

Microbial Community Analysis Using Colony Phenotype Database

Student Handout

Objectives

- Manipulate large datasets to conduct community-level ecological analyses
- Use community-level data to address questions about insect microbiomes
- Use Excel or RStudio programs to calculate community ecology variables
- Compare microbial community using community ecology variables

Introduction

Microbiomes are the communities of microbes (bacteria, viruses, fungi and archaea) living symbiotically with all metazoans. In the past decade, both interest and research on microbiomes, including their implications for human health, have increased dramatically (Christian *et al.* 2015, Costello *et al.* 2012, McFall-Ngai *et al.* 2013, The Human Microbiome Consortium 2012, Young 2016). Insects have been used as model species to study the importance of microbiomes, because of their ease of use and the fact that microbial communities play diverse roles in insects (Engel and Moran 2013).

The data that are collected in any microbiome study consists of lists of the taxonomic units identified and their abundance. The same types of data are evaluated in an ecological community analysis, but now the communities are the collections of microbes that constitute different microbiomes. The community variables, “species” richness and relative abundance, are the same and the statistical methods used to compare communities, diversity and difference indices, also are the same. Perhaps the simplest measure of community structure used by ecologists is “species” or taxon richness, a count of the number of unique taxa in a sample. However, species richness does not consider the relative abundance of species in a community. Imagine two communities with five different species. In one community, all of the species have the same relative abundance. In the other community, one species dominates comprising 95% of individuals in the community. The other four species are very rare. Based on species richness as a measure of community structure, these two communities are the same, although they are clearly very different. As a result, ecologists use other species diversity indices that consider both the number of species and the relative abundance of species in a community. Two common indices are the Simpson Index and the Shannon-Weaver Index. Communities with greater numbers of species and higher evenness (i.e., similar relative abundance of species within a community) are considered more diverse. Finally, measures of species richness and species diversity do not consider the identity of species in a community. So, communities could have the same level of species diversity, but have completely different species. Measure of community similarity, such as the Bray-Curtis Index, compare the similarity (or dissimilarity) between two communities based on the identity of species in the communities, as well as their relative abundances. For more information on indices of species diversity and measures of community similarity, see Krebs (1999).

In this study, bean beetle gut microbiome data were collected by undergraduate students using the protocols developed by Cole *et al.* (2018). Three types of data were collected: colony phenotypes from cultured bacteria, 16s rRNA gene sequencing of specific bacterial colonies, and whole community 16s rRNA gene sequencing, but we will limit our analyses to the colony phenotype and colony-based 16s data.

Microbial Community Analysis Using Colony Phenotype Database

Questions

Using data from the colony phenotype database and the analyses described below, answer the following questions.

1. Based on the diversity indices that you calculated, which treatment had the highest (lowest) diversity?
2. Does the answer depend on the measure of species (taxon) diversity that you use?
3. Is there a relationship between number of samples and taxonomic diversity? If so, what might explain this?
4. Which communities are most similar (different)?

Database description

This database contains data for the microbial community of bean beetles based on colony phenotypes of individual bacterial colonies cultured from bean beetle homogenates plated on different media. The bacteria cultured from each beetle are represented by multiple rows of data with each row representing a different combination of phenotypic characters. A unique combination of phenotypic characters is taken to represent a unique bacterial taxa.

Access the database at <https://www.beanbeetles.org/microbiome/phenotype-database-search/>.

The database allows you to limit your search by bean host species, sex, life cycle stage, media on which bacteria were grown, and colony phenotype. EMB and PEA plates select for gram negative and gram positive bacteria, respectively. So, a colony with a particular combination of phenotypic characters on an EMB plate is a different bacterial taxa than a colony with the same combination of characters on a PEA plate. The same is not true if we include blood agar or nutrient agar plates in our sample. Those plates could be considered independently.

We want to limit our search to PEA and EMB plates (or Nutrient Agar only). After limiting by media type, click "Submit."

Microbiome Phenotype Database Search

Use the form below to search the microbiome database for bean beetles based the phenotype of individual bacterial colonies cultured from bean beetle homogenates plated on different media. Selecting checkboxes will limit your search to those samples. To view all data, just click the search button.

Sample	Phenotype
Bean Type ▼	Color ▼
Life Cycle Stage ▼	Gloss ▼
Sex ▼	Form ▼
Media ▼	Elevation ▼
<input type="checkbox"/> Nutrient Agar <input type="checkbox"/> Blood Agar (fastidious bacteria) <input checked="" type="checkbox"/> EMB Agar (selective for gram negative bacteria) <input checked="" type="checkbox"/> PEA Agar (selective for gram positive bacteria)	
<input type="button" value="Reset"/> <input type="button" value="Submit"/>	

Select "EMB Agar" and "PEA Agar"

Downloading Data

While we can view the data on the website, we want to download the data to manipulate. Click the download link to download a csv file with the data. Then, re-save the file as an Excel file and rename the tab "raw data."

