

Geographical Ecology

Student Name; Z620: Quantitative Biodiversity, Indiana University

13 February, 2015

OVERVIEW

In this exercise, we will add geographical context to alpha (α) and beta (β) diversity. We will introduce Geographical Information Systems (GIS) to map and spatially examine environmental and biodiversity data. This will allow us to explore core concepts like spatial autocorrelation, aggregation, and scale dependence.

After completing this exercise you will be able to:

1. Identify primary concepts and patterns of geographical ecology
2. Examine effects of geographic distance on environmental and ecological similarity
3. Characterize aggregation of abundance across space
4. Examine the extent to which patterns of diversity depend on spatial scale
5. Use geospatial data and packages to conduct GIS operations in R
6. Use control structures such as `loops` to control how R operates on variables

Directions:

1. Change “Student Name” on line 3 (above) with your name.
2. Complete as much of the exercise as possible during class; what you do not complete in class will need to be done on your own outside of class.
3. Use the handout as a guide; it contains a more complete description of data sets along with the proper scripting needed to carry out the exercise.
4. Be sure to **answer the questions** in this exercise document; they also correspond to the handout. Space for your answer is provided in this document and indicated by the “>” character. If you need a second paragraph be sure to start the first line with “>”.
5. Before you leave the classroom, **push** this file to your GitHub repo.
6. For homework, follow the directions at the bottom of this file.
7. When you are done, **Knit** the text and code into a PDF file.
8. After Knitting, please submit the completed exercise by noon (12:00 PM) on Wednesday, February 18th by creating a **pull request** via GitHub. Your pull request should include this file (*GeographicalEcology_Exercise.Rmd*) and the PDF output of Knitr (*GeographicalEcology_Exercise.pdf*).

1.) SETUP

A. Retrieve and Set Your Working Directory

In the R code chunk below, provide the code to:

- 1) clear your R environment,
- 2) print your current working directory,
- 3) set your working directory to your “/GeographicalEcology” folder, and
- 4) load the **vegan** R package (be sure to install if needed).

```
rm(list=ls())
getwd()
setwd("~/GitHub/QuantitativeBiodiversity/Assignments/GeographicalEcology")
```

B. Load Packages

We will use the **vegan** package for biodiversity estimators and related functions.

We will also use a suite of six packages developed in R for geographical information systems (GIS). Be sure to run **install.packages("PackageName")**, if not previously installed, where **PackageName** is the name of the package you want installed. Or, run **install.packages(PackageName, type="source", dependencies=TRUE)**, if the previous command doesn't work. A pop-up may appear, asking you to install certain R packages; choose yes.

```
# install or require these packages:
# 1. sp
# 2. geoR
# 3. rgdal
# 4. raster
# 5. RgoogleMaps
# 6. maptools
install.packages("rgdal", repos="http://www.stats.ox.ac.uk/pub/RWin")
```

C. Load and Compile a Large Dataset

We will analyze environmental and bacterial community data from a survey of shallow ponds found east of Bloomington, IN. These ponds are scattered throughout Brown County State Park, Yellowood State Forest, and Hoosier National Forest. In the 1940s, Maury Reeves of the Indiana Department of Natural Resources began constructing refuge ponds for wildlife. In the summer of 2013, we visited approximately 50 of these ponds and recorded their geographic locations using a GPS unit; 'GPS' is the acronym for Global Positioning System. We sampled aspects of water chemistry, physical properties, and bacterial community composition. Let's load the environmental and site-by-species data for the refuge ponds.

Take a look at the **Environment** tab in the upper right console of RStudio. You should see there are 16,384 operational taxonomic units (OTUs) distributed across 51 sites, for which, we have 19 environmental and geographic variables recorded. These variables include elevation (m), geographical coordinates (lat-long), temperature (C), diameter (m), depth (m), redox potential (ORP), specific conductivity or SpC ($\mu\text{S}/\text{cm}$), dissolved oxygen (mg/L), total dissolved solids (g/L), salinity (p.s.u. \sim ppt), color - measured at absorbance = 660; an estimate of carbon in the water sample, chlorophyll *a* ($\mu\text{g}/\text{L}$), dissolved organic carbon (mg/L), dissolved organic nitrogen (mg/L), and total phosphorus ($\mu\text{g}/\text{L}$).

