

# Scripts, Importing Data, Simple Linear Regression, & Ripley's K in R

*Michael L. Treglia*

*Material for Lab 3 of Landscape Analysis and Modeling, Spring 2016*

*This document, with active hyperlinks, is available online at: [http://mltconsecol.github.io/TU\\_LandscapeAnalysis\\_Documents/Assignments\\_web/Assignment03\\_DataImport\\_Regress\\_RipleysK.html](http://mltconsecol.github.io/TU_LandscapeAnalysis_Documents/Assignments_web/Assignment03_DataImport_Regress_RipleysK.html)*

**Due Date: Tuesday, 9 February 2016**

**This assignment is worth 10 Points, with 1 points per question.**

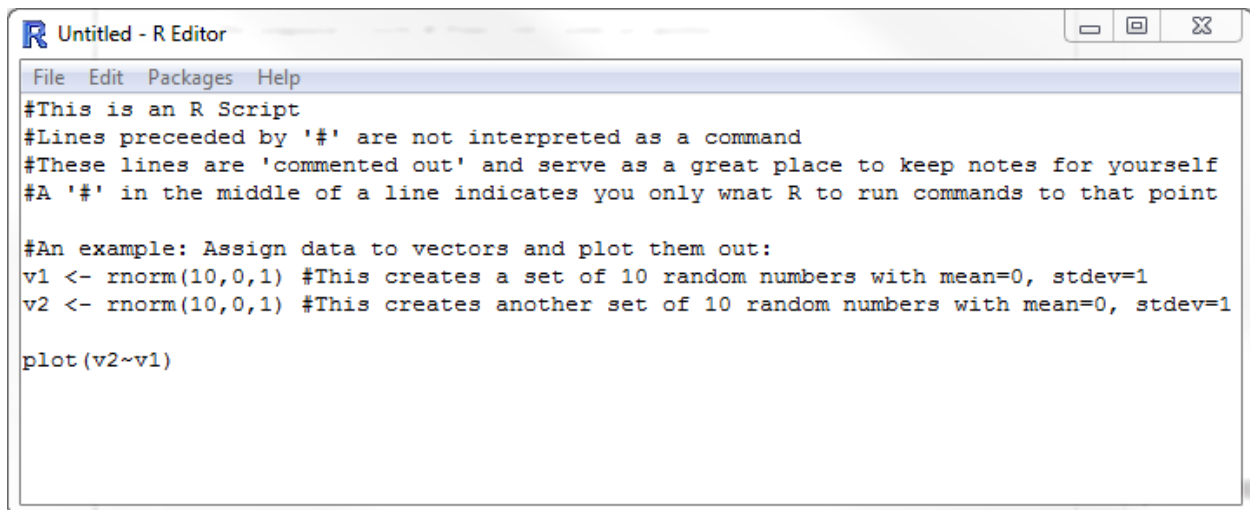
This lab is designed to further familiarize students with R for statistical computing, introducing them to point pattern analysis with Ripley's K, discussed in lecture. Simple linear regression is also highlighted in this lab, to introduce formulation of statistical tests in R.

## Scripts in R

In the first lab involving R, we simply used the Console interface to run commands, where you type things, hit enter (or return) and see an immediate result. As you observed, you can carry out the basic functions quickly and easily, though saving your commands and making your work repeatable would be difficult if restricted to that interface. Most simply, it seems you could copy and paste all of your commands into plain text editor (e.g., Notepad and Text Edit are built-in editors on Windows and Macs, respectively). But isn't there a better way?!?

**YES!!!** R has a built-in tool for this - *scripts*. Scripts are sets of code that you can re-run as you desire from a static file, either line by line, or all-at-once.

