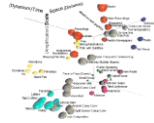


## Module 12

### Cluster Analysis



### Cluster Analysis

- Another tool we can use to simplify a very complex, multivariate database
- However, unlike factor analysis, this method operates specifically in spatial terms: group data (observations) in "space"
  - As we discussed with MDS, "space" can be geographical, or perhaps another kind of space we conceptualize (such as a "perception space" or a "similarity space")

### Cluster Analysis

- Steps in cluster analysis ("CA")
  - 1. take points, areas, or objects ("observations") and measure the "distance" between each pair
  - 2. analyze this distance data to uncover the latent grouping structure embodied in the dataset
- Some kind of measure of similarity/ dissimilarity is needed to do this analysis

### Cluster Analysis

- Purpose of CA: see trends, generate hypotheses (highly exploratory)
  - Important advantage: CA does not need normality or linearity (non-parametric), so cluster analysis can be widely used
  - There are more opportunities for use of CA than actual implementations: CA is a method to be aware of for its potential for innovation

### Cluster Analysis

- Key Idea: cluster analysis usually does not focus on geographic space
  - Clusters are often defined in non-geographic terms: "space" in some other sense
  - Focus is the creation of a classification system: clusters (in this context) = groupings
    - Groupings of people: based on health and lifestyle factors
    - Groupings of forests: based on vegetation types and climate characteristics
    - Groupings of cities: based on major industries or other socioeconomic characteristics

### Cluster Analysis

- Focus here: one specific approach to clustering called hierarchical clustering
  - Hierarchical clustering: provides information on clustering at multiple levels of complexity
    - With hierarchical clustering, possible to cluster a database into 2 groupings, 3 groupings, 4 groupings, etc. (max = # of records in dataset)
    - You don't need to know in advance how many groupings (clusters) you want to produce
    - Hierarchical clustering gives you insight to help you select how many clusters you wish to identify

### Cluster Analysis

- o **Alternative:** another approach called *k-means clustering*
  - **K-means clustering:** efficient method for producing a *specified number* of clusters
    - o "Efficient" in terms of computer run-time
    - o However, with k-means clustering you need to know how many clusters are appropriate for your dataset (or at least, how many you want to see)
    - o You could do k-means clustering multiple times to compare different levels of cluster systems, but that negates its time efficiency

### Cluster Analysis

- o Let's look at a dataset of nine values to see the basic cluster idea (hierarchical)
 

Cell Identifier	A	B	C
	50	20	18
	D	E	F
		3	34
Actual data value	G	H	I
	71	80	86

Imagine each cell value as a data observation for a given geographic area (9 observations for 9 areas in a 3x3 grid)

### Cluster Analysis

- o Let's look at a dataset of nine values to see the basic cluster idea (hierarchical)
  - CA generates a **dendrogram** chart to show the hierarchical structure in this table

A	B	C
50	20	18
D	E	F
7	3	34
G	H	I
71	80	86

These cells link first because of smallest distance (20-18=2)

These cells link next (7-3=4)

And so on ...

### Cluster Analysis

**A slightly more complex dendrogram example**

### Cluster Analysis

**And another... (they can be much bigger yet)**

### Cluster Analysis

**One last dendrogram example from research in business geography**

Joseph, Lawrence (2016) The Geographic Exposure to Lifestyles by U.S. Retail Chains. The Professional Geographer: DOI: 10.1080/00000001.2016.1140487

**Cluster Analysis**

- o **Big idea of the dendrogram:** shows the order/structure of the joins (hierarchy)
  - Most similar cells join first, less similar cells join later, until all cells are joined

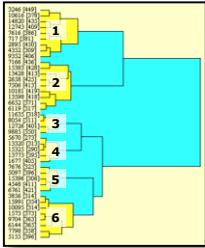
**Cluster Analysis**

- o **Key issue: how many clusters are appropriate?**
  - Let's look at a simple example to gain some insight into this question and how we can address it

**Cluster Analysis**

Here the analyst created six clusters (the six yellow areas)

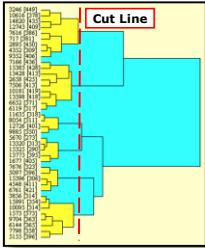
How many clusters are actually best?