In addition to this wealth of environmental data, let's add four diversity-related columns of data to **Ponds** data set. These will provide basic diversity-related variables to explore with respect to geography and environmental conditions. There will be a column for richness (S), total abundance (N), Shannon's Diversity (H), and Simpson's evenness (De).

Now we have a large compilation of geographical, environmental, and biodiversity data. Let's do some Geographical Ecology!

2.) MAP SAMPLES AND DATA

Let's visualize the spatial distribution of our samples with a basic map in RStudio. Let's generate a map of the refuge ponds using the **GetMap** function in the package **RgoogleMaps**. This map will be centered on Brown County, Indiana (39.1 latitude, -86.3 longitude).

This map displays a lot of useful information that we otherwise, would not have been aware of. For example, all points are on State or National Forest land. Likewise, the sample ponds appear to be aggregated in four or five small groups and distributed across a topographically complex area.

Despite being a fast way to contextualize our sample ponds within the broader landscape, the Google map misses a lot of information that would otherwise help us to understand the environmental and geographical factors that may coincide with our observations on diversity. Likewise, because the Google map is only an image, it doesn't contain any extractable environmental or geographic data.

For spatially explicit data on environmental and geographic features, i.e. geospatial data, we can turn to one of the many freely accessible online GIS databases and warehouses. Here, we will use the high quality geospatial data on Indiana water bodies and percent landcover. We obtained these data 'layers' from the **IndianaMap** geographical layer gallery: <http://maps.indiana.edu/layerGallery.html>.

3. PRIMARY CONCEPTS AND PATTERNS

Having imported our community and environmental data from the refuge ponds, as well as having obtained a wealth of geospatial data from online sources, we are now ready to explore primary concepts and patterns of geographical ecology.

A. Spatial Autocorrelation

When examining spatial data, it is important to check for autocorrelation not just among variables but across distance. Here, we reveal a way of detecting autocorrelation with respect to scale, that is, by using a **variogram**. Variograms are frequently used in spatial analyses and reveal the degree of spatial autocorrelation in sample data and how the autocorrelation (measured as the **semivariance**) changes over distance.

The semivariance is a measure of the dispersion of all observations that fall below the mean. In fact, if you were to find the differences between all possible points spaced a constant distance apart and then find the variance among the differences, and then divide that variance in half, you would have the semivariance. While the semivariance is similar to the variance, it only considers observations below the mean. In this case, higher semivariance means implies greater dispersion (lower autocorrelation) of values at a given spatial scale.

Let's plot the variogram for one of our environmental variables.

Question 1: How does the **semivariance** change with distance? Is there a distance (i.e. "spatial lag") where semivariance is very high or low?

Answer 1:

For a more visually informative picture, we can visualize autocorrelation across the landscape, by calculating **Moran's I**, a correlational statistic that measures autocorrelation based on feature locations and feature values. Moran's I evaluates whether the pattern expressed is clustered, dispersed, or random. The function we will use also assigns a global Moran index value, a z-score and p-value.

Using R's **raster** package, we can calculate **global** (across the landscape) and **local** (in comparison to neighbors) measures of **Moran's I**.

Question 2: Moran's global I can range from 0 to 1. What was Moran's global I? Looking at the 'landscape' of Moran's *local* I for percent landcover, make some observation about spatial autocorrelation. For example, where is it high or low?

Answer 2:

Pattern 1: Distance-decay relationship

The distance-decay relationship is the primary pattern of spatial autocorrelation, and captures the rate of decreasing similarity with increasing distance. This pattern addresses whether communities close to one

another are more similar than communities that are farther away. The distance-decay pattern can also be used to address whether near environments have greater similarity than far ones.

Let's load the **simba** package and generate distance decay relationships for bacterial communities of our refuge ponds and for some of the environmental variables we measured. Note, this analysis will span the next few code blocks.

```
# install or require the "simba" package

# transform environmental data to numeric types

# calculate the distance (Euclidean) between the plots regarding environmental variables

# transform all distance matrices into list format:
```

Now, create a data frame containing similarity of the environment and similarity of community.

Finally, let's plot the distance-decay relationships, with regression lines in red.

Let's, examine the slope of the regression lines, asking whether they are significantly different from one another.

Question 3: Are microbial communities that are closer in geographic distance also closer in compositional similarity? How about for environmental condition?

