### **Unsupervised** Learning and Clustering

- In unsupervised learning you are given a data set with no output classifications (labels)
- Clustering is an important type of unsupervised learning
  - PCA was another type of unsupervised learning
- The goal in clustering is to find "natural" clusters (classes) into which the data can be divided a particular breakdown into clusters is a *clustering* (aka grouping, partition)
- How many clusters should there be (k)? Either user-defined, discovered by trial and error, or automatically derived
- Example: Taxonomy of the species one correct answer?
- Generalization After clustering, when given a novel instance, we just assign it to the most similar cluster



CS 270 - Clustering

# Clustering

- How do we decide which instances should be in which cluster?
- Typically put data which is "similar" into the same cluster
  - Similarity is measured with some distance metric
- Also try to maximize between-class dissimilarity
- Seek balance of within-class similarity and between-class dissimilarity
- Similarity Metrics
  - Euclidean Distance most common for real valued instances
  - Can use (1,0) distance for nominal and unknowns like with k-NN
  - Can create arbitrary distance metrics based on the task
  - Important to normalize the features

# **Outlier** Handling

#### • Outliers

- noise, or
- correct, but unusual data

#### • Approaches to handle them

- become their own cluster
  - Problematic, e.g. when k is pre-defined (How about k = 2 above)
  - If k = 3 above then it could be its own cluster, rarely used, but at least it doesn't mess up the other clusters
  - Could remove clusters with 1 or few elements as a post-process step
- Absorb into the closest cluster
  - Can significantly adjust cluster radius, and cause it to absorb other close clusters, etc. – See above case
- Remove with pre-processing step
  - Detection non-trivial when is it really an outlier?
  - Unsupervised Outlier detection algorithms available

#### **Distances Between Clusters**

- Easy to measure distance between instances (elements, points), but how about the distance of an instance to another cluster or the distance between 2 clusters
- Can represent a cluster with
  - Centroid cluster mean
    - Then just measure distance to the centroid
  - Medoid an actual instance which is most typical of the cluster (e.g. Medoid is point which would make the average distance from it to the other points the smallest)

• Other common distances between two Clusters *A* and *B* 

- Single link Smallest distance between any 2 points in *A* and *B*
- Complete link Largest distance between any 2 points in A and B
- Average link Average distance between points in A and points in B

#### **Distances Between Clusters**

• Other common distances between two Clusters *A* and *B* 

- Single link Smallest distance between any 2 points in A and B
- Complete link Largest distance between any 2 points in A and  $B^{-1}$
- Average link Average distance between points in A and points in B



# **Hierarchical and Partitional Clustering**

- Two most common high-level approaches
- Hierarchical clustering is broken into two approaches
  - Agglomerative: Each instance is initially its own cluster. Most similar instance/clusters are then progressively combined until all instances are in one cluster. Each level of the hierarchy is a different clustering/grouping of clusters. (e.g. HAC)
  - Divisive: Start with all instances as one cluster and progressively divide until all instances are their own cluster.
  - You then decide which clustering you want to output.
- With partitional clustering the algorithm creates one clustering of the data (with multiple clusters), typically by minimizing some objective function (e.g. K-means)
  - Note that you could run the partitional algorithm again in a recursive fashion on all the new clusters if you want to build a hierarchy

# Hierarchical Agglomerative Clustering (HAC)

- Input is an *n* × *n* adjacency matrix giving the distance between each pair of instances
- Initialize each instance to be its own cluster
- Repeat until there is just one cluster containing all instances
  - Merge the two "closest" remaining clusters into one cluster
- Each iteration gives a potential clustering. Choose which clustering you want from this set.
- HAC varies based on "Closeness definition"
  - single, complete, average, and ward link distances common

## **Dendrogram Representation**

Item Α В С D E 3 Α 0 1 2 2 2 4 3 В 1 0 2 5 C 2 0 1 2 4 3 D 0 5 E 3 3 3 0

B

D

C

 $(\mathbf{A})$ 

E

#### Standard HAC

- Input is an adjacency matrix
- Output can include a dendrogram which visually shows clusters and merge distance



## Ward Linkage

Ward linkage measures variance of clusters. The distance between two clusters, A and B, is how much the sum of squares distance from each point to its centroid would increase if we merged them. Merge the two clusters with minimum increase.



