

# Scientific Computing for Biologists

## Lecture 9: (Dis)similarity and clustering

Instructor: Paul M. Magwene

30 October 2012

# Outline of Lecture

- Distance and dissimilarity measures
  - Quantitative data
  - Dichotomous data
  - Qualitative data
- Hierarchical clustering
- Neighbor-joining
- Minimum Spanning Tree (MST)
- K-means clustering

# Similarity/Dissimilarity

## Intuition

Similarity is a measure of “likeness” between two entities of interest. Dissimilarity is the complement of similarity.

- Dissimilarities may be converted to similarities (and vice versa) by taking any monotonically decreasing function. For example:

$$s = 1 - d_{ij} \text{ (for } 0 \leq d_{ij} \leq 1)$$

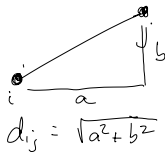
- Dissimilarities are usually in range  $0 \leq d_{ij} \leq C$  where  $C$  is the maximum dissimilarity
- Distances are one measure of dissimilarity but distances are unbounded to the right

$$d_{ij} \in [0, \infty]$$

## Dissimilarity Measures for Quantitative Data

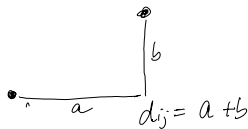
- Euclidean Distance

$$d_{ij} = \left\{ \sum_{k=1}^p (x_{ik} - x_{jk})^2 \right\}^{1/2}$$



- Manhattan (taxi-cab) distance

$$d_{ij} = \sum_{k=1}^p |x_{ik} - x_{jk}|$$



- Scaled Euclidean Distance

$$d_{ij} = \left\{ \sum_{k=1}^p w_k^2 (x_{ik} - x_{jk})^2 \right\}^{1/2}$$

where  $w_k$  are suitable weights

e.g.  $(\text{std. dev of variable } k)^{-1}$  or  $(\text{range of } k\text{th variable})^{-1}$

## Metric vs. Non-metric

A non-negative function,  $g(x, y)$ , is metric if:

i) Satisfies the triangle inequality:

$$g(x, y) \leq g(x, z) + g(y, z)$$

ii) Symmetric:

$$g(x, y) = g(y, x)$$

iii)  $g(x, y) = 0$  only if  $x = y$

Euclidean Dist. is a metric function  
(as is Manhattan distance)

## Other Quantitative Measures of Dissimilarity

- Minkowski Metric

$$d_{ij} = \left\{ \sum_{k=1}^p |x_{ik} - x_{jk}|^{\lambda} \right\}^{1/\lambda} \quad \text{for integers } \lambda$$

$\lambda=1$  is Manhattan distance,  $\lambda=2$  is Euclidean Dist.

- Canberra Metric

$$d_{ij} = \sum_{k=1}^p \frac{|x_{ik} - x_{jk}|}{(x_{ik} + x_{jk})}$$

[Accts for distance b/w.  
points & relationship to  
origin  
→ only for non-negative  
values]

- Czekanowski Coefficient

$$d_{ij} = 1 - \frac{2 \sum_{k=1}^p \min(x_{ik}, x_{jk})}{\sum_{k=1}^p (x_{ik} + x_{jk})}$$

[% dissimilarity  
over all variables]

## Quantitative Dissimilarity for Variables

Correlation provides a suitable measure of similarity

$d_{kl} = 1 - r_{kl}$  if  $r_{kl} = -1$  is taken to indicate maximum disagreement

$d_{kl} = 1 - r_{kl}^2$  is appropriate if  $r_{kl} = 1$  and  $r_{kl} = -1$  are treated equivalently (predictive power)

$$d_{kl} = 1 - \frac{\sum_{i=1}^n X_{ik} X_{il}}{\left( \sum_{i=1}^n X_{ik}^2 \sum_{i=1}^n X_{il}^2 \right)^{1/2}} \quad \leftarrow \text{uncentered correlation}$$