Answer 3:

Concept 2: Spatial Aggregation

Natural phenomena are generally clustered, i.e., spatially aggregated. Individuals, conditions, and events often occur in patches, clusters, pulses, etc. Take for example, the ponds in our sample area. A high level of aggregation would suggest that if we encounter one individual, then we are likely to encounter others of the same species nearby.

Pattern 2: Spatial abundance distribution

One of the primary patterns of spatial aggregation in ecology is the distribution of a species's abundance within a landscape, also referred to as the **species spatial abundance distribution (SSAD)**. The SSAD reveals the frequency at which we find a species at a particular abundance. In this way, the SSAD is similar to a histogram.

Here, we will examine SSADs for OTU's in the refuge pond dataset by constructing **kernel-density curves**. Kernel density curves are analogous to histograms, but avoid the arbitrary creation of bins or discrete classes.

For example, suppose we were interested in how the location of individuals varied across sites or samples. Let's simulate this by drawing values from a normal distribution, at random.

Below, we will examine the SSADs of OTUs that are randomly drawn from the refuge ponds dataset. But first, let's begin by defining a function that will generate the SSAD for a randomly drawn OTU.

Next, we will draw 4 OTUs at random and plot their SSADs. But first, we will need to introduce **while loops**.

While loops

If you have ever heard anything to the effect of, “While you’re at it, do this...”, then you are familiar with the concept of a while loop. The while loop is a type of **control flow** structure that allows us to tell a program to perform an operation *while* some condition is satisfied.

For example, we might want R to draw numbers at random until 4 numbers less than 50 have been drawn.

```
numbers = c()
while (length(numbers) < 4){ # while the counter is less than 4
  x <- runif(1, 1, 100) # draw a number at random from 1 to 100
  if (x < 50){ # if the number is less than 50...
    numbers <- c(numbers, x)
  }
}
numbers # check our numbers, each should be less than 50
```

```
## [1] 8.963865 32.382953 9.048922 13.168790
```

Having very briefly introduced while loops, and inadvertently, an `if` statement, let's write a chunk of code that will draw OTUs at random from our refuge ponds dataset, and then generate their spatial abundance distributions (i.e. SSADs).

Question 4: Is the sampled abundance for a given OTU aggregated? If so, how do you know? That is, how do you interpret the pattern in the kernel density curve? Are there many sites with low abundance and few sites with high abundance?

Answer 4:

Question 5: Each row in the site-by-species matrix represents a site. Each column represents an OTU. If the SSAD is generated by considering all rows for a single column (i.e. OTU), then what do we obtain when we consider all columns for a given row (i.e. site)? Have we examined this sort of data structure before?

Answer 5:

Concept 3: Scale-Dependence

Our idea of whether variables are spatially autocorrelated and whether the abundances of OTUs are spatially aggregated can change with aspects of spatial scale, i.e. extent and grain. **Extent** is the greatest distance considered in an observation or study. **Grain** is the smallest or primary unit by which the extent is measured.

Let's generate two random samples from a normal distribution; one sample for x-coordinates and one for y-coordinates. We'll let each x-y pair represent the location of a single individual, where all individuals belong to the same species. Then, we'll plot the spatial distribution of our randomly distributed individuals at different extents.

Question 6: What effect does changing the extent have on aggregation? Do you find this important or interesting given that a) all points were drawn from the same distribution and b) each plot contains the same points as all other plots with smaller extent?

Answer 6:

It should be clear from above, that ‘random’ does not mean absent of aggregation. In fact, most statistical distributions from which random samples can be drawn are very aggregated. That is, they have obvious modes around which most values occur.

Moving on, let’s explore the effect of changing spatial **grain**, from a fine grain to a coarse grain. We will do this while holding extent constant. We will then plot heat maps (i.e. 2D histogram) revealing the density of individuals in the landscape. Last, we will plot kernel density curves to reveal the probability that an individual chosen at random from the landscape will have come from a site with a particular abundance.

```
# install or require "gplots"
par(mfrow=c(2, 2))
```

Question 7: Beyond changing the pixilated appearance of the plots, what does changing the spatial grain mean for interpreting aggregation? Use kernel density plots to help answer this question.

Answer 7:

Question 8: How are the kernel density curves we just generated for our randomly drawn points related to the species spatial abundance distributions (SSAD) that we generated for OTUs in our refuge plots?