To create a new script in R, use the File menu at the top of your R window, and select 'New Script'. A new window will appear within R; you can type your commands here, and use [CTRL]+R on Windows or [Command]+[Enter] on Macs to run an individual line where your cursor is, or all lines that you have selected with your mouse. You can save this file as a '.R' file, which you can open up from within R, or work with it using a text editor. Below is an example Script file, opened in R, with some sample code.



```
Untitled - R Editor
File Edit Packages Help
#This is an R Script
#Lines preceded by '#' are not interpreted as a command
#These lines are 'commented out' and serve as a great place to keep notes for yourself
#A '#' in the middle of a line indicates you only want R to run commands to that point

#An example: Assign data to vectors and plot them out:
v1 <- rnorm(10,0,1) #This creates a set of 10 random numbers with mean=0, stdev=1
v2 <- rnorm(10,0,1) #This creates another set of 10 random numbers with mean=0, stdev=1
plot(v2~v1)
```

You can include “comments” in script files, which R will not run as code. To designate a comment, simply use a ‘#’ before the text, as shown above. You can do this for an entire line or for part of a line, following actual R code, as seen above. Comments are useful ways to note how a set of code works, what it does, etc. Everything that you run from the Script will be displayed in the interactive console window as it goes; comments will be displayed in the console, but again, are not interpreted as commands. You can try copying the code displayed in the preceeding image, and run it from a new script window on your own computer.

## Importing Data

Now that you have this great tool (scripts) at your fingertips to make your work repeatable and reproducible, lets try to do some more complex steps, which are normal parts of working with data in R. We’ll work with data from tree census at Harvard Forest (<http://harvardforest.fas.harvard.edu/>). We will work with a ‘.csv’ file, which stands for ‘Comma Separated Values’ - this is among the easiest of formats to import into R. If you have data you wish to work with in R but it is currently stored in a Microsoft Excel format, you can easily export it to a .csv file from Excel using the “Save As” option, and choosing ‘.csv’ as the file type.

I generally suggest manipulating your data in R, as you will do in exercise. If you save your work in a script, from data import through production of final results, you can effectively have a record of everything you did, and identify any mistakes you may have made. You cannot easily do this in spreadsheet software, where it might be easy to accidentally alter your data (e.g., sorting only one column of data instead of an entire dataset is an easy mistake to make in a spreadsheet).

For this exercise, download the Lyford Mapped Tree Plot Data from Harvard Forest, available here: <http://harvardforest.fas.harvard.edu:8080/exist/xquery/data.xq?id=hf032>. Under the “Data” section of the website, towards the top, there are two links: the first, (hf032-02) is the data - simply click to initiate the download; the second (hf032-01) is a map of plots in which data were collected. Save the dataset to a logical place that you can navigate to on your computer. If the data automatically open in Microsoft Excel, simply save it in the same format it came as (.csv). (It can be useful to browse the data in Excel, but we’ll do all data manipulation and analysis in R.)

# Harvard Forest HF032

## Lyford Mapped Tree Plot at Harvard Forest since 1969

### Related Publications

### Data

- **hf032-02**: tree data
- **hf032-01**: Lyford plot map

The citation for this dataset is: *Foster D, Barker Plotkin A, Lyford W. 1999. Lyford Mapped Tree Plot at Harvard Forest since 1969. Harvard Forest Data Archive: HF032.*

You can browse the website with the data for information about how the data were collected. In particular, look at the Detailed Metadata section. This will explain what each column in the dataset refers to, and provides a key to the abbreviations (notably different tree species) in the dataset. *This will be handy later in the assignment*

When you are working in R, it is good practice to specify the folder location where you want R to be reading data from/writing data to. If you don't do this, you may have to specify the full path for reading and writing files, rather than just the file names. The active folder that R is looking at is called the "working directory". R has a default working directory, but that is not usually where you keep your data stored. To find out the current working directory, type:

```
getwd()
```

To change it to the folder where you have your data stored you will use the 'setwd' function, and the only argument that you need to include is the desired file path in quotation marks (either double or single quotes, as long as you are consistent, will both work. For example, I want to set my working directory to 'Documents' folder on my the 'D' drive of my computer.

```
setwd("D:/Documents")
```

To get the folder path, you can simply browse the folder manager on your computer (e.g., My Computer/Windows Explorer on Windows), find where you have the data, copy the path at the top of the window, and paste it into your script window in R.

