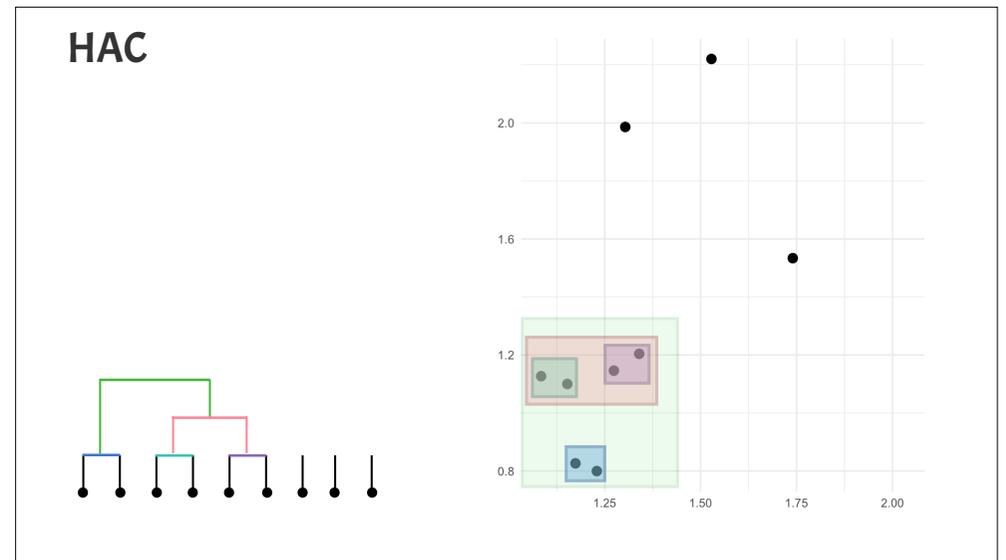
21



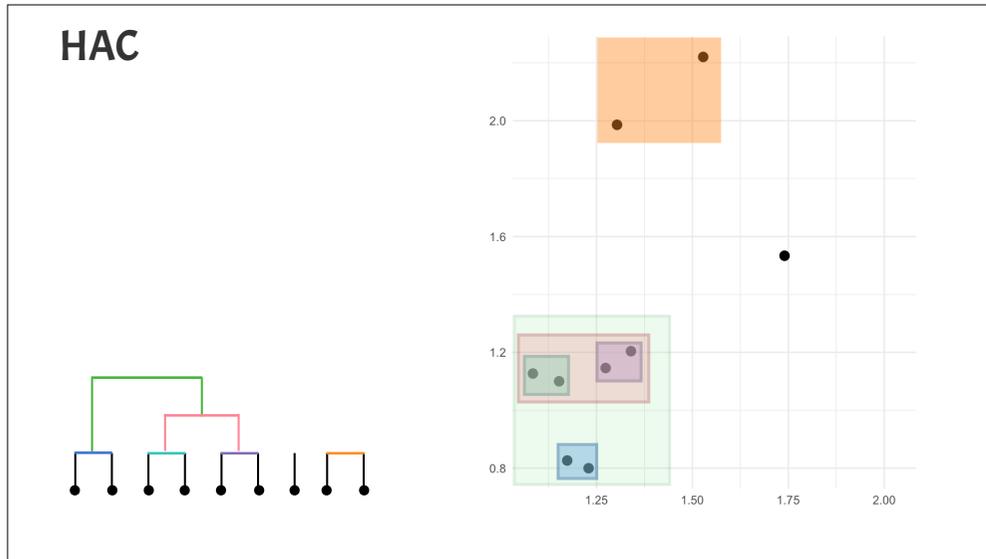
22



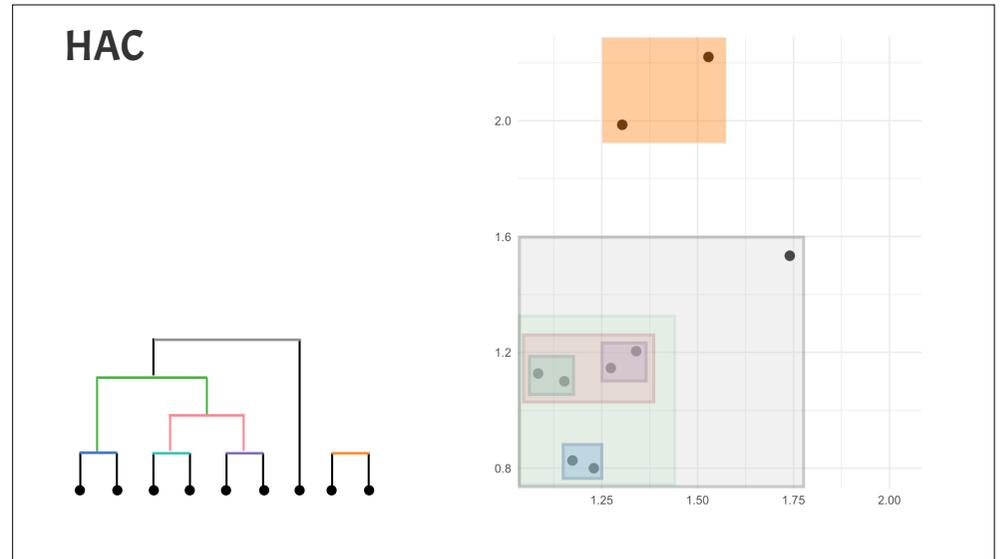
23



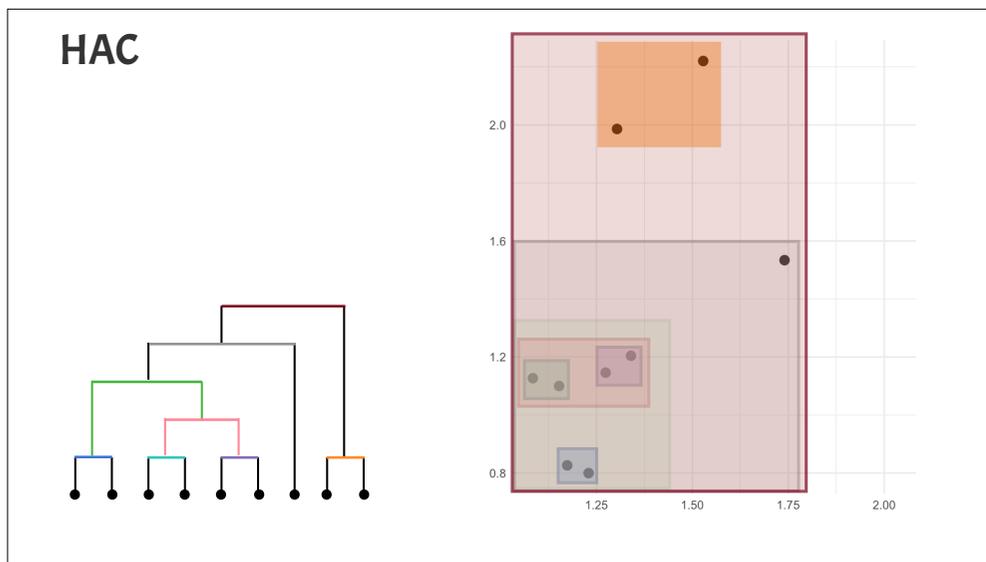
24



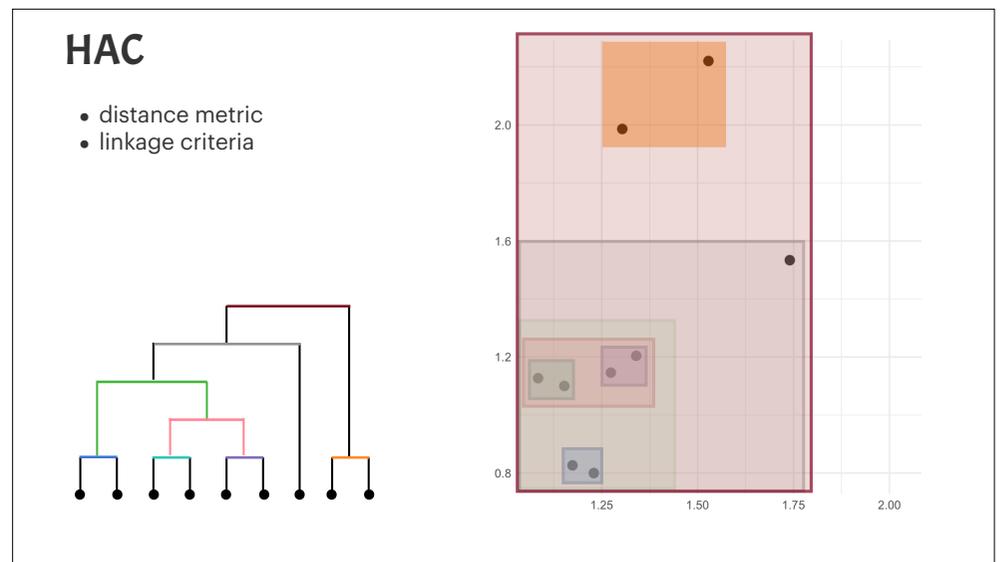
25



26



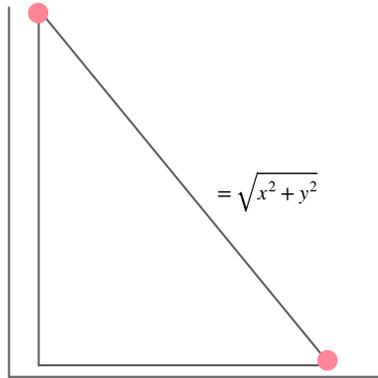
27



28

Distance Metrics: Euclidean

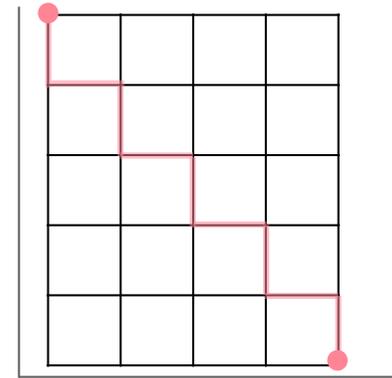
continuous data



29

Distance Metrics: Manhattan

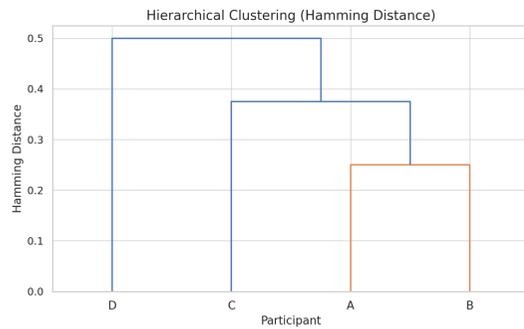
continuous data and high dimensions
(but also discrete or binary attributes)



30

Distance Metrics: Hamming

binary and categorical data



Participant	Q1	Q2	Q3	Q4
A	1	0	1	0
B	1	1	1	0
C	0	0	1	0
D	1	1	0	0