# HAC Linkage Comparisons

- Single link (nearest neighbor) can lead to long chained clusters where some points are quite far from each other
- Complete link (farthest neighbor) finds more compact clusters
- Average link Used less because have to re-compute the average each time
- Ward linkage Can often be the most suitable method for quantitative features
- Once you have distances between clusters, always merge the closest clusters in terms of the distance



CS 270 - Clustering

# Linkage Methods



# HAC \*Challenge Question\*

- For the data set below show 2 iterations (from 4 clusters until 2 clusters remain) for HAC complete link (furthest).
  - Repeat: Merge the 2 nearest clusters (using complete link distance)
  - Use Manhattan distance
  - Show the dendrogram, including properly labeled distances on the vertical-axis of the dendrogram

Pattern	x	у
a	.8	.7
b	0	0
С	1	1
d	4	4



### HAC Homework

- For the data set below show all iterations (from 5 clusters until 1 cluster remaining) for HAC single link.
  - Show work
  - Use Manhattan distance
  - In case of ties go with the cluster containing the least alphabetical instance.
  - Show the dendrogram, including properly labeled distances on the vertical-axis of the dendrogram.

Pattern	x	у
а	.8	.7
b	1	.2
С	.9	.8
d	0	.2
е	.2	.1

CS 270 - Clustering

# HAC Summary

- Complexity Relatively expensive algorithm
  - $n^2$  space for the adjacency matrix
  - $mn^2$  time for the execution where *m* is the number of algorithm iterations, since we have to compute new distances at each iteration. *m* is usually  $\approx n$  making the total time  $n^3$  (can be  $n^2 \log n$ with priority queue for distance matrix, etc.)
  - All  $k (\approx n)$  clusterings returned in one run. No restart for different *k* values. Must then decide which clustering you want.
- Divisive Starts with all the data in one cluster
  - One approach is to compute the MST (minimum spanning tree  $n^2$  time since it's a fully connected graph) and then divide the cluster at the tree edge with the largest distance similar time complexity as HAC, different clusterings obtained
  - Could be more efficient than HAC if we want just a few clusters

### Which cluster level to choose?



- Depends on goals
  - May know beforehand how many clusters you want or at least a range (e.g. 2-10)
  - Could analyze the dendrogram and data after the full clustering to decide which subclustering level is most appropriate for the task at hand
  - Could use automated *cluster validity* metrics to help
- Could do stopping criteria during clustering

#### **Automated Clustering Scores**

#### What would give us higher scores?

## **Cluster Validity Metrics - Compactness**

- One good goal is *compactness* members of a cluster are all similar and close together
  - One measure of compactness of a cluster is the sum of squares distance from each point to its cluster centroid (aka SSE)

$$Comp(C) = \mathop{\bigotimes}_{i=1}^{|X_c|} (\mathbf{c} - \mathbf{x}_i)^2$$

- where **c** is the centroid of a cluster *C*, made up of instances  $X_c$ . Lower is better.
- The overall compactness of a particular clustering is the sum of the compactness of the individual clusters
- Gives us a numeric way to compare different clusterings by seeking clusterings which minimize the compactness metric
- However, for compactness, what clustering is always best?

### **Cluster Validity Metrics - Separability**

- Another good goal is *separability* members of one cluster are sufficiently different from members of another cluster (cluster dissimilarity)
  - One measure of the separability of two clusters is their squared distance. The bigger the distance the better.
  - $dist_{ij} = (\mathbf{c}_i \mathbf{c}_j)^2$  where  $\mathbf{c}_i$  and  $\mathbf{c}_j$  are two cluster centroids
  - For a clustering which cluster distances should we compare?

### **Cluster Validity Metrics - Separability**

- Another good goal is *separability* members of one cluster are sufficiently different from members of another cluster (cluster dissimilarity)
  - One measure of the separability of two clusters is their squared distance. The bigger the distance the better.
  - $dist_{ij} = (\mathbf{c}_i \mathbf{c}_j)^2$  where  $\mathbf{c}_i$  and  $\mathbf{c}_j$  are two cluster centroids
  - For a clustering which cluster distances should we compare?
  - For each cluster we add in the distance to its closest neighbor cluster

Separability = 
$$\mathop{\overset{|C|}{a}}_{i=1}^{|C|} \min_{j} dist_{ij}(\mathbf{c}_{i},\mathbf{c}_{j})$$