Phenotype Database Search Results

[Download search results](#)

Show 10 entries

experiment_id	sample_id	sex	stage	media	color	gloss	form	elevation	CFU
EMORY_SI	EMORY_SI_1	female	adult	PEA	offwhite	matte	circular	convex	4390
EMORY_SI	EMORY_SI_2	female	adult	PEA	brown	matte	circular	raised	160000
EMORY_SI	EMORY_SI_3	male	adult	PEA	white	shiny	circular	convex	420
EMORY_SI	EMORY_SI_4	male	adult	PEA	white	shiny	circular	convex	250
EMORY_SI	EMORY_SI_5	female	adult	PEA	white	matte	circular	raised	133
EMORY_SI	EMORY_SI_6	female	adult	PEA	white	shiny	circular	raised	100000
EMORY_SI	EMORY_SI_7	male	adult	PEA	red	shiny	circular	umbonate	190
EMORY_SI	EMORY_SI_8	male	adult	PEA	red	shiny	circular	raised	100000
EMORY_SI	EMORY_SI_9	male	adult	PEA	white	shiny	circular	raised	60
EMORY_SI	EMORY_SI_10	female	adult	PEA	offwhite	shiny	circular	flat	140
EMORY_SI	EMORY_SI_11	male	adult	PEA	yellow	shiny	circular	raised	4350
EMORY_SI	EMORY_SI_12	male	adult	EMB	brown	shiny	circular	raised	20
EMORY_SI	EMORY_SI_13	male	adult	PEA	offwhite	shiny	circular	flat	10
EMORY_SI	EMORY_SI_14	female	adult	EMB	clear	shiny	circular	raised	120
EMORY_SI	EMORY_SI_15	male	adult	PEA	yellow	matte	circular	raised	200
EMORY_SI	EMORY_SI_16	female	adult	EMB	red	shiny	circular	raised	6
EMORY_SI	EMORY_SI_17	male	adult	PEA	white	shiny	circular	raised	3920
EMORY_SI	EMORY_SI_18	male	adult	PEA	white	shiny	circular	raised	2
EMORY_SI	EMORY_SI_19	female	adult	PEA	white	shiny	circular	convex	100
EMORY_SI	EMORY_SI_20	male	adult	EMB	red	shiny	circular	raised	1220
EMORY_SI	EMORY_SI_21	male	adult	PEA	orange	shiny	irregular	umbonate	2720
EMORY_SI	EMORY_SI_22	male	adult	EMB	yellow	matte	circular	raised	15000
EMORY_SI	EMORY_SI_23	female	adult	PEA	white	matte	circular	umbonate	11
EMORY_SI	EMORY_SI_24	male	adult	PEA	yellow	matte	circular	raised	830
EMORY_SI	EMORY_SI_25	female	adult	EMB	orange	shiny	circular	raised	80
EMORY_SI	EMORY_SI_26	female	adult	PEA	white	matte	circular	raised	7
EMORY_SI	EMORY_SI_27	male	adult	EMB	red	shiny	circular	flat	2800
EMORY_SI	EMORY_SI_28	female	adult	EMB	red	shiny	circular	raised	8
EMORY_SI	EMORY_SI_29	female	adult	PEA	offwhite	shiny	circular	convex	112100
EMORY_SI	EMORY_SI_30	female	adult	EMB	brown	shiny	circular	raised	400
EMORY_SI	EMORY_SI_31	female	adult	PEA	offwhite	shiny	circular	flat	1840
EMORY_SI	EMORY_SI_32	male	adult	PEA	white	shiny	circular	convex	250