## Dissimilarity for Dichotomous Data

For each pair of objects of interest form a  $2 \times 2$  Contingency table

		obj 2	
		+	-
obj 1	+	a	b
	-	c	d

$$a+b+c+d=p$$

$$\text{Simple Matching: } d_{ij} = 1 - \frac{a+d}{p} = \frac{b+c}{p}$$

$$\text{Jaccard Coefficient: } d_{ij} = \frac{b+c}{a+b+c} \quad (\text{joint absence does not contribute})$$

$$\text{Czekanowski's Coeff: } d_{ij} = \frac{b+c}{2a+b+c}$$



# Introduction to Clustering

## Goal of Clustering

- Find "natural groups" in data

→ one definition:

patches of high density surrounded by patches of lower density in the  $p$ -dimensional space defined by the variables

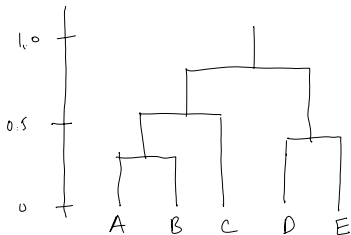


# Hierarchical Clustering

Agglomerative/Divisive methods

- In practice almost always agglomerative

For  $n$  data points define a set of  $n-1$  joins that represent groupings of objects @ different levels of similarity



## Simple Algorithm for Hierarchical Clustering

- 1) Calculate a dissimilarity matrix for the  $n$  items
- 2) Join the two nearest items,  $i$  &  $j$
- 3) Delete the  $i^{\text{th}}$  &  $j^{\text{th}}$  row and column of the dissimilarity matrix; add a new row/column \* that represents dissimilarity of new group  $(i, j)$  to all other items
- 4) Repeat from step 2 until there is a single group

## Methods of Hierarchical Clustering

The different methods are determined by the function used to determine the distance between groups

### Some Common Group Distance Criteria

Single linkage (nearest neighbor)

Complete linkage (furthest neighbor)

Group average

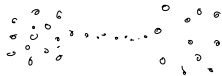
Centroid

## Single Linkage Clustering

$n_i, n_j$  are # of objects in groups  $i$  &  $j$

★  $D_{ij}$  is the smallest of the  $n_i n_j$  dissimilarities between each element of  $i$  & each element of  $j$

- Invariant under monotonic transformation of the  $d_{ij}$
- Unaffected by ties
- Provably nice asymptotic properties
- Susceptible to "chaining"



## Complete Linkage

$D_{ij}$  is the maximum of the  $n_i n_j$  dissimilarities between the two groups

→ also invariant under monotonic transformation

## Group average

$D_{ij}$  is the average of the  $n_i n_j$  dissimilarities between the two groups (UPGMA, WPGMA)

## Centroid method

$D_{ij}$  is the squared euclidean distance between the centroids of groups  $i$  &  $j$

# Hierarchical Clustering, A worked Example

	A	B	C	D	E
A	0				
B	4	0			
C	①	4	0		
D	4	2	4	0	
E	5	5	3	4	0

Single Linkage

	(A,C)	B	D	E
(A,C)				
B	4			
D	4	2	0	
E	3	5	4	0

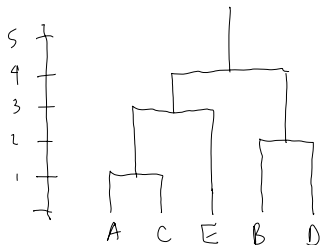
	(A,C)	(B,D)	E
(A,C)	0		
(B,D)	4	0	
E	③	4	0



Worked Example, cont.

$((A,C),E)$   $(B,D)$   
 $((A,C),E)$  0  
 $(B,D)$  4 0

Only one choice  
 $((A,C),E), (B,D)$



Single Linkage Clustering

# Neighbor Joining

Originally described by Saitou and Nei, 1987.

## Goal

Tries to create the (unrooted) tree topology with the least branch length (minimum-evolution criterion).

