Lecture 13: 11 March, 2025

Madhavan Mukund https://www.cmi.ac.in/~madhavan

Data Mining and Machine Learning January–April 2025

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

Unsupervised learning

Supervised learning requires labelled data



Unsupervised learning



- Vast majority of data is unlabelled
- What insights can you get with unlabelled data?



Unsupervised learning

- Supervised learning requires labelled data
- Vast majority of data is unlabelled
- What insights can you get with unlabelled data?

"If intelligence was a cake, unsupervised learning would be the cake, supervised learning would be the icing on the cake"

> – Yann LeCun ACM Turing Award 2018



- Customer segmentation
 - Marketing campaigns





- Customer segmentation
 - Marketing campaigns
- Anomaly detection
 - Outliers





- Customer segmentation
 - Marketing campaigns
- Anomaly detection
 - Outliers
- Semi-supervised learning
 - Propagate limited labels





- Customer segmentation
 - Marketing campaigns
- Anomaly detection
 - Outliers
- Semi-supervised learning
 - Propagate limited labels
- Image segmentation
 - Object detection





Find natural groups of data



> < ≣ >

э

- Find natural groups of data
- Define a distance measure



< E

- Find natural groups of data
- Define a distance measure
- Group together data that is close together





- Find natural groups of data
- Define a distance measure
- Group together data that is close together
- Top down
 - Partition data into clusters





- Find natural groups of data
- Define a distance measure
- Group together data that is close together
- Top down
 - Partition data into clusters
- Bottom up
 - Group items into clusters





Lecture 13: 11 March, 2025

Top down clustering

K Means Clustering

Data items are points in n dimensions

• $(x_1, x_2, ..., x_n)$





Top down clustering

K Means Clustering

- Data items are points in n dimensions
 - $(x_1, x_2, ..., x_n)$
- Partition into K clusters
 - Fix *K* in advance





Lecture 13: 11 March, 2025

Top down clustering

K Means Clustering

- Data items are points in n dimensions
 - (x_1, x_2, \ldots, x_n)
- Partition into K clusters
 - Fix K in advance
- Each cluster is represented by its geometric centre
 - Centroid, or mean
 - Hence "K means"





Lecture 13: 11 March, 2025

 Choose K points initially as random centroids



1=2

→

э

- Choose K points initially as random centroids
- In each iteration
 - Assign each point to nearest centroid
 - Recompute centroids



- Choose K points initially as random centroids
- In each iteration
 - Assign each point to nearest centroid
 - Recompute centroids



- Choose K points initially as random centroids
- In each iteration
 - Assign each point to nearest centroid
 - Recompute centroids

10

- Choose K points initially as random centroids
- In each iteration
 - Assign each point to nearest centroid
 - Recompute centroids

10 Q 0 0

- Choose K points initially as random centroids
- In each iteration
 - Assign each point to nearest centroid
 - Recompute centroids

0 10

- Choose K points initially as random centroids
- In each iteration
 - Assign each point to nearest centroid
 - Recompute centroids

10 0 **O** 0 0

- Choose K points initially as random centroids
- In each iteration
 - Assign each point to nearest centroid
 - Recompute centroids
- Termination
 - Clusters stabilize
 - Sum square distance is below threshold

э

How "tight" are the clusters?

▶ < ∃ ▶</p>

э

- How "tight" are the clusters?
- Mean squared distance from centroids inertia

$$\frac{1}{n} \sum_{j=1}^{K} \sum_{x \in C_j} distance(x, centroid_j)^2$$



- How "tight" are the clusters?
- Mean squared distance from centroids inertia

 $\frac{1}{n} \sum_{j=1}^{K} \sum_{x \in C_j} distance(x, centroid_j)^2$

 Plot inertia for different values of K and look for optimum



- How "tight" are the clusters?
- Mean squared distance from centroids inertia

 $\frac{1}{n} \sum_{j=1}^{K} \sum_{x \in C_j} distance(x, centroid_j)^2$

- Plot inertia for different values of K and look for optimum
- Can also use change in inertia threshold to stop iterations



Efficient — each iteration makes a single pass over data



- Efficient each iteration makes a single pass over data
 - Incrementally compute centroid

Disadvantages

▶ < ∃ ▶</p>

э

- Efficient each iteration makes a single pass over data
 - Incrementally compute centroid

Disadvantages

 Can only find clusters that look like ellipses



- Efficient each iteration makes a single pass over data
 - Incrementally compute centroid

Disadvantages

 Can only find clusters that look like ellipses



- Efficient each iteration makes a single pass over data
 - Incrementally compute centroid

Disadvantages

- Can only find clusters that look like ellipses
- Choice of initial random centroid matters
 - Repeat and check



- Efficient each iteration makes a single pass over data
 - Incrementally compute centroid

Disadvantages

- Can only find clusters that look like ellipses
- Choice of initial random centroid matters
 - Repeat and check



- K Means clustering can only find clusters that look like ellipses
- Instead, build clusters bottom up, by merging clusters

< 3

э

- K Means clustering can only find clusters that look like ellipses
- Instead, build clusters bottom up, by merging clusters
- Initially, each item is a singleton cluster
- At each step, merge nearest clusters