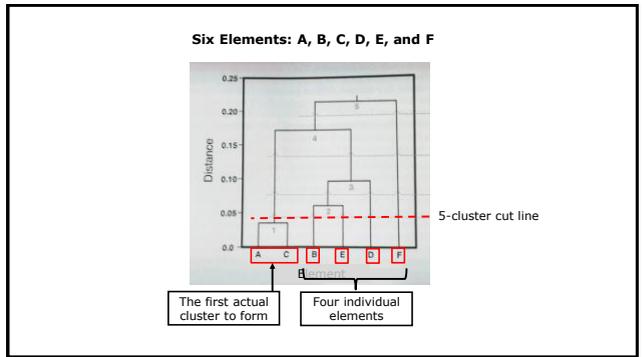
**Cluster Analysis**

The concept of a "cut line" helps here

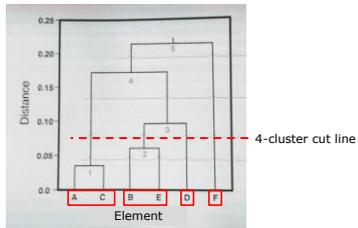
A cut line defines the level at which the user chooses to create clusters



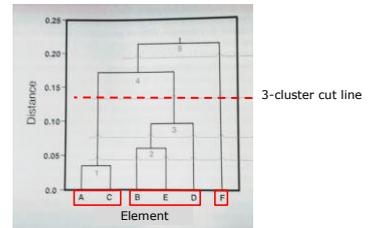
*Here's another simple example that develops the use of a cut line*



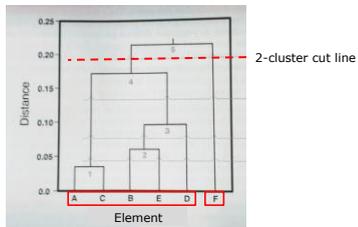
Six Elements: A, B, C, D, E, and F



Six Elements: A, B, C, D, E, and F

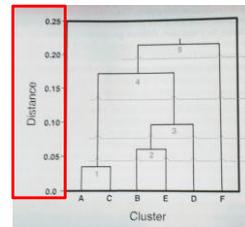


Six Elements: A, B, C, D, E, and F

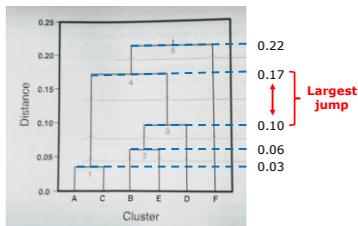


Q: Does this dendrogram give any clues as to what might be the most appropriate number of clusters?

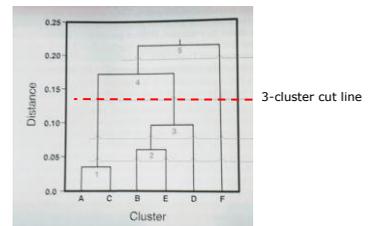
Check out the distance axis: how far each element is apart (and at what distance does each cluster form)



Q: What does the distance information here tell us?



So maybe a 3-cluster solution is best for this situation







### Cluster Analysis: The Details

- Multivariate CA Example: An Agricultural Census**

Group the Counties by Similarity of Farm Outputs

Farm Product	County W	County X	County Y	County Z
Wheat	6	5	10	8
Hay	1	2	3	4
Oats	5	5	1	2

Can't just do a simple subtraction: W-X, or X-Y

### Cluster Analysis: The Details

- Distance Metrics**

- Euclid  $d_{xy} = \sqrt{\sum_i (x_i - y_i)^2}$
- So, for example, between counties X and Y in the previous table

Farm Product	County X	County Y	X-Y	(X-Y) <sup>2</sup>
Wheat	5	10	-5	25
Hay	2	3	-1	1
Oats	5	1	4	16
<b>Total</b>				<b>42</b>

$d_{xy} = \sqrt{42} = 6.5$  So, county X is 6.5 units from county Y

### Cluster Analysis: The Details

- Distance Metrics**

- Squared Euclid
  - Same as Euclid, except squared (duh!)
  - So, in the previous example, county X is **42 units** from county Y (not 6.5)
  - Idea:** puts a much larger penalty on large distances, the groups it identifies tend to be very similar

### Cluster Analysis: The Details

- Distance Metrics**

- Manhattan (or City Block)

$$d_{xy} = \sum_i |x_i - y_i|$$

- In the previous county X/Y example,  $d_{xy} = 10$

Farm Product	County X	County Y	X-Y	X-Y
Wheat	5	10	-5	5
Hay	2	3	-1	1
Oats	5	1	4	4
<b>Total</b>				<b>10</b>