Basic algorithm:

- 1 Calculate matrix  $Q$  (next slide) from the distance matrix
- 2 Find the pair of taxa in  $Q$  with the lowest value; create a node on the tree that joins these two taxa (i.e. the closest neighbors)
- 3 Calculate the distance of each of the taxa in the pair to this new node
- 4 Calculate the distance of all taxa outside of this pair to the new node
- 5 Repeat from step 1 using the distances calculated in the previous step

## Neighbor Joining, cont.

$$Q_{ij} = (r - 2)d_{ij} - (R_i + R_j)$$

where  $r$  is the number of taxa,  $d_{ij}$  is the distance between taxa  $i$  and  $j$  and  $R_k$  is the row sum over row  $k$  of the distance matrix ( $R_k = \sum_i d_{ik}$ ).

When nodes  $i$  and  $j$  are joined they are replaced by a node,  $A$ , with distance to a remaining node  $k$  given by:

$$d_{Ak} = \frac{1}{2}(d_{ik} + d_{jk} - d_{ij})$$

# NJ example from Saitou and Nei 1987

**Table 1**  
Distance Matrix for the Tree in Figure 1

OTU	OTU						
	1	2	3	4	5	6	7
2	..	7					
3	..	8	5				
4	..	11	8	5			
5	..	13	10	7	8		
6	..	16	13	10	11	5	
7	..	13	10	7	8	6	9
8	..	17	14	11	12	10	13

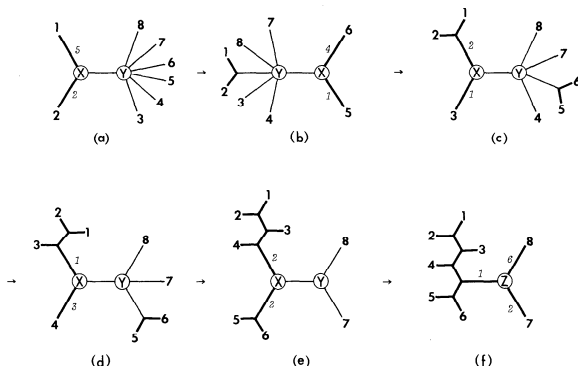


FIG. 3.—Application of the neighbor-joining method to the distance matrix of table 1. Italic numbers are branch lengths, and branches with thicker lines indicate that their lengths have been determined.

# Minimum Spanning Tree

## Goal

Construct a tree that connects all points in the data set and whose total length is minimized.

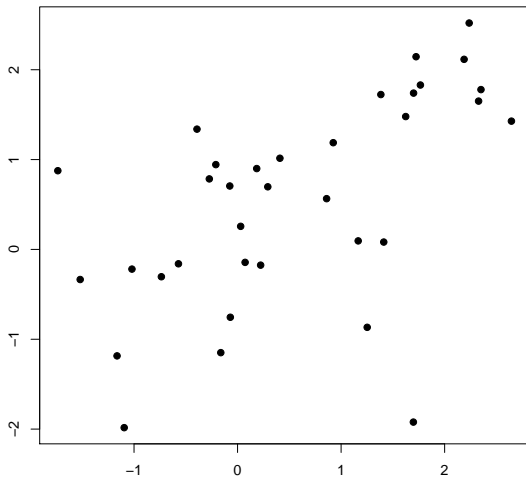
### *Statistical applications*

- highlights close neighbors in a data set
- useful check for distortions produced by projection techniques
- tests of normality

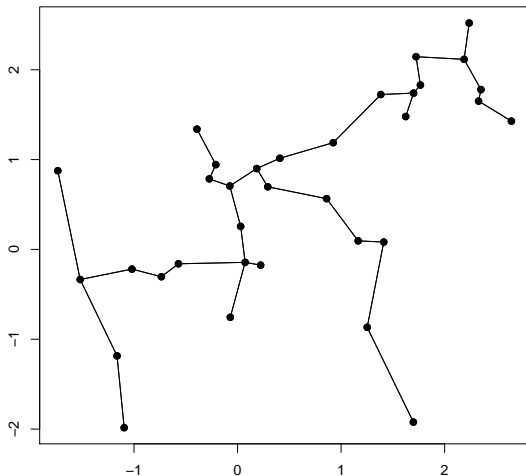
### *Other applications*

- urban planning/engineering
- circuit design

## Example Data Set



# Minimum Spanning Tree: Example



# Relationship between MST and Single Linkage Clustering

- Cut a single linkage dendrogram at height,  $\delta \rightarrow$  clusters
- Remove all edges in the MST with length  $\geq \delta \rightarrow$  subgraphs corresponding to the same clusters



