

Sept 27<sup>th</sup>

## HIERARCHICAL CLUSTERING

OUTPUT: dendrogram / tree of clusters, showing closest related points.

NOTE: we can cut at a specific threshold to generate different clusters.

- \* agglomerative  $\rightarrow$  every point in own cluster, each step, merge two closest cluster, stop when every point in own cluster
- \* divisive  $\rightarrow$  keep splitting until every point in own cluster.

## AGGLOMERATIVE CLUSTERING ALSO

1. each point in dataset be in its own cluster
2. compute the distance b/w all pair of clusters
3. merge 2 closest clusters
4. Repeat 3 and 4 until all points are in the same cluster

cluster distance options  
b/w centroid  
shortest path.

NOTE:  $d \rightarrow$  b/w points  
 $D \rightarrow$  b/w clusters.

## SINGLE LINK DISTANCE

min of all pairwise distance b/w points from one cluster and one from another cluster.

$$D_{SL}(C_1, C_2) = \min(d(p_1, p_2) \mid p_1 \in C_1, p_2 \in C_2)$$

NOTE: depends on choice d.

\*  $\checkmark$  can handle clusters of diff size

\*  $\times$  does not do well w/ overlaps

\*  $\times$  does not handle noise

## COMPLETE LINK DISTANCE

max of all pairwise distance b/w points from one cluster and one from another cluster.

$$D_{CL}(C_1, C_2) = \max(d(p_1, p_2) \mid p_1 \in C_1, p_2 \in C_2)$$

NOTE: depends on choice d.

\*  $\checkmark$  can handle noise

\*  $\times$  doesn't do well w/ overlaps

\*  $\checkmark$  does well w/ noise

## AVERAGE LINK DISTANCE

avg of all pairwise distance b/w a points from one cluster and one from another cluster.

$$D_{AL}(C_1, C_2) = \frac{1}{|C_1| \cdot |C_2|} (d(p_1, p_2) \mid p_1 \in C_1, p_2 \in C_2)$$

NOTE: depends on choice d.

## CENTROID DISTANCE

$$D_C(C_1, C_2) = d(\mu_1, \mu_2)$$

## WARD'S DISTANCE

diff b/w spread/var of points merged cluster + unmerged clusters:

$$D_{WD}(C_1, C_2) = \sum_{p \in C_2} d(p, \mu_{12}) - \sum_{p_1 \in C_1} d(p_1, \mu_1) - \sum_{p_2 \in C_2} d(p_2, \mu_2)$$

## DENSITY BASED CLUSTERING

GOAL: cluster together points that are densely packed together

define density: given fixed radius  $\epsilon$ , if it has more than a min number of points, we can categorize it as dense.

DBSCAN - looks at points at core + perimeter (border) or noise point  
 $\epsilon$ -neighborhood of core  
in  $\epsilon$  neighborhood w/min pts  
neither of previous.

create clusters by connecting data points

\* handles oddly shaped clusters

## DBSCAN Algo

given:  $\epsilon$ , min-pt

1. Find  $\epsilon$  neighborhood of each pt
2. label core if  $\geq$  min-pt
3. NOT core, but in neighborhood = border
4. else noise
5. for each core, assign to same cluster all core points in its neighborhood
6. assign border points to nearby clusters.

$\times$  limit is varying density

$\times$  created of same density.

$\times$  highly problematic in high-dim space.

\* can identify clusters of diff shape/size

\* resistant to noise