	A	B	C	D
A	0	1	2	2
B	1	0	2	1
C	2	2	0	3
D	2	1	3	0

31

Distance Metrics: Cosine

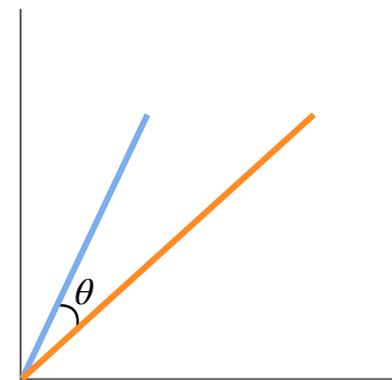
measures the angular difference between two vectors in a multi-dimensional space
word (or other) count data

$$\text{cosine similarity}(x, y) = \frac{x \cdot y}{\|x\| \|y\|} = \cos(\theta)$$

where:
 $x \cdot y = \sum_{i=1}^n x_i y_i$ is the dot product of the two vectors

$\|x\| = \sqrt{\sum_{i=1}^n x_i^2}$ is the Euclidean norm (length) of vector x

$\|y\| = \sqrt{\sum_{i=1}^n y_i^2}$ is the Euclidean norm (length) of vector y



32

Linkage Criteria

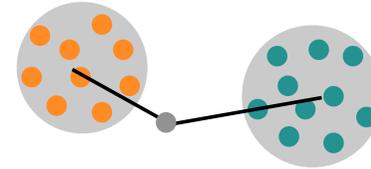
how compare different clusters?