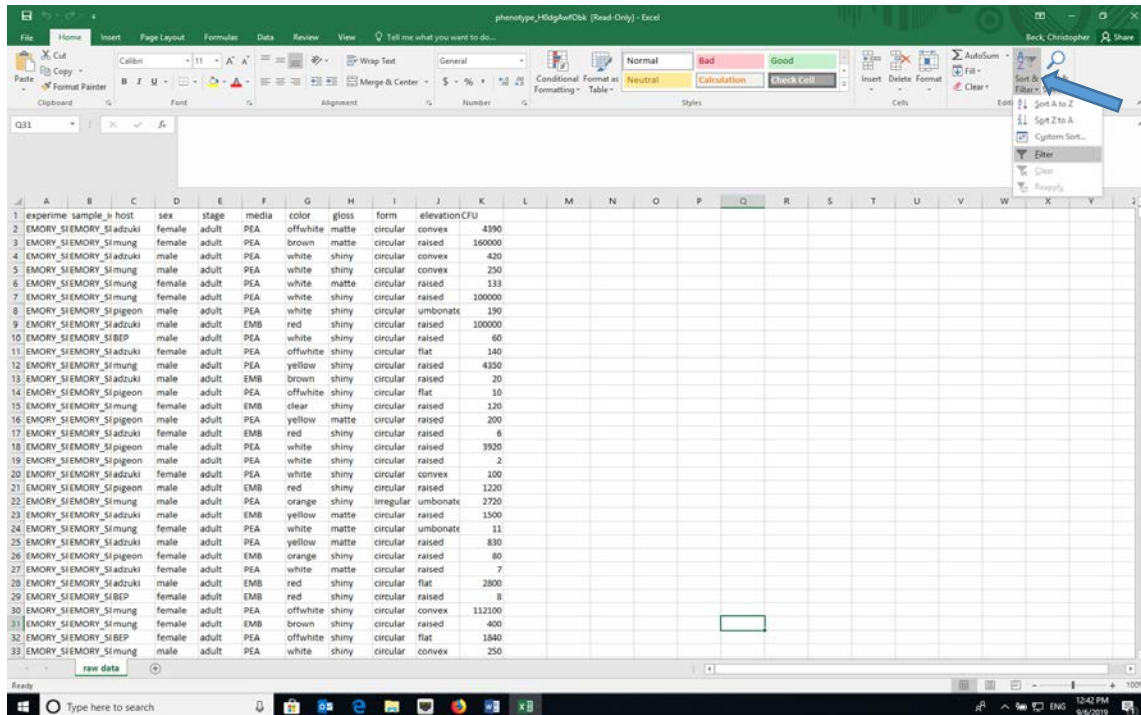
phenotype_HRdgAw0bk (Read-Only) - Excel

Double click and relabel tab as "raw data"

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z
1	experime	sample_id	host	sex	stage	media	color	gloss	form	elevation	CFU															
2	EMORY_SI	EMORY_SI_1	adzuki	female	adult	PEA	offwhite	matte	circular	convex	4390															
3	EMORY_SI	EMORY_SI_2	adzuki	female	adult	PEA	brown	matte	circular	raised	160000															
4	EMORY_SI	EMORY_SI_3	adzuki	male	adult	PEA	white	shiny	circular	convex	420															
5	EMORY_SI	EMORY_SI_4	adzuki	male	adult	PEA	white	shiny	circular	convex	250															
6	EMORY_SI	EMORY_SI_5	adzuki	female	adult	PEA	white	matte	circular	raised	133															
7	EMORY_SI	EMORY_SI_6	adzuki	female	adult	PEA	white	shiny	circular	raised	100000															
8	EMORY_SI	EMORY_SI_7	pigeon	male	adult	PEA	white	shiny	circular	umbonate	190															
9	EMORY_SI	EMORY_SI_8	adzuki	male	adult	EMB	red	shiny	circular	raised	100000															
10	EMORY_SI	EMORY_SI_9	BEP	male	adult	PEA	white	shiny	circular	raised	60															
11	EMORY_SI	EMORY_SI_10	adzuki	female	adult	PEA	offwhite	shiny	circular	flat	140															
12	EMORY_SI	EMORY_SI_11	adzuki	male	adult	PEA	yellow	shiny	circular	raised	4350															
13	EMORY_SI	EMORY_SI_12	adzuki	male	adult	EMB	brown	shiny	circular	raised	20															
14	EMORY_SI	EMORY_SI_13	pigeon	male	adult	PEA	offwhite	shiny	circular	flat	10															
15	EMORY_SI	EMORY_SI_14	adzuki	female	adult	EMB	clear	shiny	circular	raised	120															
16	EMORY_SI	EMORY_SI_15	pigeon	male	adult	PEA	yellow	matte	circular	raised	200															
17	EMORY_SI	EMORY_SI_16	adzuki	female	adult	EMB	red	shiny	circular	raised	6															
18	EMORY_SI	EMORY_SI_17	pigeon	male	adult	PEA	white	shiny	circular	raised	3920															
19	EMORY_SI	EMORY_SI_18	pigeon	male	adult	PEA	white	shiny	circular	raised	2															
20	EMORY_SI	EMORY_SI_19	adzuki	female	adult	PEA	white	shiny	circular	convex	100															
21	EMORY_SI	EMORY_SI_20	pigeon	male	adult	EMB	red	shiny	circular	raised	1220															
22	EMORY_SI	EMORY_SI_21	adzuki	male	adult	PEA	orange	shiny	irregular	umbonate	2720															
23	EMORY_SI	EMORY_SI_22	adzuki	male	adult	EMB	yellow	matte	circular	raised	15000															
24	EMORY_SI	EMORY_SI_23	adzuki	female	adult	PEA	white	matte	circular	umbonate	11															
25	EMORY_SI	EMORY_SI_24	adzuki	male	adult	PEA	yellow	matte	circular	raised	830															
26	EMORY_SI	EMORY_SI_25	pigeon	female	adult	EMB	orange	shiny	circular	raised	80															
27	EMORY_SI	EMORY_SI_26	adzuki	female	adult	PEA	white	matte	circular	raised	7															
28	EMORY_SI	EMORY_SI_27	adzuki	male	adult	EMB	red	shiny	circular	flat	2800															
29	EMORY_SI	EMORY_SI_28	BEP	female	adult	EMB	red	shiny	circular	raised	8															
30	EMORY_SI	EMORY_SI_29	adzuki	female	adult	PEA	offwhite	shiny	circular	convex	112100															
31	EMORY_SI	EMORY_SI_30	adzuki	female	adult	EMB	brown	shiny	circular	raised	400															
32	EMORY_SI	EMORY_SI_31	adzuki	female	adult	PEA	offwhite	shiny	circular	flat	1840															
33	EMORY_SI	EMORY_SI_32	adzuki	male	adult	PEA	white	shiny	circular	convex	250															

Data manipulation

1. The raw data have some missing values for phenotypic characters or values of zero for CFU (colony forming units) (which represents missing data for CFU). These rows need to be deleted. The easiest way to do this is with the Filter function in Excel. In the “raw data” worksheet, turn on filtering by selecting Filter under Sort and Filter.