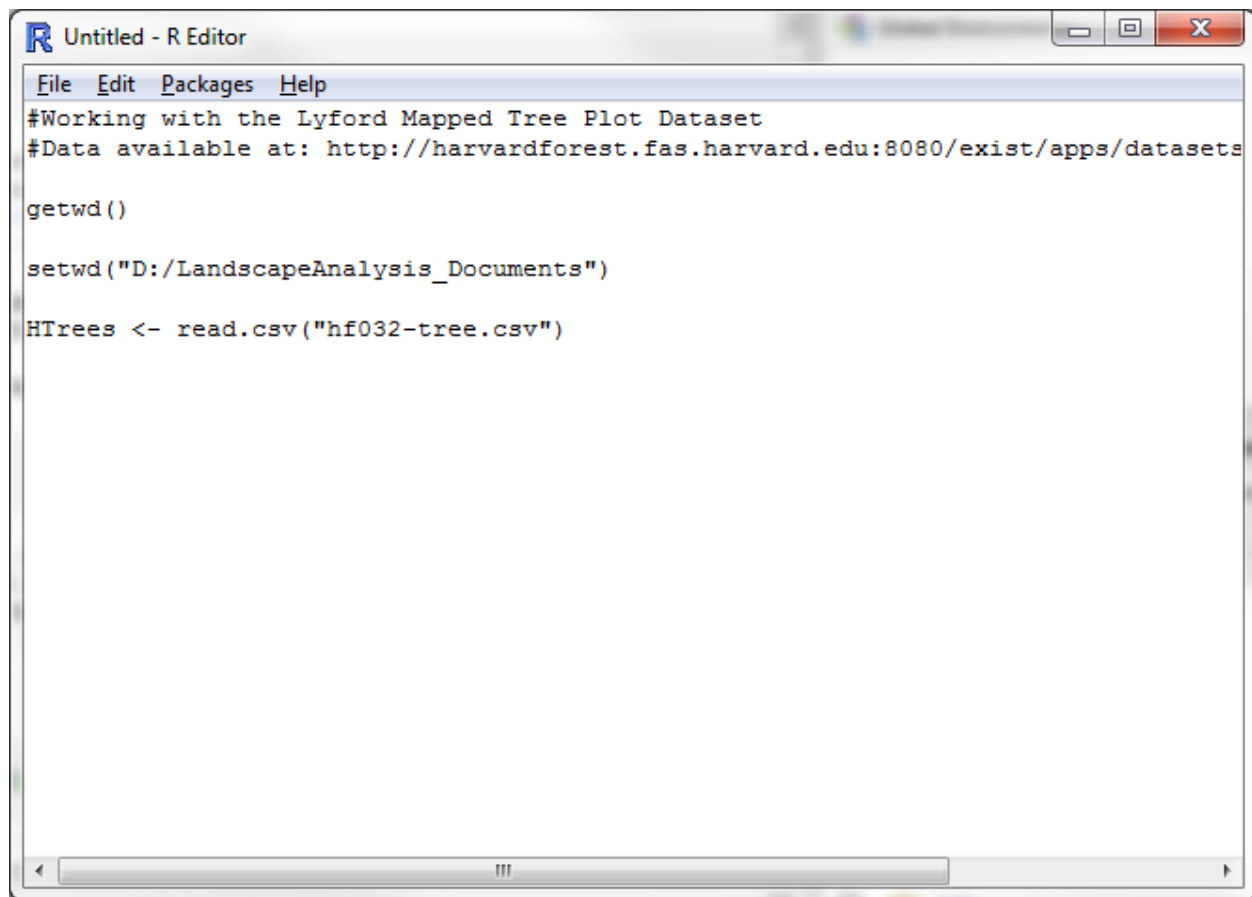
*\*If you are a Windows user, you will have to make a small edit: you need to change the backslashes (') in the file path to forward slashes (/) or double backslashes (\\) for your path to be read correctly.*

The conventions for navigating through the file system in this context is similar to on a unix/unix-like system. To direct the computer to the directory one level up, use '..', and to get to the home directory, use '~'.