33

Linkage Criteria

how compare different clusters?

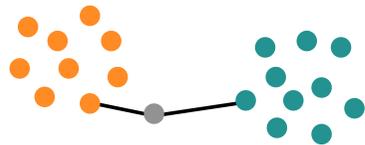


the average of each clusters ("centroid")

34

Linkage Criteria

how compare different clusters?

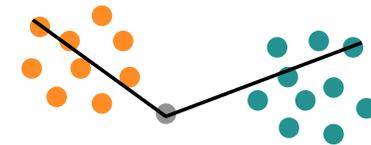


the closest point in each cluster ("single linkage")

35

Linkage Criteria

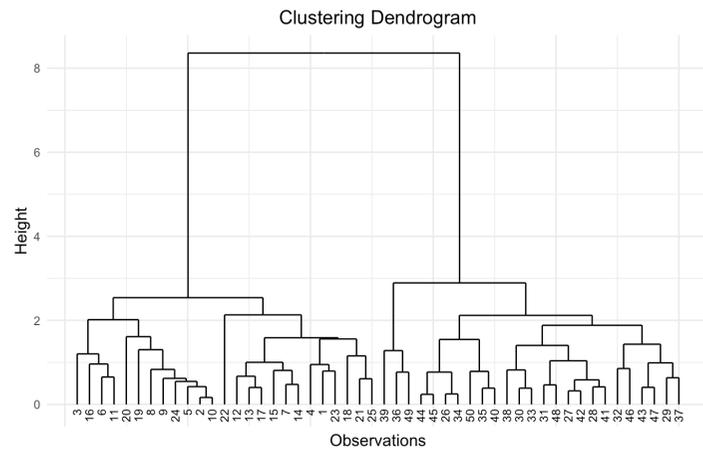
how compare different clusters?



the furthest point in each cluster ("complete linkage")

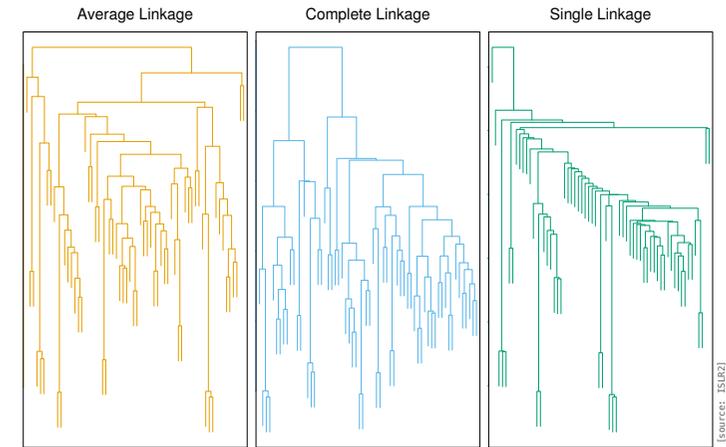
36

Reading a Dendrogram



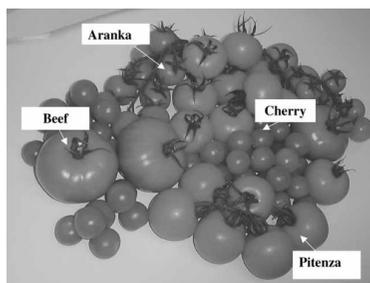
37

Balanced Clusters and Density

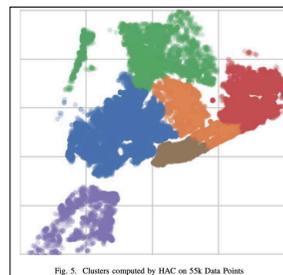


38

Applications



[source]



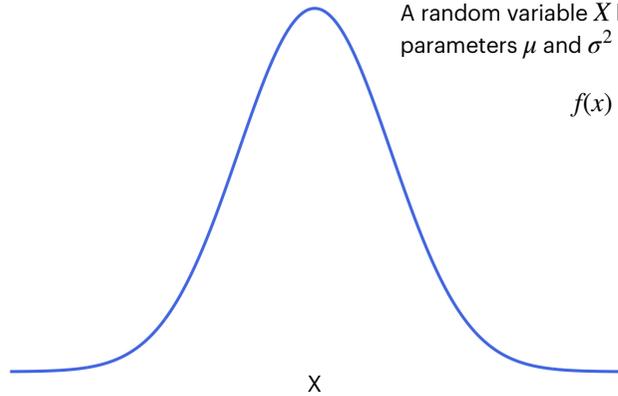
[source]

39

Gaussian Mixture Models (GMMs)

40

Normal (Gaussian) Distribution

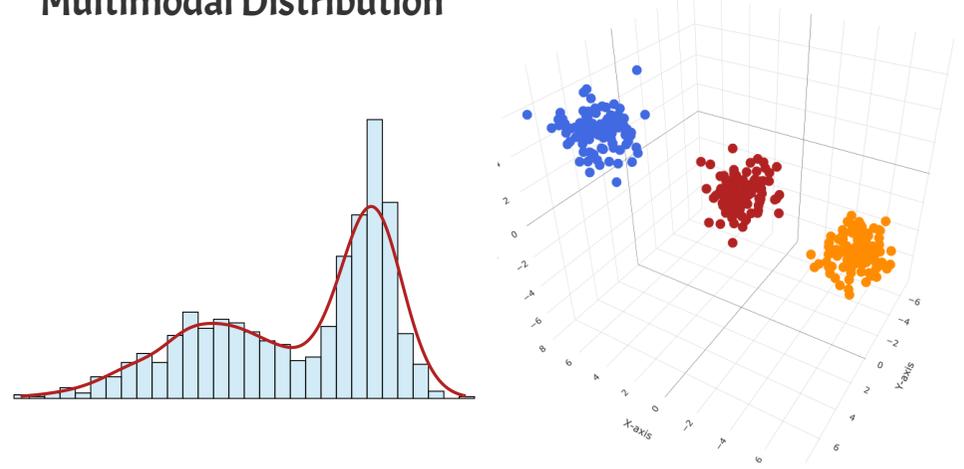


A random variable X has normal distribution with parameters μ and σ^2 if it has the following pdf:

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

41

Multimodal Distribution



42

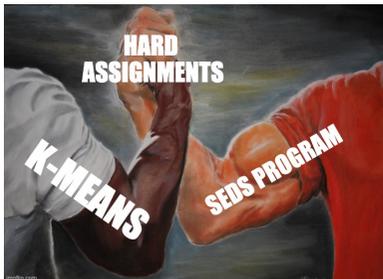
GMM

K-Means

- Hard assignment
- All variances are the same
- Roughly the same number of data points

GMM

- Soft (probabilistic) assignment
- Variances can be different
- Explicitly models number of data points



43

Recall K-Means Algorithm



- $$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$$
1. Choose k random points as cluster centers
 2. For each data point, assign it the cluster whose centroid is the closest
 3. Using these assignments, recalculate the centers
 4. Iterate from step (2) until **convergence**:
 - cluster membership does not change
 - center only changes very very little

$$\frac{dJ}{d\mu_k} = 2 \sum_{n=1}^N r_{nk} (x_n - \mu_k) = 0 \rightarrow \mu_k = \frac{\sum_{n=1}^N r_{nk} x_n}{\sum_{n=1}^N r_{nk}} = \frac{1}{N_k} \sum_{n=1}^N r_{nk} x_n$$

44

GMM: EM Algorithm

1. Choose k random points to be cluster centers (or estimate using k-means...)
2. For each data point, calculate the **probability** of belonging to each cluster
3. Using these probability weights, recalculate the **means + variances** (and weights)
4. Repeat 2 and 3 until **distributions converge**

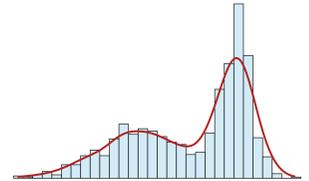
45

Multimodal Distribution

$$p(x) = \sum_{k=1}^K w_k \mathcal{N}(x | \mu_k, \Sigma_k)$$

probability of being in group k

likelihood of seeing x in group k



46

Posterior Probabilities

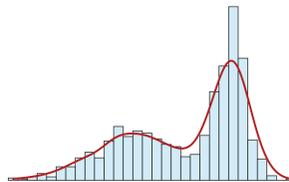
$$p(x) = w_k \mathcal{N}(x | \mu_k, \Sigma_k)$$

prior probability of being in group k

likelihood of seeing x in group k

$$p(\text{cluster } k | x) = \frac{w_k \mathcal{N}(x | \mu_k, \Sigma_k)}{\sum_{j=1}^K w_j \mathcal{N}(x | \mu_j, \Sigma_j)}$$

posterior probability of being in cluster k



47

Maximum Likelihood Estimation

$$p(x) = \sum_{k=1}^K w_k \mathcal{N}(x | \mu_k, \Sigma_k)$$

$$p(\mathbf{X} | \mathbf{w}, \mu, \Sigma) = p(x_1, x_2, \dots, x_n | \mathbf{w}, \mu, \Sigma) = \prod_{n=1}^N \sum_{k=1}^K w_k \mathcal{N}(x_n | \mu_k, \Sigma_k)$$

$$\log p(\mathbf{X} | \mathbf{w}, \mu, \Sigma) = \sum_{n=1}^N \log \left[\sum_{k=1}^K w_k \mathcal{N}(x_n | \mu_k, \Sigma_k) \right]$$

Goal: choose w, μ, Σ that maximizes the log likelihood



48

GMM: EM Algorithm

1. Choose k random points to be cluster centers (or estimate using k-means...)
2. For each data point, calculate the **probability** of belonging to each cluster
3. Using these probability weights, recalculate the **means + variances** (and weights)
4. Repeat 2 and 3 until **distributions converge**

49

The E-step in EM Algorithm

Responsibilities are the posterior probability of a data point being in cluster

$$p(\text{cluster } k | x) = \frac{w_k \mathcal{N}(x | \mu_k, \Sigma_k)}{\sum_{j=1}^K w_j \mathcal{N}(x | \mu_j, \Sigma_j)}$$

prior probability of being in group k likelihood of seeing x in group k
posterior probability of being in cluster k

How Likely is the cluster?
Many/few data points there?

How well does that data fit with that cluster?

$$r_{nk} = \frac{w_k \mathcal{N}(x | \mu_k, \Sigma_k)}{\sum_{j=1}^K w_j \mathcal{N}(x | \mu_j, \Sigma_j)}$$

normalize to get a probability

Responsibility is high if the data point is likely to belong to that cluster rather than other clusters

this is soft assignment

50

GMM: EM Algorithm

1. Choose k random points to be cluster centers (or estimate using k-means...)
2. For each data point, calculate the **probability** of belonging to each cluster
3. Using these probability weights, recalculate the **means + variances** (and weights)
4. Repeat 2 and 3 until **distributions converge**

51

The M-step in EM Algorithm

Via MLE we get the following estimates:

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} x_n \quad \Sigma_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (x_n - \mu_k)(x_n - \mu_k)^T \quad w_k = \frac{N_k}{N} = \frac{1}{N} \sum_{n=1}^N r_{nk}$$

the higher the responsibility of a data point for a cluster is, the more influence it has on what the mean and variance is

Note: $N_k = \sum_{n=1}^N r_{nk}$ is now based on soft assignments now

if data points are unlikely to belong to cluster k , the N_k small,
if data points are likely to belong to cluster k , then N_k large