Then, for each of the phenotypic characters, click the down arrow in the column heading and unselect “Blank” at the bottom of the list and unselect “0” in the CFU column.

The top screenshot shows an Excel spreadsheet with a table containing data for 33 rows. The columns are labeled A through Y. The data includes various attributes such as 'sex', 'stage', 'media', 'color', 'gloss', 'form', 'elevat', and 'CPU'. A 'Filter by Color' dialog box is open, showing a list of colors with checkboxes. The 'color' column is selected, and the 'Filter by Color' dialog box is open, showing a list of colors with checkboxes. The 'color' column is selected, and the 'Filter by Color' dialog box is open, showing a list of colors with checkboxes.

The bottom screenshot shows the same Excel spreadsheet, but with a different 'Filter by Color' dialog box open. This dialog box is for the 'CPU' column, showing a list of CPU values with checkboxes. The 'CPU' column is selected, and the 'Filter by Color' dialog box is open, showing a list of CPU values with checkboxes.

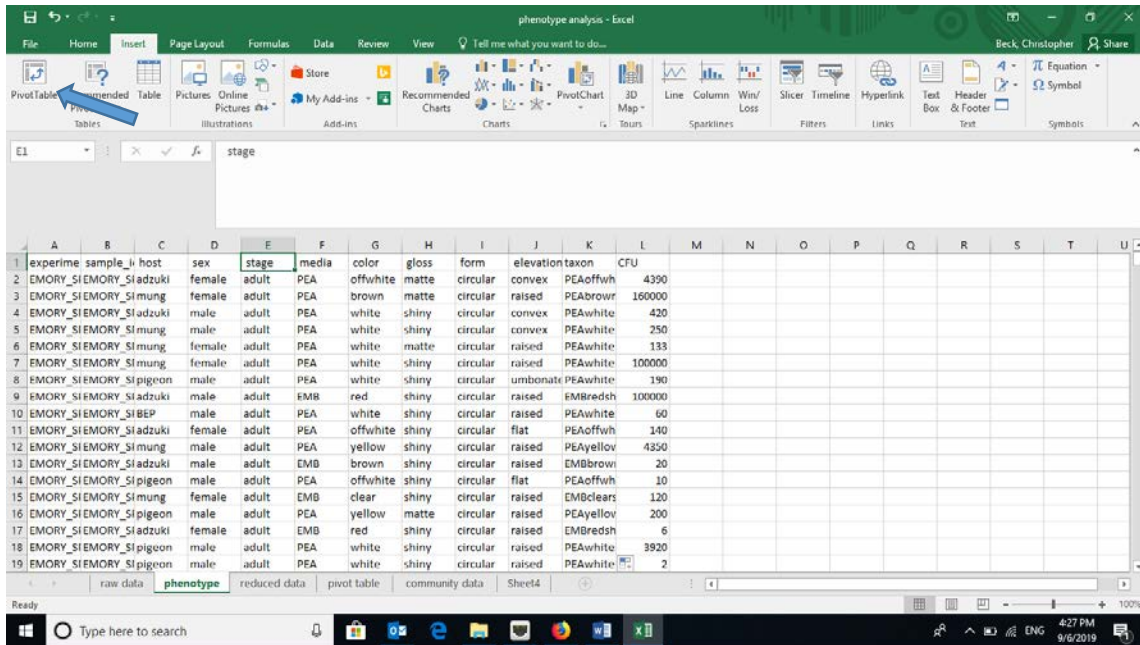
Next, select all (CTRL-A or CMD-A) and copy and paste into a new worksheet and name that sheet “phenotype”. The “phenotype” worksheet represents the cleaned raw data. If you will be analyzing the data in R later, copy and paste a second time into another blank sheet and name that sheet “community”.

- With our new “phenotype” dataset, we need to define a bacterial “taxon” based on the combination of media and the four phenotypic characters (color, gloss, form and elevation). One way to do this is to create a “taxon” name by concatenating the media and the four different phenotypic traits. You can do this using the CONCATENATE function in Excel (=CONCATENATE(F2,G2,H2,I2,J2)). After you create the “taxa” names, you might want to select the column and then re-paste it in the same column by pasting values (using paste special) to remove the formula.