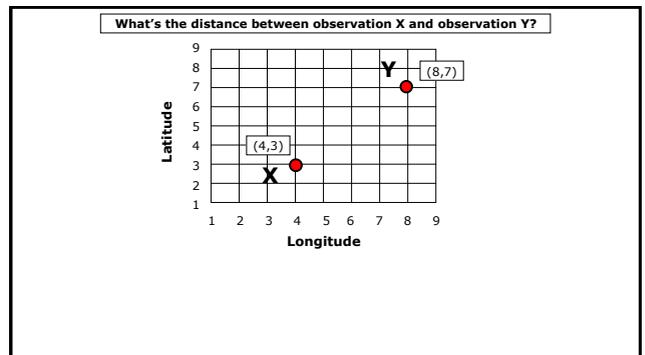
### Cluster Analysis: The Details

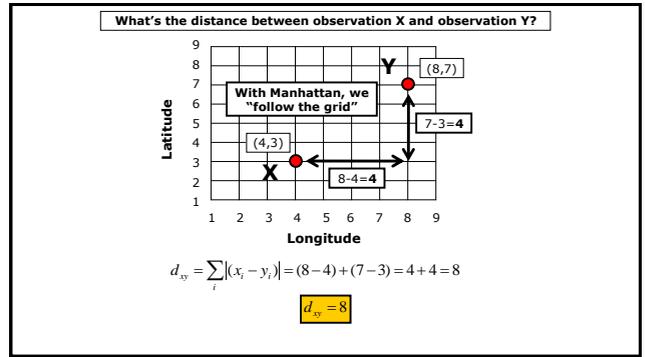
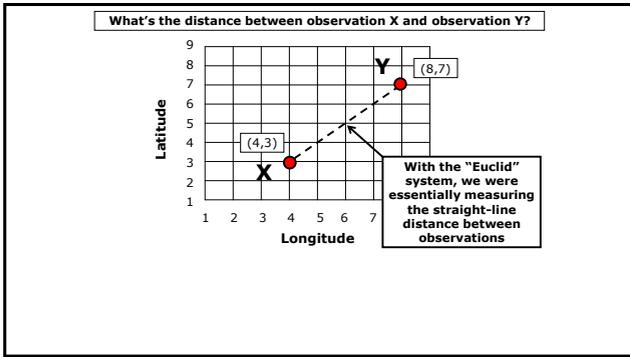
- Distance Metrics**

- Manhattan (or City Block)

$$d_{xy} = \sum_i |x_i - y_i|$$

- Why "Manhattan"?
- Think of a spatial example of clustering: cluster observations based on their spatial coordinates
  - Two variables for each observation: longitude and latitude

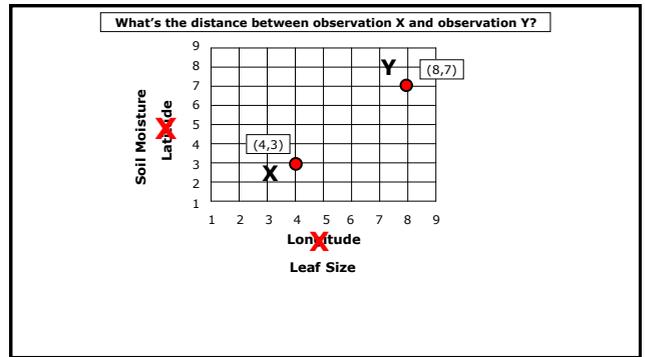




What's the distance between observation X and observation Y?

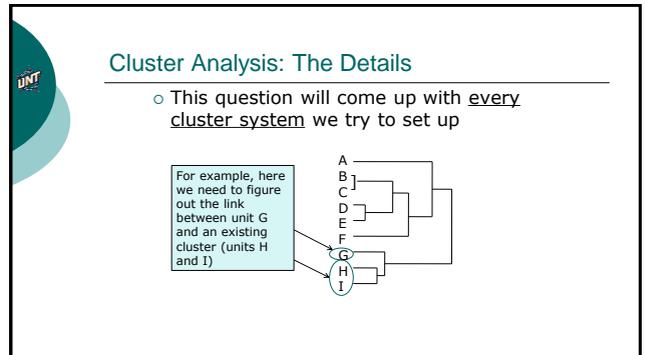
Now that we've defined this distance measurement system, we could cluster **any number of points** (observations) using their "Manhattan" differences as the starting point

Also note that we can use **any number and type of variables** for our distance calculations, not just the **two** longitude and latitude variables used here



**Cluster Analysis: The Details**

- o **Another question:** following from what we just did
  - We now know how to calculate distances between two different, individual units (like counties or census tracts)
  - How do we calculate distances between an individual unit and a cluster that's already been created?
  - And, how do we calculate distance between two existing clusters?