Once you are within your desired working directory, you can view all of the files and directories contained in that directory:

```
dir()
```

To actually read the focal dataset into R, we will use the 'read.csv' function, with the only argument as the file name for the dataset you've saved. This is a specific form of a function 'read.table', which can import various types of text files, but 'read.csv' is set to use .csv files by default. For different datasets you may need to use different arguments, explained in the Help for 'read.csv'. We will assign the imported data to a sensible name - 'HTrees'. If doing this in a script window, you should have something similar to the following image.



At this point, you can inspect the dataset structure, and quickly look at the first and last few lines to get a feel for what the data look like.

**Question 1:** *How many different types of tree species are listed in this dataset? (I want you to figure this out in R, and for full credit you must indicate what function you used).*

## Subsetting Data

The following website may aid in the next few steps: <http://www.ats.ucla.edu/stat/r/modules/subsetting.htm>

This dataset is large, with many different species in a variety of plots from many years. We will only focus on the plot "SE11", and the species *Acer rubrum* and *Quercus alba*, from the year 1991.

**Question 2:** *take a subset of the data so you have a dataset with the properties just described. To show that you've done, simply provide the commands you used. You can do this in one or multiple steps. For consistency with the rest of this lab, assign the final subsetting dataset the name 'HTreesSub'.*

If you look at the structure of this dataset (HTreesSub), you will notice the number of levels listed for factors is the same as for the original dataset, but if you look at the actual data, it should only have *Acer rubrum* and *Quercus alba*. You can clean this up a bit using the 'droplevels' function:

```
HTreesSub <- droplevels(HTreesSub)
```

Given that our focal year is 1991, and a lot of the columns in the dataset have information about other years, we can discard those columns. There are a few ways to do this, but here we'll use a "column index" method. In R, every row and column has an numerical index (row i, column j), starting at 0. We can specifically call on these in functions and we refer to indices in R using 'datasetName[rowNumber,columnNumber]'. For example, if you want to see what is in the 2nd row and 3rd column of your subsetted dataset, you can use

```
HTreesSub[2,3]
```

To quickly identify which column corresponds to which index number, you can use the function 'names'

```
names(HTreesSub)
```

Then, use indices to define a dataset with only the columns you want. In this case We will the following variables:

- treeid
- xsite
- ysite
- species
- dbh91

**Question 3:** *Create another subset of the data with only the variables listed above, from HTreesSub. Name this new subset "HTreesSub91"; provide the code you used to do this.*

## Linear Regression

In class we discussed linear regression; lets give this a try in R. Perhaps we hypothesize that for some reason, there is a significant effect of x-position (on a cartesian plane) on the diameter at breast height of the trees. Perhaps we are aware of a historic gradient of farming in this area, which has a lasting effect on trees that would cause this. We can test for the presence of such a relationship using a linear regression. Note this is a very basic example with a relatively simple method, but should familiarize you with the way model equations work in R.

It can be helpful to plot the data in advance and visually inspect the data.

**Question 4:** *Plot the diameter of trees along the y-axis and the x-position of trees along the x-axis; paste the result into a document with your answers.*

The function to apply a simple linear regression is 'lm' (i.e.,linear model). In functions like this, we need to define an equation, just as we would if writing out a regression equation. The function, for a simple case like this, is structured as 'dependentVariable~independentVariable'. We can also specify our dataset in a separate argument of the function. Thus, in our case, we would use:

```
dbhXpos <- lm(dbh91~xsite, data=HTreesSub91)
```

The above command creates an object called 'dbhXpos' that contains all of the information directly calculated from the regression. Note that there were some missing values for dbh91 in this dataset - the default way of dealing with these observations is simply to ignore those rows of data.