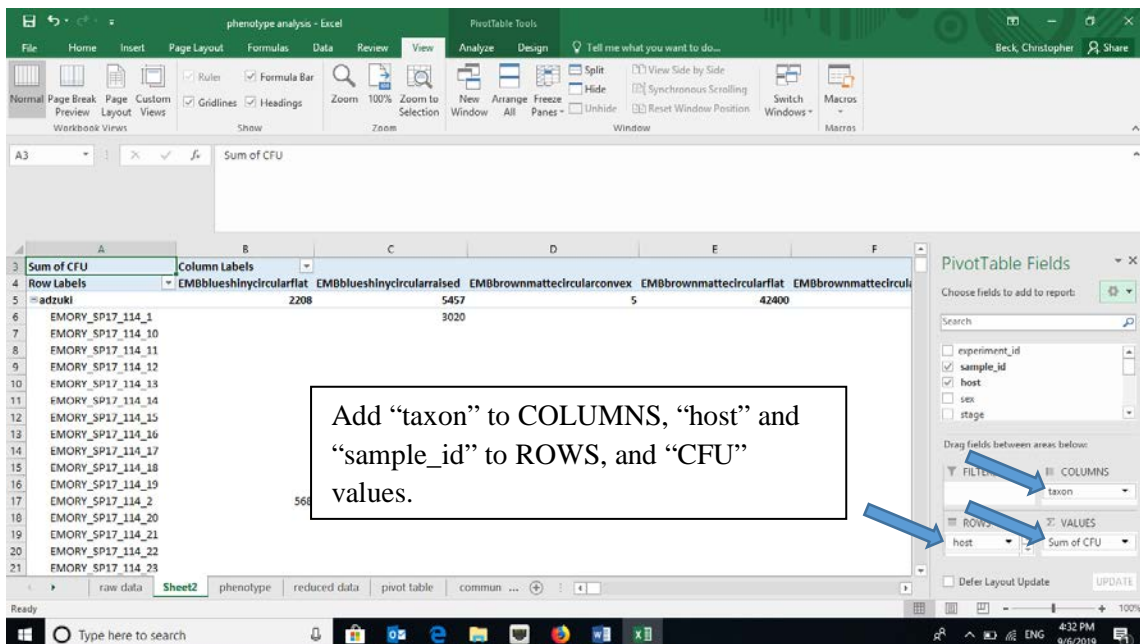
Insert a new column for “taxon” then in the first blank cell of the column insert the formula =CONCATENATE(F2,G2,H2,I2,J2). Then, copy the formula down for the rest of the column.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U
1	experime	sample_i	host	sex	stage	media	color	gloss	form	elevation	taxon	CFU									
2	EMORY_SI	EMORY_SI	adzuki	female	adult	PEA	offwhite	matte	circular	convex	=CONCATENATE(F2,G2,H2,I2,J2)	150000									
3	EMORY_SI	EMORY_SI	adzuki	female	adult	PEA	brown	matte	circular	raised		250									
4	EMORY_SI	EMORY_SI	adzuki	male	adult	PEA	white	shiny	circular	convex		133									
5	EMORY_SI	EMORY_SI	mung	male	adult	PEA	white	shiny	circular	convex		100000									
6	EMORY_SI	EMORY_SI	mung	female	adult	PEA	white	matte	circular	raised		190									
7	EMORY_SI	EMORY_SI	mung	female	adult	PEA	white	shiny	circular	raised		100000									
8	EMORY_SI	EMORY_SI	pigeon	male	adult	PEA	white	shiny	circular	umbonate		60									
9	EMORY_SI	EMORY_SI	adzuki	male	adult	EMB	red	shiny	circular	raised		140									
10	EMORY_SI	EMORY_SI	IBEP	male	adult	PEA	white	shiny	circular	raised		4350									
11	EMORY_SI	EMORY_SI	adzuki	female	adult	PEA	offwhite	shiny	circular	flat		20									
12	EMORY_SI	EMORY_SI	mung	male	adult	PEA	yellow	shiny	circular	raised		10									
13	EMORY_SI	EMORY_SI	adzuki	male	adult	EMB	brown	shiny	circular	raised		120									
14	EMORY_SI	EMORY_SI	pigeon	male	adult	PEA	offwhite	shiny	circular	flat		200									
15	EMORY_SI	EMORY_SI	mung	female	adult	EMB	clear	shiny	circular	raised		6									
16	EMORY_SI	EMORY_SI	pigeon	male	adult	PEA	yellow	matte	circular	raised		3920									
17	EMORY_SI	EMORY_SI	adzuki	female	adult	EMB	red	shiny	circular	raised		2									
18	EMORY_SI	EMORY_SI	pigeon	male	adult	PEA	white	shiny	circular	raised											
19	EMORY_SI	EMORY_SI	pigeon	male	adult	PEA	white	shiny	circular	raised											