We would like to find clusterings where separability is maximized
However, separability is usually maximized when there are are very few clusters

squared distance amplifies larger distances

- We want techniques that find a balance between inter-cluster similarity and intra-cluster dissimilarity
- Silhouette is one good popular approach
- Scores any clustering with an arbitrary number of unique clusters. Clustering can come from any clustering algorithm.
- a(i) = average dissimilarity of instance *i* to all other instances in the cluster to which *i* is assigned Want it small
  - Dissimilarity could be Euclidian distance, etc.
- b(i) = the smallest average dissimilarity of instance i to instances in other clusters Want it large
- b(i) is smallest for the best different cluster that i could be assigned to the best cluster that you would move i to if needed



$$s(i) = \begin{bmatrix} 1 - a(i) / b(i) & \text{if } a(i) < b(i) \\ 0 & \text{if } a(i) = b(i) \\ 1 & b(i) / a(i) - 1 & \text{if } a(i) > b(i) \end{bmatrix}$$

- *s*(*i*) is 0 when *i* is right on the border between two clusters
- *s*(*i*) is negative when *i* probably belongs in another cluster
- By definition, s(i) = 0 if it is the only node in the cluster
- The quality of a single cluster can be measured by the average silhouette score of its members, (close to 1 is best)
- The quality of a total clustering can be measured by the average silhouette score of all the instances
- To find best clustering, compare total silhouette scores across clusterings with different *k* values and choose the highest

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

 $-1 \pm s(i) \pm 1$ 

Visualizing Silhouette



Silhouette analysis for KMeans clustering on sample data with n\_clusters = 2

Width is quality of cluster, height is size of cluster Dashed line is average silhouette score for clustering

CS 270 - Clustering





• Best case graph for silhouette?

- Best case graph for silhouette?
- Clusters are wide Scores close to 1
- Not many small silhouette instances
- Depending on your goals:
  - Clusters are similar in size
  - Cluster size and/or number are close to what you want

- Could just use total silhouette average to decide best clustering but best to do silhouette analysis with a visualization tool and use score along with other aspects of the clustering
  - Cluster sizes
  - Number of clusters
  - Shape of clusters
  - Etc.
- Note when task dimensions are > 3 (typical and no longer visualizable for us), silhouette graph still easy to visualize
- $O(n^2)$  complexity due to b(i) computation
- There are other cluster metrics out there
- These metrics are rough guidelines and should be "taken with a grain of salt"

#### Silhouette Homework

Assume a clustering with {a,b} in cluster 1 and {c,d,e} in cluster 2. What would the Silhouette score be for a) each instance, b) each cluster, and c) the entire clustering. d) Sketch the Silhouette visualization for this clustering. Use Manhattan distance for your distance calculations.

Pattern	x	у
а	.8	.7
b	.9	.8
С	.6	.6
d	0	.2
е	.2	.1

#### K-means

- Perhaps the most well-known clustering algorithm
  - Partitioning algorithm
  - Must choose a *k* beforehand
  - Thus, typically try a spread of different k's (e.g. 2-10), run multiple times, and compare to decide which is the best clustering
    - Could use cluster validity metrics (e.g. Silhouette) to help in the decision
- 1. Randomly choose k instances from the data set to be the initial k centroids
- 2. Repeat until no (or negligible) more changes occur
  - a) Group each instance with its closest centroid
  - b) Recalculate the centroid based on its new cluster
- Time complexity is O(mkn) where *m* is # of iterations and space is O(n), both much better than HAC time and space  $(n^3 \text{ and } n^2)$

# K-means Example



#### **K-means** Continued

- Type of EM (Expectation-Maximization) algorithm, Gradient descent
  - Can struggle with local minima, unlucky random initial centroids, and outliers
    - K-medoids finds medoid (median) centers rather than average centers and is thus less effected by outliers
  - Local minima, empty clusters: Can just re-run with different initial centroids
    - Could compare different solutions *for a specific k value* by seeing which clusterings minimize the overall SSE to the cluster centers (i.e. compactness), or use silhouette, etc.
    - And test solutions with different *k* values using Silhouette or other metric
- Can do further refinement of HAC results using any *k* centroids from HAC as starting centroids for *k*-means

## **\*\* K-means Challenge Question \*\***

- For the data below, show the centroid values and which instances are closest to each centroid *after* centroid calculation for two iterations of K-means using Manhattan distance
- By 2 iterations I mean 2 centroid changes after the initial centroids
- Assume k = 2 and that the first two instances are the initial centroids

Pattern	<i>x</i>	У	Iteration	Centroid 1 and instances	Centroid 2 and instances
a	3	1	0		
b	4	0	0		
	0	0	1		
С	0	0	2		
d	0	1	Z		

1.Repeat until no more changes occur
 a)Group each instance with its closest centroid
 b)Recalculate the centroid based on its new cluster