To inspect everything stored in the regression object, simply look at the structure as you would any other object. It looks complex - you can take time and try to interpret it, but we'll use the 'summary' command to get some useful information.

```
summary(dbhXpos)
```

You can see the F-statistic, and the p-value at the bottom of this output. Additionally, there is an output of 'r-squared', which is the proportion of variance in the dataset explained by the regression (calculated as Regression Sum of Squares / Total Sum of Squares)

Using the plot function on the lm object (i.e., `plot(dbhXpos)`) will produce some diagnostic plots for inspecting the residuals, which can be useful in examining whether the data actually meet assumptions of the test.

If you want to view the regression line on the plotted data, plot the data as for Question 4. Then, immediately after, use (without the '#'):

```
abline(dbhXpos)
```

If you want to adjust the appearance of line, look into the arguments 'lwd' and 'lty'. Check the help/google around for more informatino about 'abline'.

## Ripley's K Analyses

### R packages

Though R has great functionality on it's own, a major strength of the software is the ability of others to write and contribute packages that add valuable functions. There are packages for a huge range of analyses, well beyond ecological applications. Some packages come with R, such as 'MASS' and 'lattice', whereas others have to be installed.

Upon starting R, only a limited set of pacakges are loaded by default - there are a few reasons for this, and here are a couple of them. First, multiple packages, developed independently, might have the same name for different functions, thus loading them at the same time can cause confusion (by default, the function in the most-recently loaded package takes priority, though a function for a specific loaded package can be called using `'packageName::function(arguments)'`). Another reason for only opening a limited number of packages upon starting R is the potential for lots of memory to be used; sometimes large sets of code underly R functions, which may occupy more temporary computer memory that you may need for working with large files.

Each package has a website hosted by CRAN (Comprehnsive R Archive Network), with at least basic documentation of commands. For example, the site for spatstat is <http://cran.r-project.org/web/packages/spatstat/index.html>. Some packages also have 'vignettes', which go into detail on accomplishing some particular tasks.

If you are trying to figure out what R packages might be useful for a particular field or realm of analysis, you may find it helpful to browse the 'Task Views' pages, available at: <https://cran.r-project.org/web/views/>. Clicking on a relevant link there will take you to a page that highlights some R packages and contextualizes how they may be used.

To see if a package is already installed, you can simply try to load it it using the 'library' function. For example,

```
library(sp)
```

If a package is not installed, R will indicate that in the console that it is not available. If the package is available and opens, R will typically indicate what version of the package it is, and sometimes include other details.

If a package you need is not installed, use the 'install.packages' function:

```
install.packages("spatstat")

#To install multiple packages at once:
install.packages(c("spatstat", "maptools"))
```

When using this command, R will ask you to specify what repository or mirror to use. These are essentially servers around the world, at various institutions, with all of the R code available for download. Typically it makes sense to pick one from nearby for the quickest download times - in Tulsa I usually chose a repository in Kansas or Texas as adjacent states.

Sometimes functions within a package call on functions from other packages - this is called dependency; when you install one package, R will typically install other packages necessary for the specified packages to work. Check out the help for 'install.packages' for more details.

If the 'spatstat' package is already installed, try 'maptools', 'raster', or 'vegan', which are some other useful packages.

In the next steps of this exercise, we will use the following packages:

- spatstat
- sp
- maptools

Make sure they are all installed and loaded.