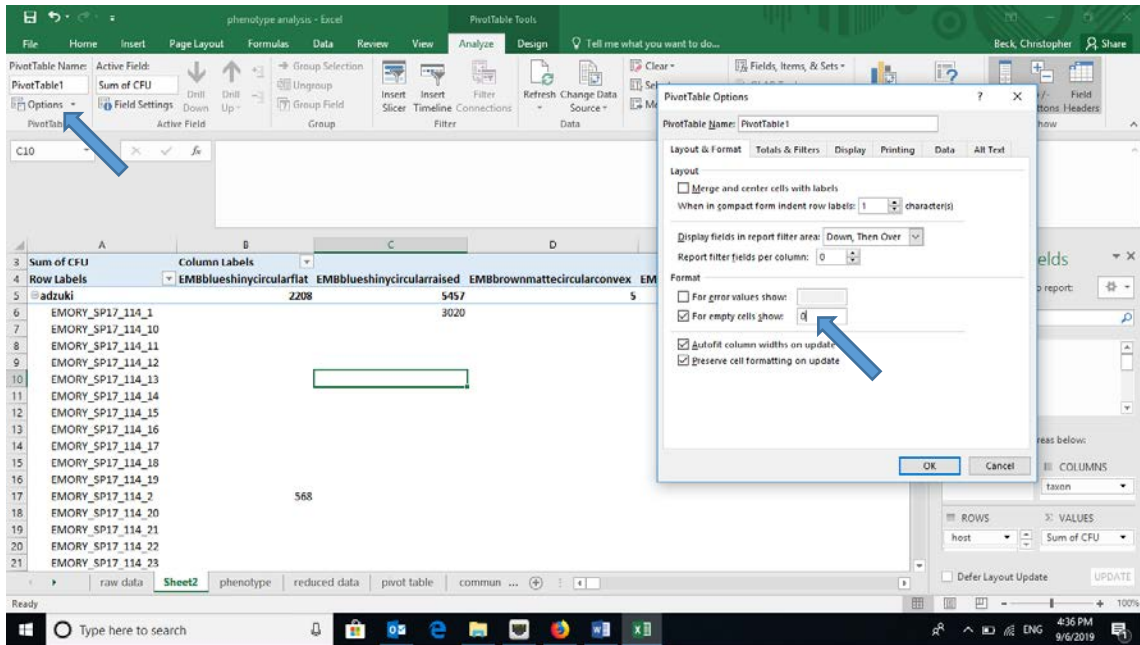
- Now, we need to consolidate the data for each host species, each sex, or the combination of the two by the bacterial taxa. The easiest way to do this is with the Pivot Table function in Excel.
- When clicked on a cell within the data, create a Pivot Table (Insert -> Pivot Table OR Data -> Summarize with Pivot Table, depending on your version of Excel). Make sure that the data source includes the top row, which has the column headings.



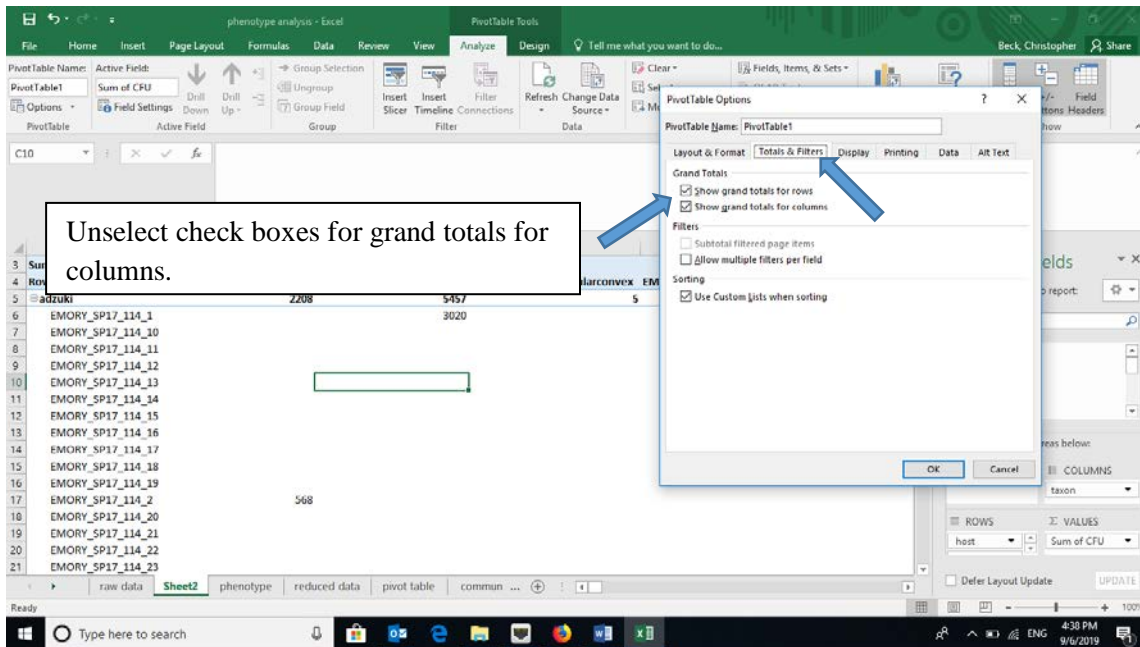
Set the host and sample_id as the rows as these represent the treatment and individual communities, respectively. The new “taxon” column should be the columns in the pivot table. The Values should be a SUM of the CFU (colony-forming units, a measure of density), which will be shown as “SUM of CFU” using the Options menu.