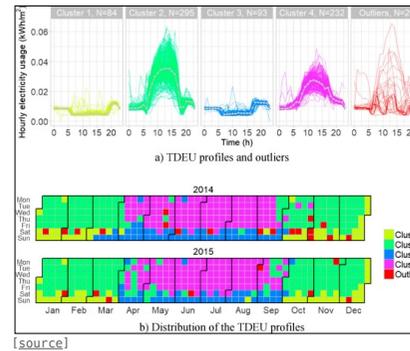
52

Take Aways

- GMM does **soft assignment**, every data point belongs to every cluster with some probability
- Data points that are more likely to be in a cluster have **more influence** over its parameters
- GMM uses the EM algorithm to iteratively update the cluster distributions:
 - first assign a responsibility to each data point (**E-step**)
 - then using them to calculate weighted means and variances for each cluster (**M-step**)
- Responsibilities measure **the probability of a data point being in each cluster** (technically the posterior probability).
- Responsibilities contain information about how common a cluster is as well as **the likelihood of a data point belonging to that cluster**

53

Applications



[source]

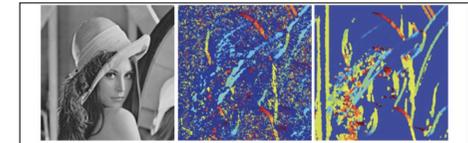
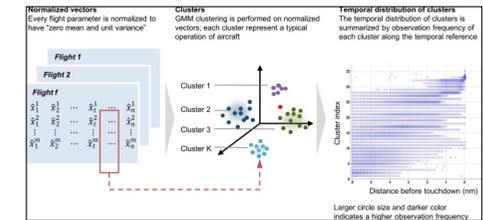


Fig. 1. Illustration of clustering of patches in the PLE method for the Lena image. LEFT: Original image; RIGHT: Clustered image; The pixels in the same color indicate that 8×8 patches around them are in the same cluster. It can be seen that patches from different parts of image are grouped into one cluster [17].

[source]



[source]

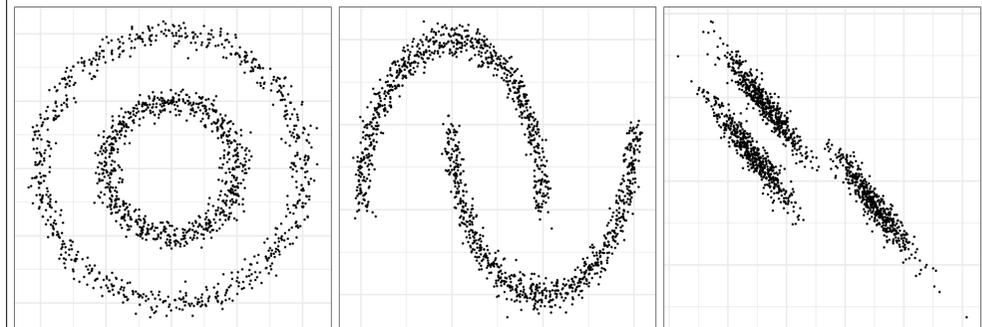
54

DBSCAN

Density Based Spatial Clustering of Applications **with Noise**

55

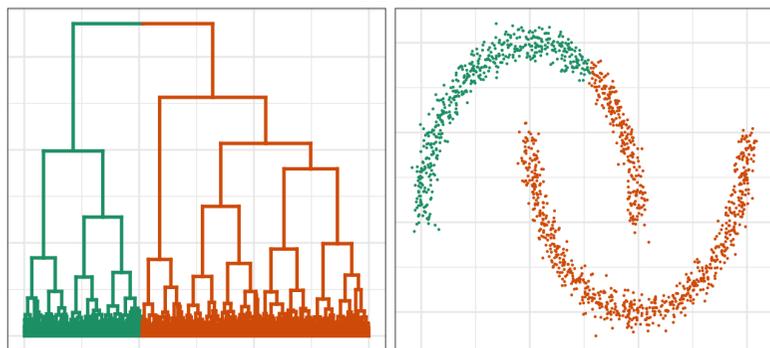
Example Data Structures



56

Where Others Fail...

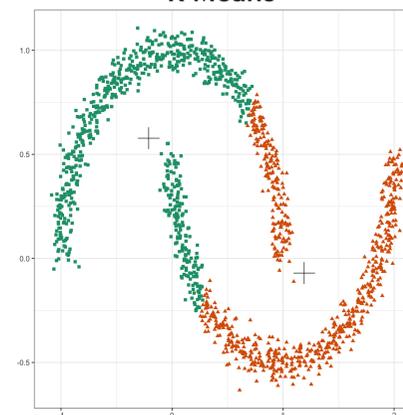
HAC



57

Where Others Fail...

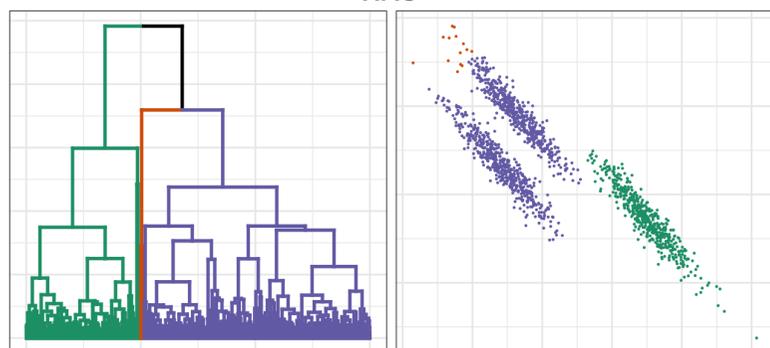
K-Means



58

Where Others Fail...