## Converting the Data into Spatial Objects

Our focal dataset is currently stored as a 'dataframe', or a table with multiple types of data in different columns (some columns have numeric data, while some have factor data). Though we know from the metadata that the xsite and ysite variables have refer to coordinates within the sampling plots, R has no way of knowing this is the case - to R, they are simply additional variables. For some analyses this is okay, but for others, the data need to be defined as a spatial dataset. We can simply do this as follows for points with x/y coordinates:

```
coordinates(HTreesSub91) <- c('xsite', 'ysite')
```

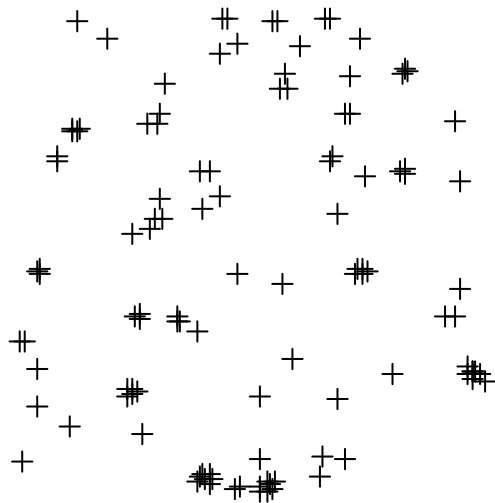
**Question 5:** *After the above step, what class of object is 'HTreesSub91'?*

We have not specified a coordinate reference system. However, there is none for these data - they are projected according to a simple cartesian coordinate system in the study plots. If the data were in a standard coordinate system and we wanted to include them in a map, this would be necessary and we will cover this later in the semester.

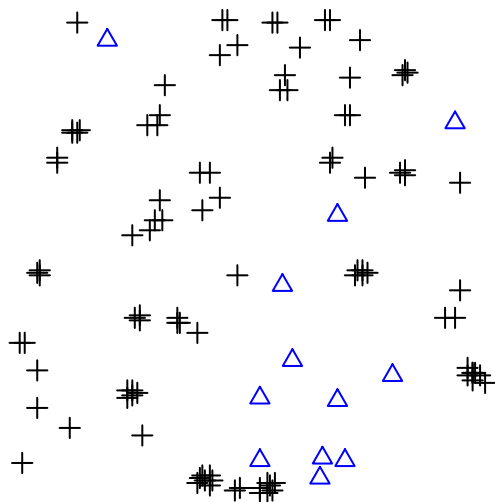
## Plotting Spatial Data

As with simple linear regressions, it is often informative to view our data as we analyze it. The 'plot' function is actually very broad, and its use varies with the type of object you are working with. You can view the location of all the trees in the focal plot in a simple graphic with:

```
plot(HTreesSub91)
```



**Question 6:** Building on what you've learned, using a simple plotting function and data subsets, plot *Acer rubrum* and *Quercus alba* on the same figure in different symbols. You may need to use two commands, and the arguments 'pch' and 'col'. Here's a hint: look into the function 'points'. You may need to browse resources given in the Lab 2 handout, *Getting Familiar with R*. Don't worry about making a legend, but you should be able to keep track of which points represent which tree species for yourself. The result should look something like what is depicted below; copy and paste your figure into your answer document and provide your code for the plot.



It is definitely possible to make nicer-looking maps in R, but for now, something simple as shown above is sufficient.

**Question 7:** Based on inspection of your previous figure, do you suspect that one species exhibits a significant spatial pattern at certain scales? If so, which species, what pattern (clustered or dispersed), at what scale (small/local or broad/landscape), and why?

## Ripley's K Analyses

Though plotting the arrangement of the two focal tree species might help us guess which trees are clumped or dispersed at different scales, we need to use statistical inference to draw strong conclusions. As we covered in lecture, Ripley's K analysis is perfect for this.