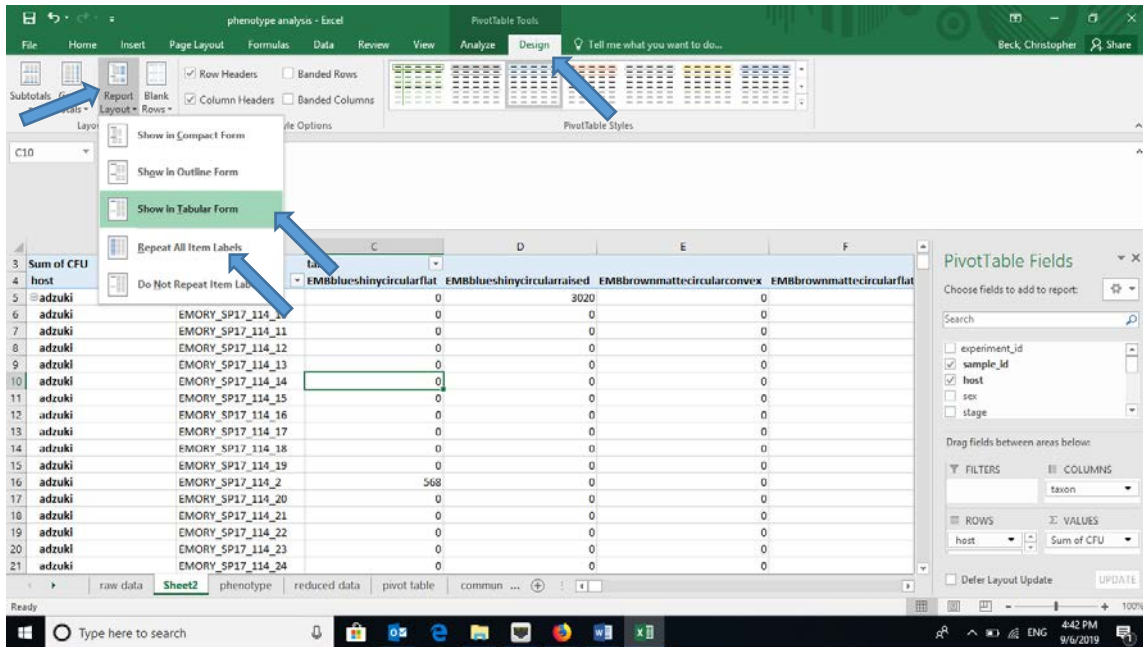
5. You can add zeros to all of the empty cells in the Pivot Table using the Options menu.



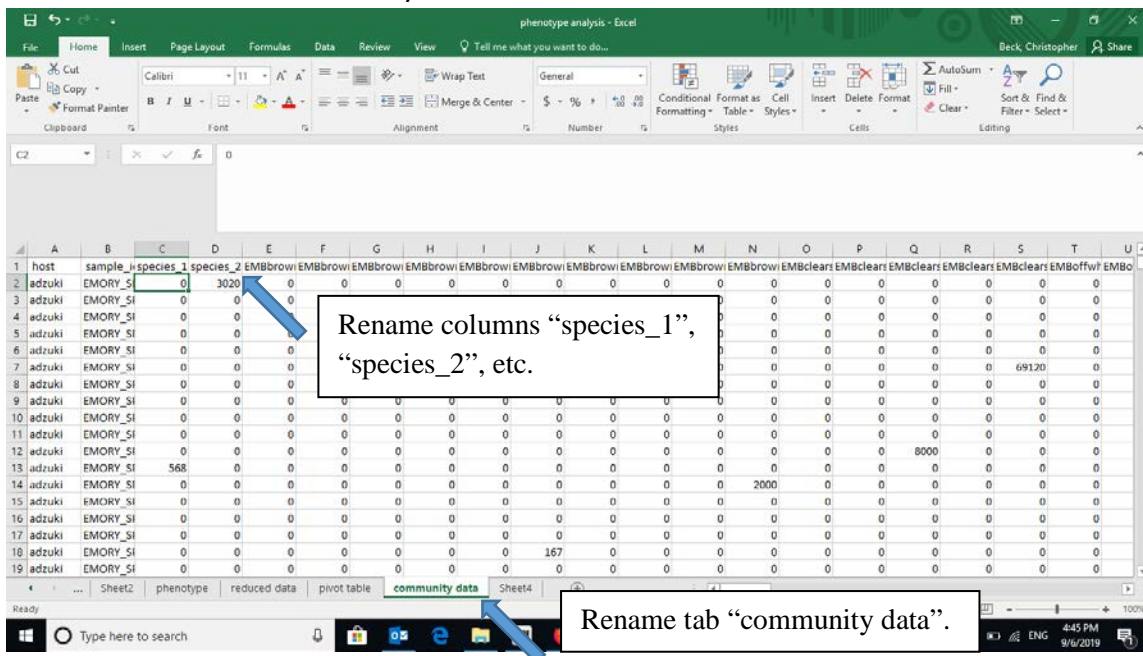
- Remove the Grand totals for Columns using the Options menu or the Design tab depending of your version of Excel.



- To get the treatment data to repeat for each sample, in the Design tab, select "Report Layout" and choose "Show in Tabular form" and "Repeat All Item Labels".



- Copy and paste (as values) the pivot table to a new worksheet and remove any extra rows at the top. Each of the columns in this new worksheet represents a unique bacterial taxa. The exact phenotype doesn't matter, so we are going to rename them species_1, species_2, Title this worksheet tab "community data".



Calculating diversity indices

- Species (taxon) richness – the number of unique species (taxa) in a sample

- a. Although you could manually count the number of cells with values greater than zero for each treatment, using the COUNTIF formula in Excel is easier (e.g., =COUNTIF(range,">0") where "range" is the range of cells in the spreadsheet containing the data, such as "C2:EW2").

phenotype analysis - Excel

FileHomeInsertPage LayoutFormulasDataReviewViewTell me what you want to do...