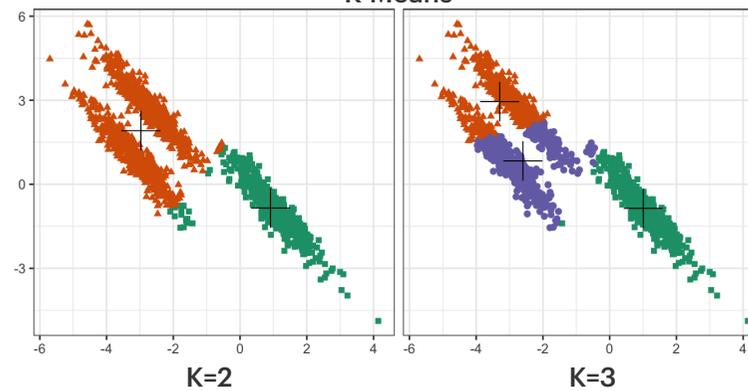
HAC



59

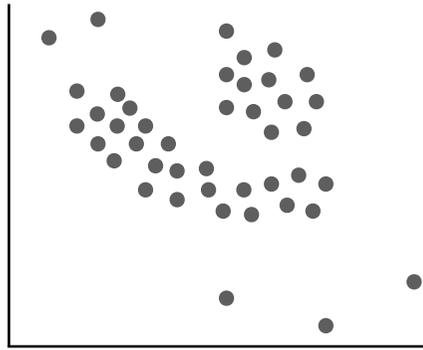
Where Others Fail...

K-Means



60

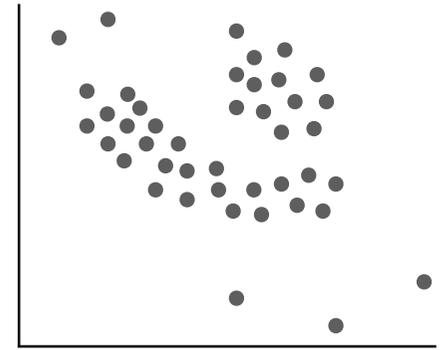
The Algorithm



61

Hyperparameters

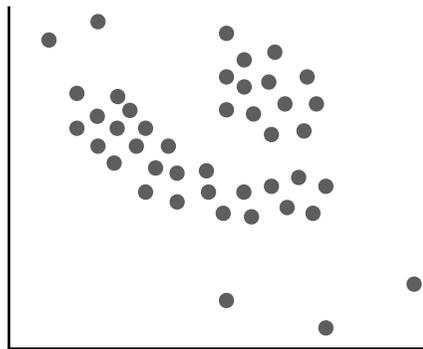
1. Distance metric
2. Epsilon (ϵ)
3. Minimum Points ($minpts$)



62

Hyperparameters

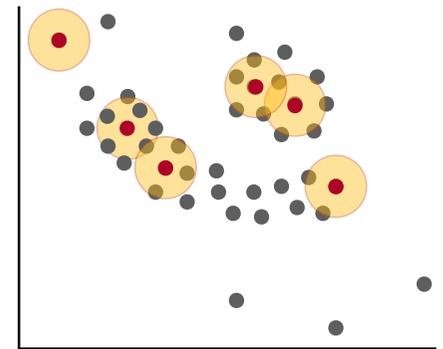
1. Distance metric
2. Epsilon (ϵ) 
3. Minimum Points ($minpts$)



63

Hyperparameters

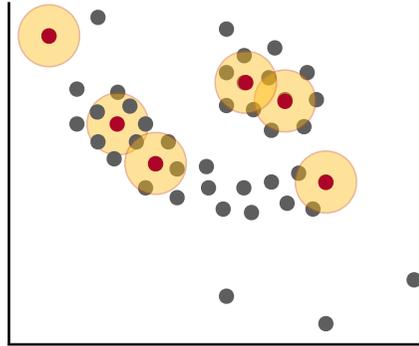
1. Distance metric
2. Epsilon (ϵ)
3. Minimum Points ($minpts$)



64

Hyperparameters

1. Distance metric
2. Epsilon (*eps*)
3. Minimum Points (*minpts*)
minimum number of points within *eps* distance of it in order to be considered dense

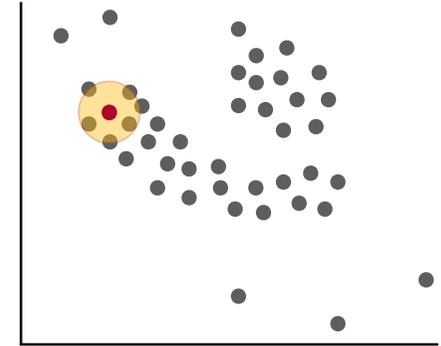


65

The Algorithm

Core Point

A point is a core point if it has **at least** *minpts* neighbors within *eps* distance of itself.



66

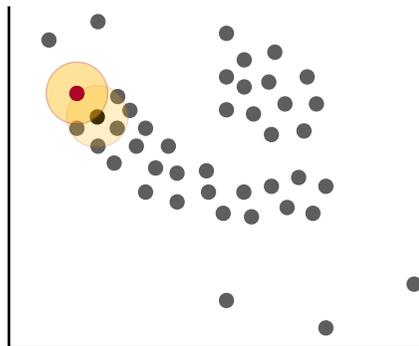
The Algorithm

Core Point

A point is a core point if it has **at least** *minpts* neighbors within *eps* distance of itself.

Border Point

A point **without** at least *minpts* neighbors within *eps* distance of itself, **but** is a neighbor of a core point.



67

The Algorithm

Core Point

A point is a core point if it has **at least** *minpts* neighbors within *eps* distance of itself.

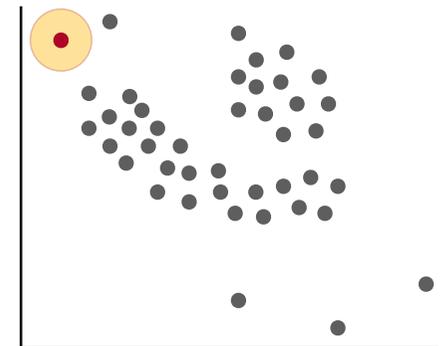
Border Point

A point **without** at least *minpts* neighbors within *eps* distance of itself, **but** is a neighbor of a core point.