Answer 8:

Primary Concept 3: Spatial Accumulation

So far, we have discussed spatial autocorrelation and aggregation as core concepts of geographical ecology. Likewise, we have introduced and examined primary patterns for both of those concepts. Here, we introduce another core concept, accumulation across space. It may seem self-evident that, if starting from the scale of a single individual and increasing our sample area, that we will inevitably encounter more species, OTUs, or other taxa.

For example, suppose we replicate our above random sampling strategy of drawing x-y coordinates from a normal distribution. But, instead of drawing just one sample representing one species, we will draw 50 samples, with each representing a species with 1000 individuals.

```
# initiate the plot

# while our community has less than 100 species

# choosing the mean, standard deviation, and colors at random
```

Having generated a simulated landscape occupied by 50 species having 1000 individuals apiece, we can examine how richness can accumulate with area. Let’s begin by picking a corner at random and then accumulating area.

```
# while the spatial extent on the x and y is less than or equal to 100
# Set richness to be zero
# for each species in the community
# assign the x coordinates
# assign the y coordinates
# assign the species name
```

```
# for each pair of xy coordinates in xy.coords
# if the individual is within our current spatial extent...
# then the species occurs there
# break out of the last for loop because we're only considering
# incidence, and not abundance.
# In other words, if the species occurs once, that's good enough

# increase the extent multiplicately, but slowly
```

Having generated our primary vectors, **S.list** and **A.list**, we can analyze how species richness scales with area. In short, we can analyze one of ecology's oldest and most intensively studied patterns, the **Species-Area Relationship**.

Pattern 3: Species-area relationship (SAR)

The fact that we accumulate species, and likewise increase richness, with increasing area is far from interesting. In fact, we just showed that we can expect this as a result of random sampling. What is interesting is the rate at which the accumulation of taxa occurs. Arrhenius (1921) first described the general form of the *species-area relationship (SAR)* as a power-law: $S = cA^z$ where S is species richness and A is area.

Power-laws reveal how one quantity scales with another, most often across orders of magnitude. Arrhenius's formula predicts a rate of increase in richness that is approximately linear in log-log space. That is, $\log(S) = c + z\log(A)$, where z is the scaling exponent.

Question 9: The authors of your assigned reading revealed that the exponent of the SAR may be influenced by geographic, ecological, and evolutionary factors. But, what in general, is the value of z?

Answer 9:

Question 10: What was the slope of the species-area relationship for our randomly assembled community? Is this similar to the slopes you encountered in the reading?

Answer 10:

Question 11: We could use this 'random placement' approach to model how many ecological phenomena might occur via random sampling. What other spatial aspects of alpha and beta diversity could we address? Suggest at least 3.

Answer 11:

7) HOMEWORK

1.) Complete the in-class exercise and the homework, Knit to a pdf, and submit a pull request before noon on Wednesday, February 18th, 2015.

2.) Each refuge pond has an associated diameter. Build the species-area relationship for the refuge pond dataset using the following recipe:

1. Using the formula for the area of a circle, calculate the area for each pond.
2. Randomly choose one pond and obtain its area and richness.

3. Choose two ponds at random and obtain their combined richness and their summed area. Do not simply sum the richnesses of the two sites, as this will result in double-counting species.
 4. Choose three, four, five, etc. ponds at random, repeating the above steps each time. You will eventually work your way up to 51 ponds. **You will need to use loops for this.**
 5. At this point you should have two vectors, one for richness (S) and one for area (A).
 6. Plot the SAR and estimate its slope.
 7. REPEAT steps 2 through 6 one thousand times, adding each new SAR to the same plot. Once again, you will need to use loops, as above.
 8. In a second plot, generate a kernel density curve for the slopes of the SARs.
-
- 3.) Draw several general conclusions from your analyses in question #2.
 - 4.) Which environmental and diversity variables reveal positive spatial autocorrelation?
 - 5.) A. How many OTUs are present at more than 10 sites? How many OTUs only occur at one site?
 - 6.) In considering total abundances (N) among the refuge ponds, we are really only considering the number of detected 16S rRNA reads for any given OTU. Find the mode of the SSAD for each OTU that is present at 10 or more sites; it's difficult to generate an informative histogram for less than 10 sites. Then, generate a kernel density curve for these modes, revealing the pattern of modal abundance across these more common OTUs. Draw general conclusions about trends across the refuge pond OTUs from your results.