CS 270 - Clustering

### **\*\* K-means Challenge Question \*\***

- For the data below, show the centroid values and which instances are closest to each centroid *after* centroid calculation for two iterations of K-means using Manhattan distance
- By 2 iterations I mean 2 centroid changes after the initial centroids
- Assume k = 2 and that the first two instances are the initial centroids

Pattern	x	y
а	3	1
b	4	0
С	0	0
d	0	1

Iteration	Centroid 1 and instances	Centroid 2 and instances
0	$3, 1 \{a, c, d\}$	4,0 { <i>b</i> }
1	$1, 2/3 \{c, d\}$	4,0 { <i>a</i> , <i>b</i> }
2	$0, .5 \{c, d\}$	3.5, .5 { <i>a</i> , <i>b</i> }
3	$0, .5 \{c, d\}$	3.5, .5 { <i>a</i> , <i>b</i> }

#### K-means Homework

- For the data below, show the centroid values and which instances are closest to each centroid *after* centroid calculation for two iterations of K-means using Manhattan distance
- By 2 iterations I mean 2 centroid changes after the initial centroids
- Assume k = 2 and that the first two instances are the initial centroids

Pattern	x	у
а	.9	.8
b	.2	.2
С	.7	.6
d	1	6
е	.5	.5

# **Other Unsupervised Models**

- Vector Quantization Discretize into codebooks
- K-medoids
- Conceptual Clustering (Symbolic AI) Cobweb, Classit, etc.
  - Incremental vs Batch
- Density based clustering, DBSCAN, OPTICS
- Outlier detection models
- Special models for large data bases  $-n^2$  space?, disk I/O
  - Sampling Bring in enough data to fill memory and then cluster
  - Once initial prototypes found, can iteratively bring in more data to adjust/fine-tune current prototypes as desired
  - Linear algorithms

# Semi-Supervised Learning Examples



Combine labeled and unlabeled data with assumptions about typical data to find better solutions than just using the labeled data

#### Summary

- Standard clustering highly used
- Can also use clustering as a discretization technique on continuous data for many other models which favor nominal or discretized data
  - Including supervised learning models (Decision trees, Naïve Bayes, etc.)
- With so much (unlabeled) data out there, opportunities to do unsupervised learning are growing
  - Semi-Supervised learning is becoming very important
  - Use unlabeled data to augment the more limited labeled data to improve accuracy of a supervised learner
- Deep Learning Unsupervised training of early layers is an important approach in some deep learning models

# **Clustering Project**

• Last individual project

# Neural Network Clustering



- Single layer network
  - Bit like a chopped off RBF, where prototypes become adaptive output nodes

y

 $\bigotimes_{\gamma}$ 

Х

- Arbitrary number of output nodes (cluster prototypes) User defined
- Locations of output nodes (prototypes) can be initialized randomly
  - Could set them at locations of random instances, etc.
- Each node computes distance to the current instance
- Competitive Learning style winner takes all closest node decides the cluster during execution
- Closest node is also the node which usually adjusts during learning
- Node adjusts slightly (learning rate) towards the current example

# **Neural Network Clustering**

y

 $\otimes$ 

Х



- Could start with more nodes than probably needed and drop those that end up representing none or few instances
  - Could start them all in one spot However...
- Could dynamically add/delete nodes
  - Local vigilance threshold
  - Global vs local vigilance
  - Outliers

# **Example Clusterings** with Vigilance

Example ART II Classifications





# **Self-Organizing Maps**

- Output nodes which are close to each other represent similar classes – Biological plausibility
- Neighbors of winning node also update in the same direction (scaled by a learning rate), as the winner
- Self organizes to a topological class map (e.g. vowel sounds)
  - Can interpolate, k value less critical, different 2 or 3-dimensional topologies



CS 270 - Clustering

# Association Analysis – Link Analysis

- Used to discover relationships/rules in large databases
- Relationships represented as *association rules* 
  - Unsupervised learning, can give significant business advantages, and also good for many other large data areas: astronomy, etc.
- One example is *market basket analysis* which seeks to understand more about what items are bought together
  - This can then lead to improved approaches for advertising, product placement, etc.
  - Example Association Rule:  $\{Cereal\} \Rightarrow \{Milk\}$

Transaction ID and Info	Items Bought
1 and (who, when, etc.)	{Ice cream, milk, eggs, cereal}
2	{Ice cream}
3	{milk, cereal, sugar}
4	{eggs, yogurt, sugar}
5	{Ice cream, milk, cereal}