

ASSIGNMENT 2

(Due on Tuesday, February 15)

As stressed in Assignment 1, your approach to this assignment should be to write a *short report* on your findings for each of these three “studies”. Since this is your first *analytical* assignment, it should also be stressed that each study should be a *pedagogical study* (as in the Example Assignment) that illustrates the application of the relevant analytical method(s). So for example, when applying the *Clark-Evans Test* below, it is not enough to simply display the output from various software programs. You should give a short development of what this method tests, and exactly how it does so. What assumptions are involved? What test statistics are used? What are the advantages and possible limitations of this procedure? What information does it convey about the data set being studied?

(1) In this study you will focus on the type of point pattern analysis done in class for the “Redwood Seedling” data. Here you will use the “Iowa County Seat” data which is displayed in the ARCMAP file, **T:\ese502\arcview\projects\Iowa\Iowa_county_seats.mxd**.¹ Observe first that, unlike the Redwood Seedlings, this point pattern appears to be uniformly dispersed (much like the “Cell Center” data discussed in class). If you open the ‘Iowa counties’ layer, you will see the actual county boundaries. In fact the regularity of this grid-like pattern of county boundaries strongly contributes to the uniformity of county seats. But there is more to the story. Notice that especially for interior counties (away from the Missouri and Mississippi Rivers), each seat is near the center of its county. There are several references here that you should look at when interpreting this pattern. These references offer alternative explanations for landscapes like Iowa:

- The first is the classic discussion of “Central Place Theory” in the *Economics of Location* reference by August Lösch. Using the beer-producing area of Southern Germany as his key example, Lösch describes an idealized agricultural economy in which a hierarchy of towns, or “places”, is regularly spaced on a hexagonal grid of small farms. [Pay particular attention to the highlighted sections]. This is actually quite similar to the case of Iowa, especially during the 1800s. (You can view Iowa as a “Löschian landscape” by closing the ‘Iowa counties’ layer, and opening the ‘Voronoi_cells’ layer above it.) A more general discussion of human settlement patterns is given in the *Locational Analysis* reference by Peter Haggett.
- The second pair of references are by Michael Dacey (*Dacey_1* and *Dacey_2*). Here Dacey offers an alternative probabilistic model of “place” patterns that are more regular than Poisson randomness. This is applied directly to county seats in Iowa (in both references). [Pay particular attention to the highlighted summary sections, *Pattern Summaries*, p.564 in *Dacey_1*, and *Evaluation and Interpretation of Results*, p.540 in *Dacey_2*. See also my notes, *DACEY_model.pdf*, located in the *extra_materials* folder of the class directory, <http://www.seas.upenn.edu/~ese502/lab-content>]

¹ This data is taken from the National Association of Counties website, <http://www.naco.org>.

With this background, we turn now to a point-pattern analysis of Iowa:

(a) First you will construct nearest-neighbors to each point using MATLAB:

1. Load the workspace **Iowa.mat** into MATLAB, and observe that the matrix, **L**, contains the locations of the 99 county seats (with coordinates in State Plane miles).
2. Construct the vector, **nn_dist**, of (first) nearest-neighbor distances in a manner similar to Problem 1(b) of Assignment 1. Use the commands:

```
» OUT = neighbors(L,1);  
» nn_dist = OUT(:,4);  
» nn_dist(1:2,:)
```

The last command is a check, and should yield the values 20.98 and 20.448. (For example the nearest neighbor of county seat 1 is about 21 miles away.)

3. Now save this result to your own directory (say **S:\home**) by the command

```
» save 'S:\home\nn_dist.txt' nn_dist -ascii
```

(b) Next, you will import this data to JMP and examine its properties.

1. Start by opening JMP:

[JMP 7] Use the option “Text Import Preview” to import the file, **nn_dist.txt**, as in Problem 1(b).2 of Assignment 1.

[JMP 8] Reset preferences as in Problem 1(b).2 of Assignment 1 and import the file, **nn_dist.txt**, as a Text file. (Eliminate the empty “Column 1” if necessary.)

Relabel the data column as “nn_dist”. (Check to be sure the first two values are as above.) You will now use this to construct a file which is parallel to **Redwoods_data.jmp**, used for the analysis of the Redwood Seedlings in class.