Noise

A point **without** at least *minpts* neighbors within *eps* distance of itself, **and is not a** neighbor of a core point.

these three types of points will define our clusters

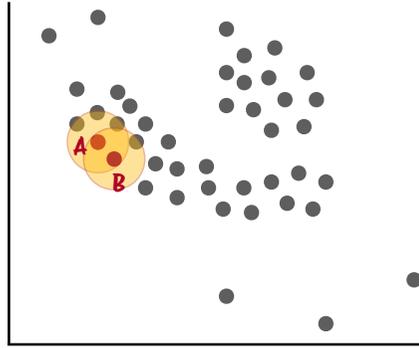


68

The Algorithm

Directly density reachable

Point A is directly density reachable from a **core-point** B if it is in the neighborhood of B.

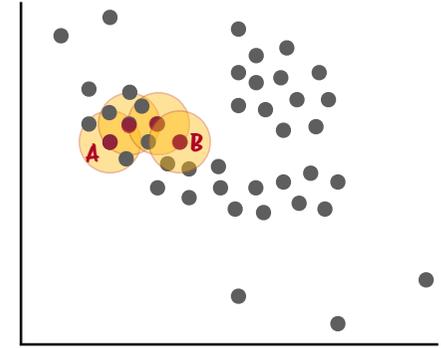


69

The Algorithm

Density reachable

Point A is density reachable from a core-point B if there are a **chain of points that are directly density reachable** from B to A.

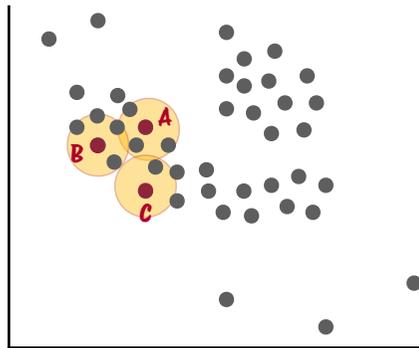


70

The Algorithm

Density connected

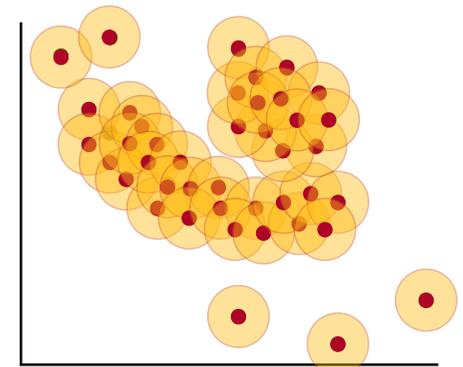
Point B and C are density connected if they are both density reachable from a third point A.



71

The Algorithm

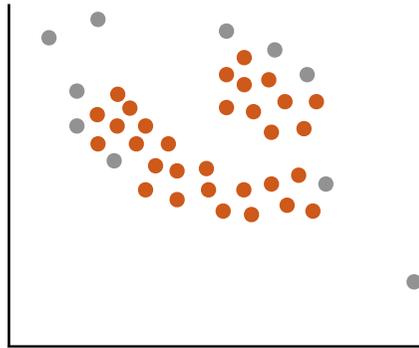
1 Find all core points.



72

The Algorithm

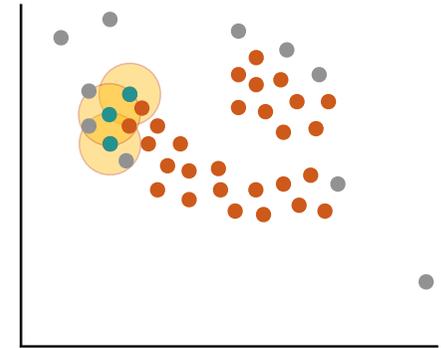
- 1 Find all core points.



73

The Algorithm

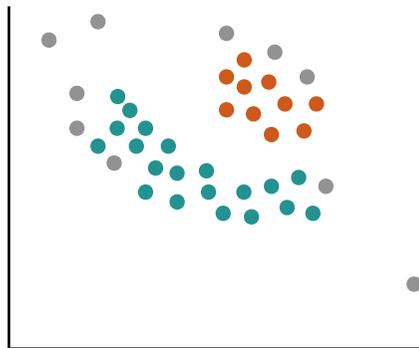
- 1 Find all core points.
- 2 Pick random core point, find other core points that are density reachable from it and assign to cluster.



74

The Algorithm

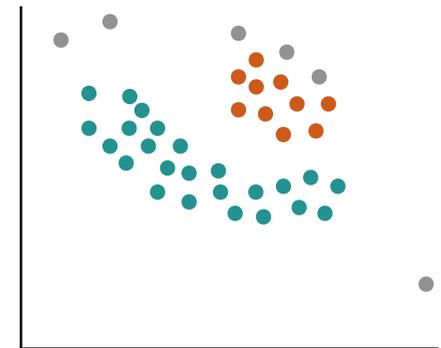
- 1 Find all core points.
- 2 Pick random core point, find other core points that are density reachable from it and assign to cluster.



75

The Algorithm

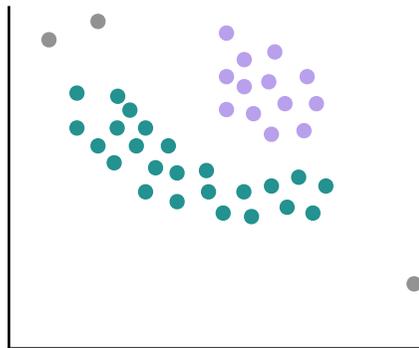
- 1 Find all core points.
- 2 Pick random core point, find other core points that are density reachable from it and assign to cluster.
- 3 Add border points and cluster one is done.



76

The Algorithm

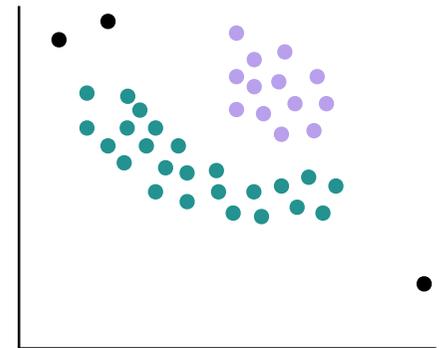
- 1 Find all core points.
- 2 Pick random core point, find other core points that are density reachable from it and assign to cluster.
- 3 Add border points and cluster one is done.
- 4 Repeat from (2) until all clusters detected.