We will first adjust the dataset to only include the positional information (now stored as a separate entity within the object), and the tree species.

```
RipleyTrees <- HTreesSub91[,2]
```

Lets first look at each tree species separately to identify whether individual trees of a single species tend to be clustered or dispersed with respect to one another. We will create two subsets - one for each species.

```
acru <- subset(RipleyTrees, species=="ACRU") #Acer rubrum subset  
qual <- subset(RipleyTrees, species=="QUAL") #Quercus alba subset
```

Now we need to turn each of these into ‘point pattern’ objects. To learn more about this class of objects, you can use `help(ppp)`.

```
acru.ppp <- as(acru, "ppp")  
qual.ppp <- as(qual, "ppp")
```

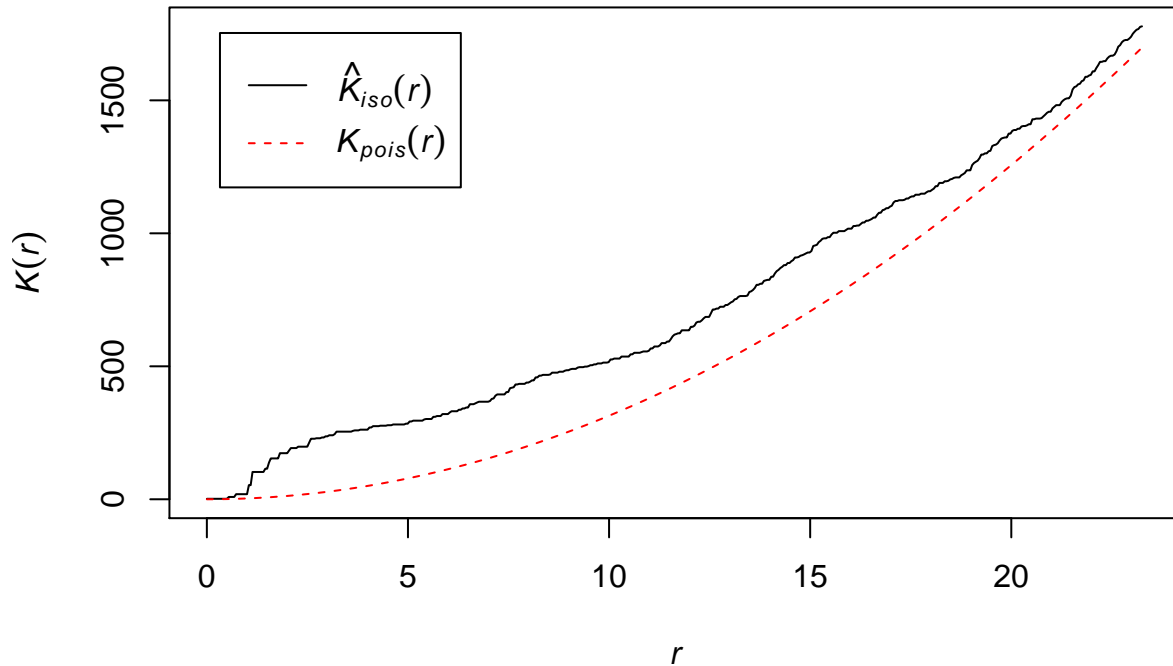
Ripley’s K is computed using the function ‘Kest’ in the spatstat package. We can simply run the function on each of the two tree species and plot the results within a single line of code. We want to use some sort of edge correction as we discussed in class - the help file for Kest has some explanation and citations for a number of different options. The ‘best’ option should generally be appropriate. In your interpretation, remember that the observed K function indicates clustering when it is above the theoretical curve, and segregation when it is below the theoretical curve. The theoretical curve may be labeled ‘Ktheo’ or ‘Kpois’ in the output.

```
plot(Kest(acru.ppp, correction="best"))  
plot(Kest(qual.ppp, correction="best"))
```

We can also achieve the same thing by assigning the ‘Kest’ part of the above code to an object, and then plotting it. For example, with *Acer rubrum*:

```
K.acru <- Kest(acru.ppp, correction="best")  
plot(K.acru, main="Ripley's K function for Acer rubrum")
```