2. First add a second column, “rand_relabel”, and use the formula: **Random** → **Col Shuffle**. After clicking **Apply** and **OK**, you should see a random relabeling of the 99 row numbers in this column.
3. Next, add a new column labeled “Sample”, and as in Problem 1(b).3 of Assignment 1, construct the subscripted variable:

$nn_dist_{rand_relabel}$

This yields a random shuffling of the nearest-neighbor-distance values. The first 30 elements of this column then constitute the desired subsample of the nearest-neighbor-distance values to be tested.

4. Next, open the JMP file, **CE_Tests.jmp**, and copy-and-paste the first 30 values from “Sample” into the “nn_dist” column of this file.
 5. To determine the appropriate area, open the ARCMAP file, **Iowa.mxd**, and observe from the attribute table of the layer, “Iowa_Boundary” that the last field contains the area of Iowa in square miles (56269 *sq.ml.*). This is the value to be used here. To enter this value into **CE_Tests.jmp**, right click on the “area” column, select **Column Info** → **New Property** → **Formula**, and in the “Formula” box to the right, click **Edit Formula**. Now the **Calculator** window will open, and you can type in the area value.
 6. Next, set the “n” value equal to 99 using the same procedure. (Remember that the *full* sample size is used to estimate point density.) The spreadsheet should now fill in the rest of the values.
 7. Observe that the value “mu” gives the theoretical mean nearest-neighbor distance predicted by the *CSR theory*. Compare this value with the sample mean calculated in the column “s_mean”. What can you conclude from this comparison?
- (c) Finally, to carry out the **Clark-Evans test**, you can simply use the **P-value** columns of the table.
1. For a one-tailed test of “dispersion” versus the null hypothesis of CSR, the appropriate P-value is given by “P-Val Disp”. Interpret this result.
 2. How does this result compare with your findings in 1(b).7 above?
- (d) Now you will repeat this analysis in MATLAB.
1. First load **Iowa.mat** into the workspace (found in **T:\ese502\matlab**). You will see one matrix, **L**, listing the locations of county seats. (*Note:* At this point you might also try the command `>> voronoi(L(:,1),L(:,2))`; in MATLAB to see how it compares with ARCMAP).
 2. Use the program **ce_test.m** (as done in class for **redwoods.mat**) with locations, **L**, area = 56269, and test = 2 (one-tailed test of dispersion):

```
>> ce_test(L,56269,30,2);
```

Compare the resulting **P-value** with that in 1(c).1 above.

3. Finally, use the program **ce_test_distr.m** (as in class) with number of tests, $N = 1000$.

```
>> ce_test_distr(L,56269,30,2,1000);
```

How do these results add to what you have already learned above?

Conclusions. In your concluding discussion of this study, consider how your findings relate to the Central Place Theory of Lösch. Also, how do they relate to the more statistical theory of Dacey?

- (2) In this study (based on Exercise 3.3 in B&G), you will analyze the “Volcanoes in Uganda” data in the ARCMAP file, **T:\ese502\arcview\projects\Volcanoes\Volcanoes.mxd**. There is a key paper that you should look at before analyzing this data, namely the “Volcanoes in Uganda” paper by Tinkler in the Reference Materials. [Pay particular attention to the *Introduction* and *Conclusions* sections.]

- The present volcano data is taken from this study, and it should provide you with useful background material. Try to compare your results with his basic findings about clustering of volcanoes in the Bunyaruguru field.
- Pay particular attention to the notion of “areal volcanism” and how it relates to possible clustering of volcanoes.
- Also look at some of the additional PDF reference files that have been placed in **T:\ese502\arcview\projects\Volcanoes**.

(a) First you will export the point data from ARCMAP to MATLAB. To do so:

1. Open EXCEL (**Programs** → **Office Suites** → **MS Office** → **Excel**) and then open the data base (.dbf) file, **..\Volcanoes\Volcano_pts.dbf** (which is a component of the shape file for **Volcano_pts**).
2. In EXCEL delete the first row of the file (containing the column headings “X” and “Y”) and save the file in your home directory as **volcano_pts.txt** using the save option: “Text (Tab delimited)(*.txt)”.
3. In MATLAB click: **File** → **Import Data** and load the file **volcano_pts.txt**. You will first see the data appear in the **Import Wizard** window. Click **Next** and **Finish**. The data will then appear in the **Workspace** window.