77

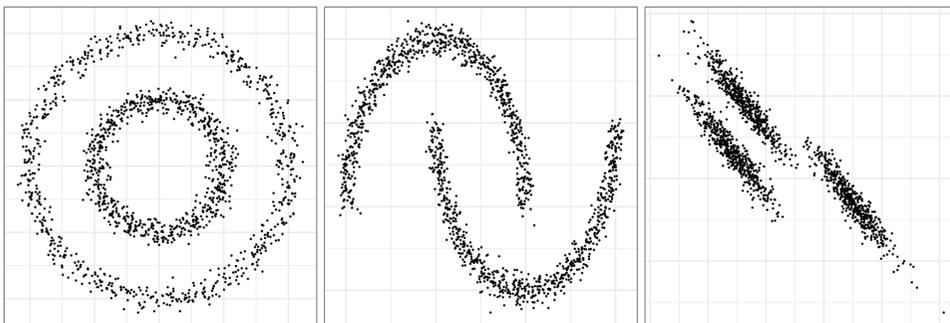
The Algorithm

- 1 Find all core points.
- 2 Pick random core point, find other core points that are density reachable from it and assign to cluster.
- 3 Add border points and cluster one is done.
- 4 Repeat from (2) until all clusters detected.
- 5 Remaining points are noise.



78

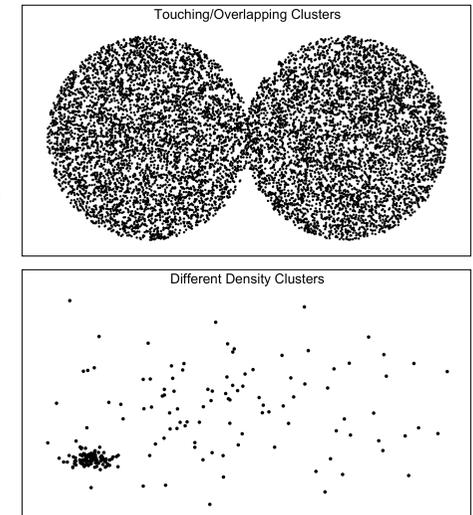
Example Data Structures



79

Where DBSCAN fails...

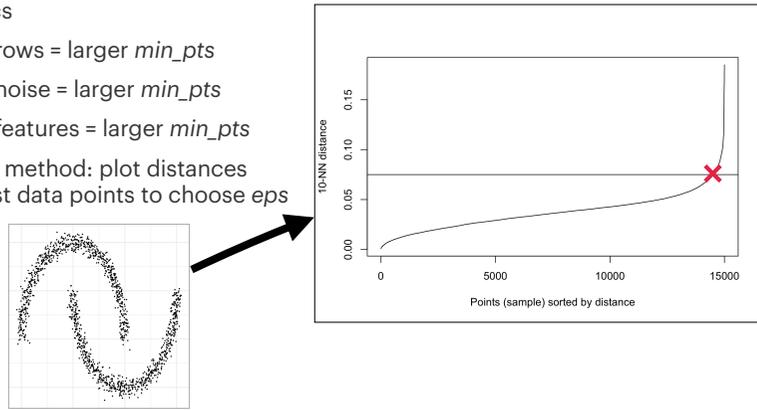
- Less effective on high dimensional data
- Overlapping/touching clusters
- Clusters have different densities



80

Hyperparameter Tuning

- Domain Knowledge + Distance Metrics
- More rows = larger *min_pts*
- More noise = larger *min_pts*
- More features = larger *min_pts*
- Elbow method: plot distances against data points to choose *eps*



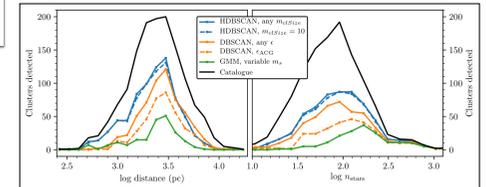
81

Applications

Figure 9 presents a sample image that was segmented using DBSCAN. In the figure, the individual clusters are regrouped together forming close to the original image. It can be observed that the pixels of similar color are clustered together.



[source]

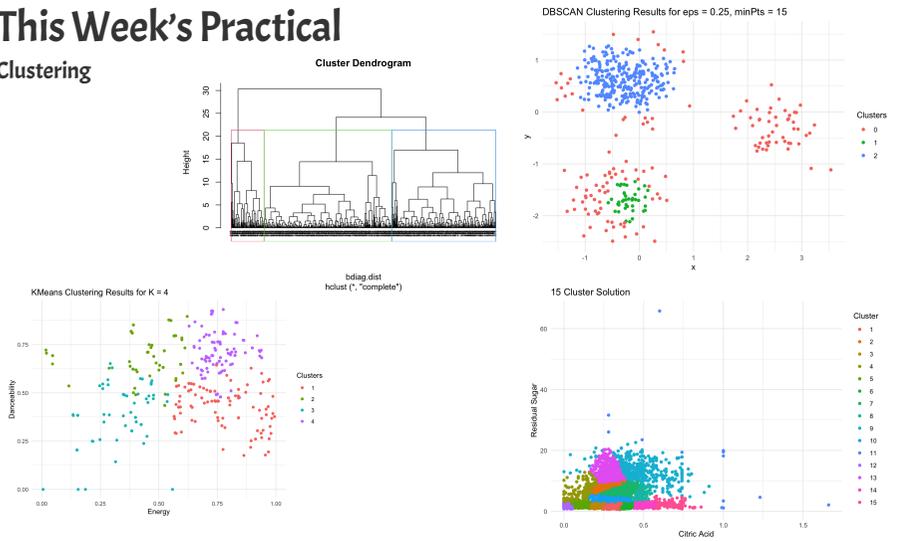


[source]

82

This Week's Practical

Clustering



83