- K Means clustering can only find clusters that look like ellipses
- Instead, build clusters bottom up, by merging clusters
- Initially, each item is a singleton cluster
- At each step, merge nearest clusters
- Can represent process using a tree dendrogram
- Choose appropriate level in dendrogram for final clustering



To merge clusters, define distance between clusters

< E

э

To merge clusters, define distance between clusters

- Single link: distance between closest points
 - Creates chain effect



To merge clusters, define distance between clusters

- Single link: distance between closest points
 - Creates chain effect
- Complete link: maximum of pairwise distances



To merge clusters, define distance between clusters

- Single link: distance between closest points
 - Creates chain effect
- Complete link: maximum of pairwise distances
- Average link: mean of pairwise distances





To merge clusters, define distance between clusters

- Single link: distance between closest points
 - Creates chain effect
- Complete link: maximum of pairwise distances
- Average link: mean of pairwise distances
- All require $O(n^2)$ computation expensive



- How to identify odd shaped clusters?
- Cluster group of points that are "close together"
- Identify "dense" neighbourhoods
- How do we formalize this?



- How to identify odd shaped clusters?
- Cluster group of points that are "close together"
- Identify "dense" neighbourhoods
- How do we formalize this?



- Construct a small ball around each point, radius *Eps*
- Identify a threshold for neighbours within ball, *MinPts*



- Construct a small ball around each point, radius *Eps*
- Identify a threshold for neighbours within ball, *MinPts*



< E

- Construct a small ball around each point, radius *Eps*
- Identify a threshold for neighbours within ball, *MinPts*
- Core point has at least MinPts neighbours inside Eps ball



- Construct a small ball around each point, radius *Eps*
- Identify a threshold for neighbours within ball, *MinPts*
- Core point has at least MinPts neighbours inside Eps ball
- Connect each core point to all its neighbours
 Directed edges

- Construct a small ball around each point, radius *Eps*
- Identify a threshold for neighbours within ball, *MinPts*
- Core point has at least MinPts neighbours inside Eps ball
- Connect each core point to all its neighbours
- Border points attached to core points but not core themselves



- Construct a small ball around each point, radius *Eps*
- Identify a threshold for neighbours within ball, *MinPts*
- Core point has at least MinPts neighbours inside Eps ball
- Connect each core point to all its neighbours
- Border points attached to core points but not core themselves
- Noise isolated, disconnected points



- Formally, edges from core points to neighbours define a directed graph
- Border points are part of this graph, but cannot add edges to extend the graph
- Discard the edge directions
- Connected components are clusters





Implementation of density based
 Clustering available in Python and R

▶ < ∃ ▶</p>

э

Dbscan

- Implementation of density based
 pclustering available in Python and R
- Smaller value of *Eps* subdivides into small clusters

eps=0.05, min_samples=5



Dbscan

- Implementation of density based pclustering available in Python and R
- Smaller value of *Eps* subdivides into small clusters
- Larger *Eps* groups larger clusters

eps=0.05, min_samples=5



eps=0.20, min_samples=5



Lecture 13: 11 March, 2025

Outliers are anomalous values

э

- Outliers are anomalous values
- K Means lie outside natural clusters, far from all centroids



- Outliers are anomalous values
- K Means lie outside natural clusters, far from all centroids
 - But outliers can distort the clustering process



- Outliers are anomalous values
- K Means lie outside natural clusters, far from all centroids
 - But outliers can distort the clustering process
- Density based clustering not connected to any core point
 - But density is applied uniformly





Lecture 13: 11 March, 2025

- Outliers are anomalous values
- K Means lie outside natural clusters, far from all centroids
 - But outliers can distort the clustering process
- Density based clustering not connected to any core point
 - But density is applied uniformly
- How to identify outliers before clustering?





Lecture 13: 11 March, 2025

 An outlier is less dense than its nearest neighbours

▶ < ∃ ▶</p>

э

- An outlier is less dense than its nearest neighbours
- But difference in density may be local



- An outlier is less dense than its nearest neighbours
- But difference in density may be local
- A distance metric to eliminate o₂ could make all of C₁ outliers



- An outlier is less dense than its nearest neighbours
- But difference in density may be local
- A distance metric to eliminate o₂ could make all of C₁ outliers
- C_1 has 400 points, C_2 has 100 points
- Larger distance would make all of C₂ outliers with respect to C₁



 For clustering, we defined a radius *Eps* and looked for *MinPts* neighbours within that ball



э

- For clustering, we defined a radius *Eps* and looked for *MinPts* neighbours within that ball
- Instead, fix *MinPts* and find smallest ball with that many neighbours



- For clustering, we defined a radius *Eps* and looked for *MinPts* neighbours within that ball
- Instead, fix *MinPts* and find smallest ball with that many neighbours
- Compare radius(p) with radius of its neighours



- For clustering, we defined a radius *Eps* and looked for *MinPts* neighbours within that ball
- Instead, fix *MinPts* and find smallest ball with that many neighbours
- Compare radius(p) with radius of its neighours
- A is an outlier because its radius is much more than that of its neighbours



Local outlier factor LOF(p)

Mean radius of MinPts-neighbours(p) radius(p)

- The smaller this ratio, the more likely that p is an outlier
- Comparison is local to neighbourhood, so this can deal with different densities across range of data