4. Now click: **File** → **Save Workspace** and save the workspace in your home directory as **volcanoes.mat**. This is the workspace you will use for the remainder of the analysis.
 5. To add the **area** of Uganda to this workspace:
 - (i) Open the boundary file **..\Volcanoes\Uganda_bnd.dbf** in EXCEL and right click on the numerical “AREA” value to copy it. (This can also be done using the attribute table for **ugand_bd** in ARCMAP.)
 - (j) In MATLAB type: **» area =** , and the paste in the value. (You could of course copy it by hand, but it is useful to be able to copy-and-paste sections of data.)
 - (k) Now save the workspace again. (There is a *correct* version of this workspace, **volcanoes.mat**, in **T:\ese502\matlab**. So you may wish to compare your version with this workspace to be sure you have done everything correctly.)
- (b) Next you will construct a version of the boundary file for Uganda that is suitable for MATLAB analysis.
1. First open **..\arcview\Volcanoes** again and copy the three (shape) files, **ugand_bd.shp**, **ugand_bd.shx**, **ugand_bd.dbf**, into your home directory.
 2. To transform this shape file into a text-format boundary file suitable for MATLAB, you will use the DOS program, **shp2bnd.exe**. To learn about this program [and also the reverse program, **bnd2shp.exe**, that converts text-format boundary files into shape files] go to the directory **T:\ese502\extra_materials**. If you are working in the MUSA Lab across campus, open the PDF file, **BATCHFILE_for_SHP2BND** All others should open the PDF file, **BOUNDARY_FILES**.
 3. If you are working in the MUSA Lab, then the instructions in the batchfile above will produce a text-format boundary file, **ugand_bd.bnd**, in your directory. Open this file in NOTEPAD to be sure it is correct. The first two lines should be “ 1 396” and “25740 41890”. (Don’t worry if the first row doesn’t line up.) Now change the extension from **.bnd** to **.txt** so that MATLAB can find it more easily, and proceed to step 7 below.
 4. If you are using a computer in the SEAS Lab, proceed to step 5. If you are copying files to use on a stand-alone computer, you can run DOS programs directly on your computer. So copy the three files **SHP2BND.exe**, **BND2SHP.exe**, **BOUNDARY_FILES** from the directory **T:\ese502\extra_materials** into your home directory. (You will not use

BND2SHP.exe in this assignment, but you may use it later.) On your computer, place these files in the same directory as the input data: **ugand_bd.shp, ugand_bd.shx, ugand_bd.dbf** . Now open a DOS window (click: **Programs** → **Command Prompt**), navigate to this directory, and proceed to step 6 below.

5. To open a DOS window in the Lab, click: **Programs** → **Command Prompt**. In the DOS window, navigate to your directory containing the input data: **ugand_bd.shp, ugand_bd.shx, ugand_bd.dbf** . (If the DOS command prompt is **C:\user>** then return to **C:** by typing **cd..**) If your home directory is **S:\home** then at the command prompt, **C:** , type: **S:** , hit return, and type **cd home**.

6. Now start at step (2) of the instructions in **BOUNDARY_FILES** for using **SHP2BND** (“shape-to-boundary”). The command line should read:

```
>> shp2bnd ugand_bd
```

This will produce a text-format boundary file, **ugand_bd.bnd**, in your directory. Open this file in NOTEPAD to be sure it is correct. The first two lines should be “ 1 396” and “25740 41890”. (Don’t worry if the first row doesn’t line up.) Now change the extension from **.bnd** to **.txt** so that MATLAB can find it more easily.

7. Finally, load the file **ugand_bd.txt** into the MATLAB workspace **volcanoes.mat** (**File** → **Import Data** → **S:/home/ugand_bd.txt**). Again, be sure that it appears in the **Workspace** window, and save the workspace.

(c) You are now ready to analyze this volcano point pattern in MATLAB.

1. First, plot out a **simulation envelope** for **99** random patterns by using the commands:

```
» rand('seed',23456);  
» k_function_sim(volcano_pts,area,40,1,99,ugand_bd);
```

The first line sets the random seed value (so that you all should get the same output).

- (i) Before proceeding further, save this figure to your home directory by clicking: **File** → **Save As** (in the Figure window) and saving the figure as **k_function.fig**.

- (j) What can you conclude from this graphical output? Do the results make sense when you compare them with the mapped point pattern?
- (k) If you would like to see some random distributions of volcanoes, write

```
» poly = ugand_bd(2:end,:);  
» pt_in_poly_plot(120,poly);
```

The first line creates a pure polygon by stripping the information row from the boundary file. The small program in the second line then generates 120 random points in **poly**. (Be sure the Figure window is closed so that you don't superimpose this plot with the envelope above.) If you want to print the results, export the figure to WORD, and use the "finishing" procedures discussed in (4) below.