## Ripley's K function for *Acer rubrum*



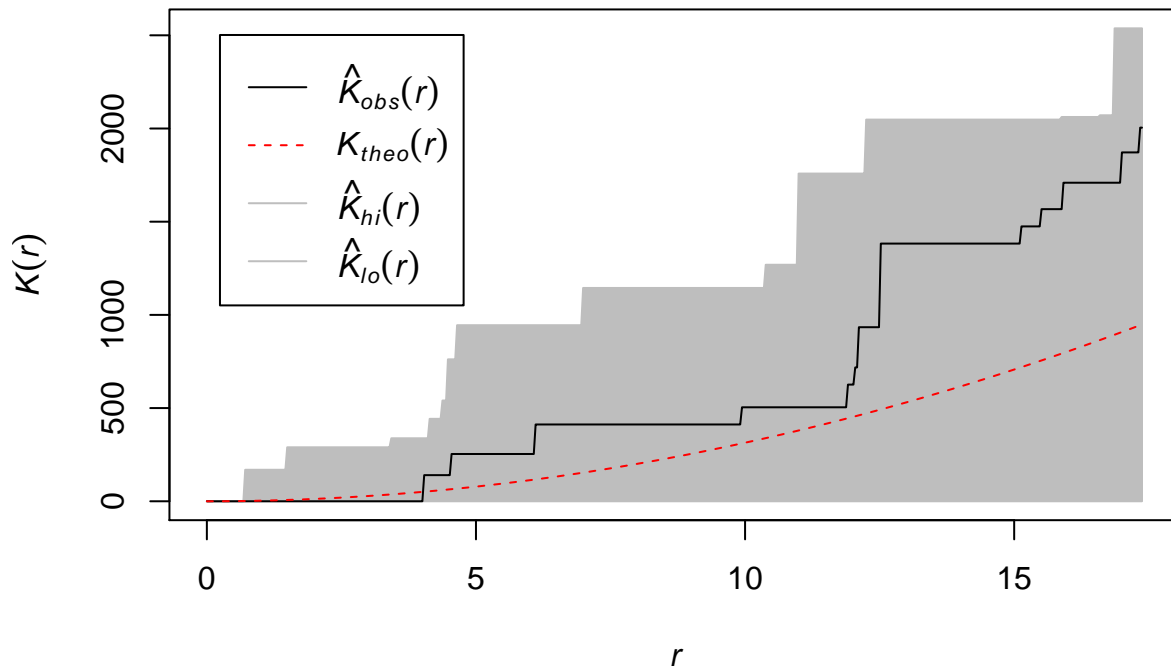
*Question 8: Describe how the individuals of each tree species are clustered, dispersed, or random across a range of scales. In addition, plot the K functions computed no edge correction - how would your interpretation of the results change for each species? You don't need to show the final plots, but rather describe the differences you notice.*

To determine the significance of the patterns in K functions, it is useful to use simulations or randomizations of the data to attain confidence intervals for complete spatial randomness. We can use the 'envelope' function for this - the gray area will represent the 95% confidence intervals for complete spatial randomness based on these simulations of data.

```
plot(envelope(acru.ppp, Kest, correction="best", nsim = 100)) #100 simulations
plot(envelope(qual.ppp, Kest, correction="best", nsim = 100)) #100 simulations
```

For *Quercus alba* your output figure and in the R console should be something like this:

## Ripley's K for Quercus alba with 95% Confidence Intervals of CSR



*Question 9: Create the envelopes with 1000 simulations and explain how this increase in simulations appears to affect the confidence intervals.*

The following code produces Ripley's K output for both species together- it shows whether individuals of one species tend to cluster near trees of another species. Again, higher positive values of Kest indicate clustering, and lower negative values indicate segregation among events, this time by type (i.e, tree species).

```
RipleyTrees.ppp <- as(RipleyTrees,"ppp")
plot(envelope(RipleyTrees.ppp, Kcross, i="ACRU", j="QUAL",nsim=100))
```

*Question 10: Does there appear to be a relationship between the two tree species in this focal plot? Or do they follow Complete Spatial Randomness? If there is a significant pattern at any scale, describe it.*