### Cluster Analysis: The Details

- o This question will come up with every cluster system we try to set up

And here we need to figure out the link between the cluster of units B and C and the cluster of units D and E

### Cluster Analysis: The Details

- o **Methods of Defining Distances When Clusters are Involved**
- o 1. Single Linkage (Nearest Neighbors)
  - The smallest distance between a cluster and a cell, or a cluster and another cluster

So, the single linkage distance from "cell G" to "cluster H-I" is **9**

### Cluster Analysis: The Details

- o **Methods of Defining Distances When Clusters are Involved**
- o 1. Single Linkage (Nearest Neighbors)
  - The smallest distance between a cluster and a cell, or a cluster and another cluster

**Q:** What are the single linkage distances here (F to B-C, and F to D-E), and which cluster does F link to?

### Cluster Analysis: The Details

- o **Methods of Defining Distances When Clusters are Involved**
- o 1. Single Linkage (Nearest Neighbors)
  - The smallest distance between a cluster and a cell, or a cluster and another cluster

We can use this methodology to calculate and then compare distances for multiple clusters

### Cluster Analysis: The Details

- o **Methods of Defining Distances When Clusters are Involved**
- o 2. Complete Linkage (Furthest Neighbors)
  - The largest distance between a cluster and a cell, or a cluster and another cluster
  - Pretty much the same idea as what we just went through, except substitute "largest" for "smallest" in your distance calculations

### Cluster Analysis: The Details

- o **Methods of Defining Distances When Clusters are Involved**
- o 2. Complete Linkage (Furthest Neighbors)
  - However: you still cluster based on smallest distances
  - "Largest" only applies to the distance calculation, not the clustering

### Cluster Analysis: The Details

- o **Methods of Defining Distances When Clusters are Involved**
- o 2. Complete Linkage (Furthest Neighbors)
  - Below, what are the complete linkage distances between cell F and the two clusters?

### Cluster Analysis: The Details

- o **Methods of Defining Distances When Clusters are Involved**
- o 3. Centroid Linkage
  - The distance between a cell and a cluster is the difference between the cell value and the average of the cluster
  - The distance between two cells is the difference between their averages

### Cluster Analysis: The Details

- o **Methods of Defining Distances When Clusters are Involved**
- o 3. Centroid Linkage
  - So, for example

### Cluster Analysis: The Details

- o **Methods of Defining Distances When Clusters are Involved**
- o 4. Average Linkage (Between Groups)
  - Examine all pairs of points (cell-cluster, cluster-cluster) in the distance calculation
  - Reduces to the same as centroid linkage when just dealing with a cell-cluster distance calculation

### Cluster Analysis: The Details

- o **Methods of Defining Distances When Clusters are Involved**
- o 4. Average Linkage (Between Groups)

$$d = \frac{\sum diff}{n}$$

- $d$  = the total calculated inter-cluster distance (the actual average linkage)
- $diff$  = individual cell pair differences
- $n$  = number of cell pairs

### Cluster Analysis: The Details

- o **Methods of Defining Distances When Clusters are Involved**
- o 4. Average Linkage (Between Groups)
  - Example

**Cluster Analysis: The Details**

- **Methods of Defining Distances When Clusters are Involved**
- 4. Average Linkage (Between Groups)
  - Example

**Cluster Analysis: The Details**

- **Methods of Defining Distances When Clusters are Involved**
- 4. Average Linkage (Between Groups)
  - Example

**Cluster Analysis: The Details**

- **Methods of Defining Distances When Clusters are Involved**
- 4. Average Linkage (Within Groups)
  - Our last method for discussion is the same as what we just discussed, except that the within group option also takes into account intra-group distances in the calculation

**Cluster Analysis: The Details**

- **Methods of Defining Distances When Clusters are Involved**
- 4. Average Linkage (Within Groups)
  - Symbolically, this means

**Cluster Analysis: The Details**

- **Methods of Defining Distances When Clusters are Involved**
- In general, average linkage methods are best because they use all the data points
  - However, there is a place for all methods in the appropriate circumstance

**Cluster Analysis: The Details**

- **Methods of Defining Distances When Clusters are Involved**
- Other methods you may run across include Ward's, median
  - Some restrictions exist on the use of certain distance methods, depending on the nature of your database
  - But for the ones you've just seen here, no big problems