- 2. Next construct a more detailed **P-value plot** for a new set of **99** random patterns. Here you will use the same random seed, **23456**, now specified in the information structure, **info**. This can be done with the pair of commands:

```
» info.seed = 23456;  
» k_count_plot(volcano_pts,99,40,1,ugand_bd,info);
```

The use of the same random seed will ensure that the *same* data is generated to allow a better comparison of these two tests.

- (i) Save this figure to your home directory as **k_count.fig**.
 - (j) Now (with **k_count.fig** still open) open **k_function.fig** for purposes of comparison, by clicking: **File** → **Open** (in the main MATLAB menu) and navigating to your home directory.
- 3. How do these results compare with those above? Is there any new information here?
 - 4. For discussion purposes, you should include a print of these figures in your assignment report. A simple procedure is the following (using **k_count.fig** as an example):

- (i) First click: **File** → **Export** (in the Figure window) and export the file to your home directory in **Enhanced Metafile (.emf)** format.
- (j) Open WORD. [The following instructions are for **WORD 2003**]. click: **Insert** → **Picture** → **From File** , and insert the file **k_count.emf** .

- (k) Click on the image and resize it by dragging the corner with the mouse.
- (l) To position the image, click: **Format** → **Picture** (or **Object**) and in the **Layout** folder set **Wrapping Style** = “Behind Text”. You can then drag the image where you like.
- (m) Finally, you can (if you like) add additional labels by using **text boxes**:
 - a. Be sure the Drawing toolbar in WORD is open (**View** → **Toolbars**), and click on the **Textbox** tool. Then click (approximately) where you would like it to appear.
 - b. You can edit the Textbox by right clicking on the boundary of the box and selecting **Format Text Box**.
 - c. You can also use the mouse to drag the box and resize it (by again right click on the boundary).

Conclusions: In your concluding discussion of this study, consider how these results relate to the work of Tinkler. Are they consistent with his findings?

(3) In this study (based on Exercises 4.1 and 4.5 in B&G) you will analyze the “Myrtle Disease” data in the ARCMAP file: **T:\ese502\arcview\projects\Myrtles\Myrtles.mxd**. This data, which can also be found in the MATLAB workspace **myrtles.mat**, shows part of a grove of Myrtle trees in southern Tasmania that has been inflicted with a disease known as “Myrtle Wilt”.

- The data set is taken from Jillian Packham’s Ph.D. Dissertation (listed as *Myrtle Wilt* in the course Reference Materials). I have scanned portions of this work that will provide you with useful background material. [Pay particular attention to the highlighted portion of the discussion in *Section 2.5*.]
- Two shorter documents, *myrtle_wilt_dissertation.pdf* and *myrtle_wilt_explanation.pdf* are included in **T:\ese502\arcview\projects\Myrtles**. Additional material can be found by Googling “Myrtle Wilt”, etc.

The object of this study is to analyze the spatial relations between diseased and healthy myrtle trees (diseased trees are shown in red). In particular, the task is to determine whether these spatial relations are consistent with contagion of Myrtle disease, or whether they are more consistent with random infections. A *contagion hypothesis* might be that the diseased trees appear in “clumps” or “patches”. With respect to spatial dependencies *between* populations, this suggests that at scales where clumpiness is observed, healthy and diseased trees should show some degree of “repulsion”. Similarly, this hypothesis also suggests that *within* the diseased-tree population there should be some degree of “clustering” at these scales.

(a) To address these issues, you will first use both the “random shift” and “random permutation (relabeling)” tests to examine possible spatial dependencies between diseased and healthy trees. The locations of diseased and healthy trees are given in the matrices **L1** and **L2**, respectively, and the dimensions of the box shown in **Myrtles.mxd** are given as (**xmin,xmax,ymin,ymax**). In the following steps, be sure to save and print out all plots from MATLAB as in Problem 2 above.

1. Remember that (to avoid neighborhoods overlapping themselves on the torus) the maximum distance for shifts cannot exceed half the minimum of the width and height of the box. In this case,

$$\min\{(\mathbf{xmax} - \mathbf{xmin}),(\mathbf{ymax} - \mathbf{ymin})\}/2 = 46$$

So vector **D** contains a set of 11 representative distances between 1 meter and 45 meters to be used for the analysis. Here, you will use the **seed value**, 30456, to initialize the random number generator in the program. Hence, using the program **k12_shift_plot.m** for the *random shift test*, the command:

```
» k12_shift_plot(L1,L2,xmin,xmax,ymin,ymax,99,D, 30456);
```

generates a plot of the **P-values** for the distances in **D** based on 99 random shifts of the healthy trees (**L2**) relative to the diseased trees (**L1**).