# A Generic MST Algorithm

**Input:** dissimilarity matrix,  $D$ , between each object (point) of interest

- 1 Create a graph,  $G$ , where  $V = \{v_1, \dots, v_n\}$  and  $E = \{\}$  ( $E$  initially empty)
- 2 Find the smallest dissimilarity,  $d_{ij}$  where  $(i,j)$  is not in  $E$ .
- 3 Add  $(i,j)$  to  $E$  if  $(i,j)$  does not create a cycle
- 4 Repeat from step 2 until every vertex is included in at least one edge

Not particularly efficient algorithm, but simple. More efficient algorithms for finding MSTs include Kruskal's Algorithm and Prim's algorithm.

# Applications of the MST

MST tends to highlight close neighbors; can be used to look for distortions associated with projections to lower dimensional spaces.

## Using the MST to look for Projection Distortion

- Calculate the MST based on dissimilarity in a high-dimensional space
- Draw the MST edges among points in the projection space (e.g. MDS or PCA)
- MST edges that cross highlight geometric relationships among points that are not well represented by the projection

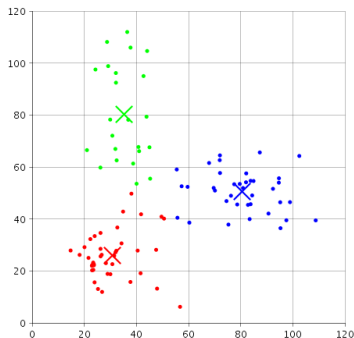
# K-mean Clustering

## General idea

Assign the  $n$  data points (or  $p$  variables) to one of  $K$  clusters to as to optimize some criterion of interest.

- The most common criterion to minimize is the sum-of-squares from the group centroids.

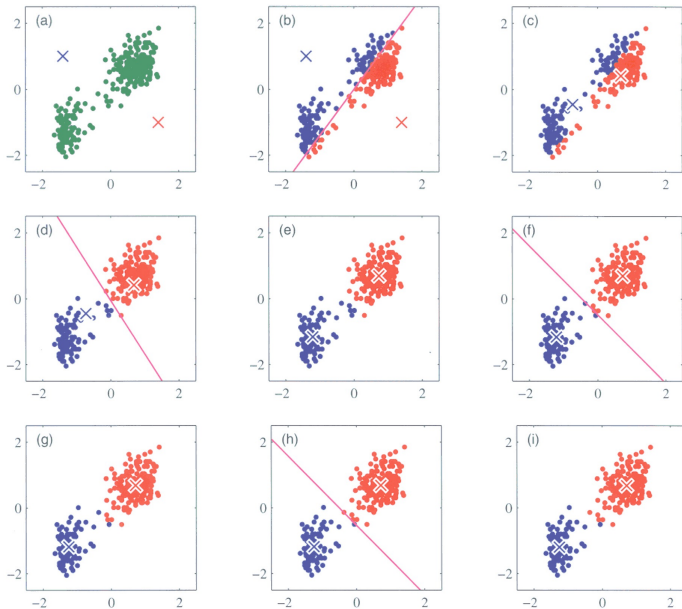
$$V = \sum_{i=1}^k \sum_{j \in g_i} |x_j - \mu_i|^2$$



# Simple algorithm for K-means clustering

- 1 Decide on  $k$ , the number of groups
- 2 Randomly pick  $k$  of the objects to act as the initial centers
- 3 Assign each object to the group whose center it is closest to
- 4 Recalculate the  $k$  centers as the centroids of the objects assigned to them
- 5 Repeat from step 3 until centroids no longer move (convergence)

# Illustration of K-means algorithm



## Things to note re:K-means clustering

- The algorithm described above does not necessarily find the global optimum
- The algorithm is sensitive to choice of initial cluster center; k-means is often run multiple-time with different initial centers to insure inferred clusters are robust.