2. What do these results tell you about the spatial relations between healthy and diseased trees? Are they consistent with the *contagion hypothesis* above?
3. Now try the *random permutation test* procedure, using the program **k12_perm_plot.m**. Here the relevant command is:

```
» k12_perm_plot(L1,L2,99,D, 30456);
```

4. How do these new results differ from those in (2) above? Why might this be so? Do they add support to the *contagion hypothesis*?
5. If you want to see some randomly relabeled patterns in MATLAB, you can do so as follows:

- a) First open **k12_perm_plot** and observe that the last two items of the output structure, **OUT**, contain the coordinates of a randomly relabeled pattern. To obtain this output, rerun **k12_perm_plot** as follows:

```
» OUT = k12_perm_plot(L1,L2,1,D,4345);
```

Here we only need *one* simulation to get the desired output. Notice also that I have chosen a different random seed.

b) Now save the new pattern by writing

```
» S1 = OUT.S1; S2 = OUT.S2;
```

c) To plot the results, write:

```
» plot(S1(:,1),S1(:,2),'r.', 'MarkerSize',20); hold on;
» plot(S2(:,1),S2(:,2),'k.','MarkerSize',15);
```

Here 'r.' means “plot as red points” ('k' denotes black). Also 'MarkerSize',20 means “make points of size 20”.

d) Finally, if you want to compare the test results for this random case, (S1,S2) with those of (L1,L2) above, use the command:

```
» k12_perm_plot(S1,S2,99,D,30456);
```

This will test the new pattern using the same seed as in (3) above.

6. Notice that in both the plots for **k12_shift_plot** and **k12_perm_plot** there appears to be some significant *attraction* between these populations at very small distances [$D(2) = 3$]. To check whether this makes sense, try the following commands in MATLAB (descriptions on the right):

» L = [L1;L2];	[stacks points to form single vector]
» n = size(L,1);	[counts number of rows in L]
» DIST = distance_mat(L);	[computes all pairwise distances for L]
» DD = DIST + 10*eye(n);	[replaces diagonal zeros with 10 (> 3)]
» D0 = DD(1:49,:);	[restricts DD to the L1 rows]
» D3 = (D0 <= 3);	[binary form of D0 with “1” for distances <= 3]
» rows = find(sum(D3'));	[finds rows with at least one distance <= 3]
» D3 = D3(rows,:);	[restricts rows of D3 to rows]
» v1 = sum(D3');	[counts all neighbors with distance <= 3]
» D3 = D3(:,50:end);	[restricts D3 to the L2 columns]
» v2 = sum(D3');	[counts only L2 neighbors with distance <= 3]
» MAT = [rows',v1',v2']	[collects all results in a matrix]

The last command (with *no semicolon*) then displays a matrix, **MAT**, in which the first column lists the row numbers of those *diseased* trees (**L1** trees) having at least one neighbor within a distance of 3 units. For each of these trees, the second column then shows exactly how many such neighbors there are. Finally, the last column shows how many of these are *healthy* trees (**L2** trees). If you have written these commands correctly, the first row of **MAT** should be [8 2 2], indicating that the diseased tree in row 8 of **L** has two neighbors within three units, and that both of these neighbors are healthy trees. [You can find

each of these trees in ARCMAP by noting that FID numbers (which start at zero) are always one less than the corresponding row numbers in MATLAB. So the tree in row 8 has FID number 7 in the Attribute Table for Diseased Trees.]

Based on the results in **MAT**, comment on whether the above significance result appears to be consistent with this information. (*Hint*: If the distribution were random, then (as a crude estimate) one would expect about 58.5% [= $100\{69/(49+69)\}$] of these neighbors to be healthy trees. Compare this with the observed fraction).

(b) Finally, you will consider possible spatial relations *within* the diseased population (**L1**) relative to the healthy population (**L2**), again by means of “random permutations”:

1. Using the program, **k2_diff_plot.m**, with the same data and a new **random seed**, 65846, the relevant command is now:

» **k2_diff_plot(loc,n1,99,D,65846);**

2. Do these results lend support to the *contagion hypothesis*? Be specific in your answer.

Conclusion: In your concluding discussion of this study, consider how your results relate to the findings of Packman (and the earlier work of Elliot that she discusses).