CutCopyFormat PainterClipboardFontAlignmentNumberStylesCellsEditing

Calibri11FontWrap TextGeneralConditional FormattingFormat as TableCell StylesInsertDeleteFormatFillSort & FilterFind & Select

Countif=-COUNTIF(C2:EW2,>0)

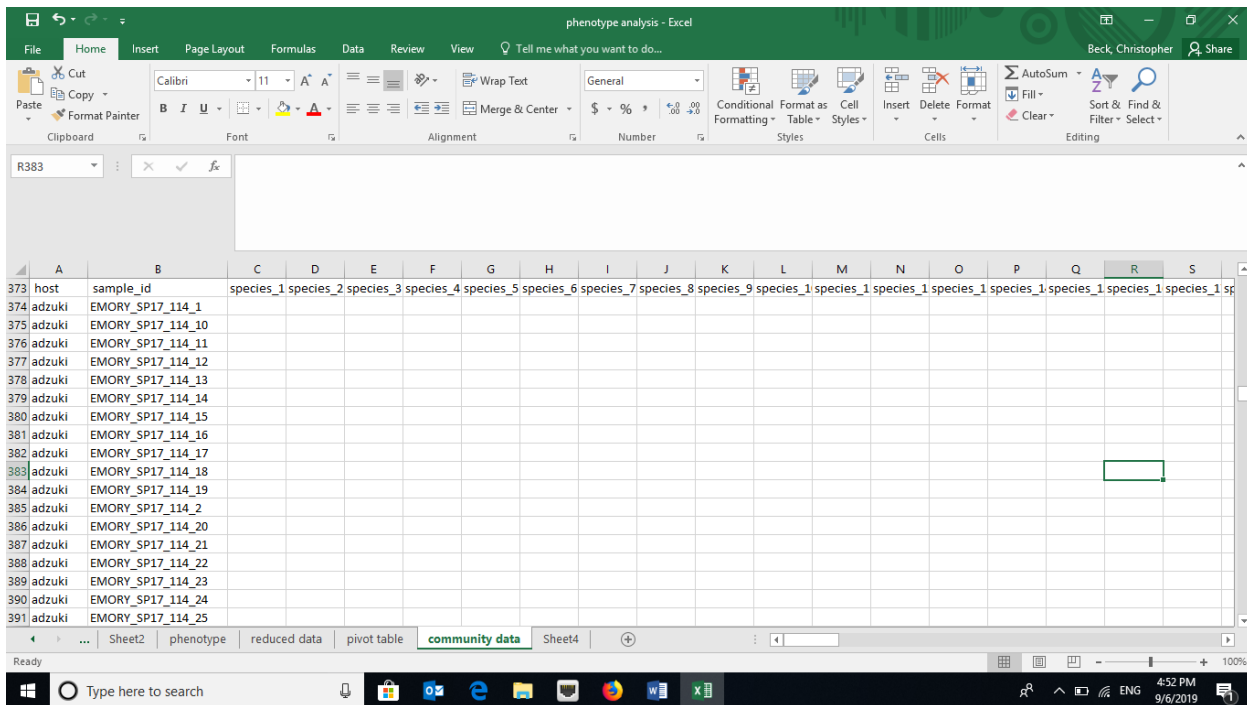
	EK	EL	EM	EN	EO	EP	EQ	ER	ES	ET	EU	EV	EW	EX	EY	EZ	FA	FB	FC	FD	FE
1	PEAwhite	PEAwhite	PEAyellow	PEAyellow	PEAyellow	PEAyellow	PEAyellow	PEAyellow	PEAyellow	PEAyellow	PEAyellow	PEAyellow	PEAyellow	Abundance	Richness						
2	0	0	0	0	0	0	0	0	0	0	0	0	0	7714	=COUNTIF(C2:EW2,>0)						
3	700	0	0	0	0	400	0	0	0	0	0	0	0	4000	=COUNTIF(range, criteria)						
4	247	0	0	0	0	0	0	0	0	0	0	0	0	34247	3						
5	0	0	0	0	0	0	0	0	0	0	0	0	0	361	1						
6	0	0	0	0	0	0	0	0	0	0	0	0	0	52	2						
7	81	0	0	0	0	0	0	0	0	0	0	0	0	117417	4						
8	0	0	0	0	0	0	0	0	0	0	0	0	0	4147	2						
9	0	0	0	0	0	0	0	0	0	0	0	0	0	100	1						
10	0	0	0	0	0	0	0	0	0	0	0	0	0	2040	2						
11	0	0	0	0	0	0	0	0	0	0	8220	0	0	16391	5						
12	0	0	0	0	0	0	0	0	0	0	0	0	0	13014	3						
13	0	0	0	0	0	0	0	0	0	0	0	0	0	2222	4						
14	0	0	0	0	0	0	0	0	0	0	0	0	0	2774	4						
15	0	0	0	0	0	0	0	0	0	0	0	0	0	276720	6						
16	0	0	0	0	0	0	0	0	0	0	0	0	0	30615	3						
17	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1						
18	0	0	0	0	0	0	0	0	0	0	0	0	0	349	3						
19	0	0	0	0	920	0	0	0	0	0	0	0	0	1494	3						

Sheet2phenotypereduced datapivot tablecommunity dataSheet4

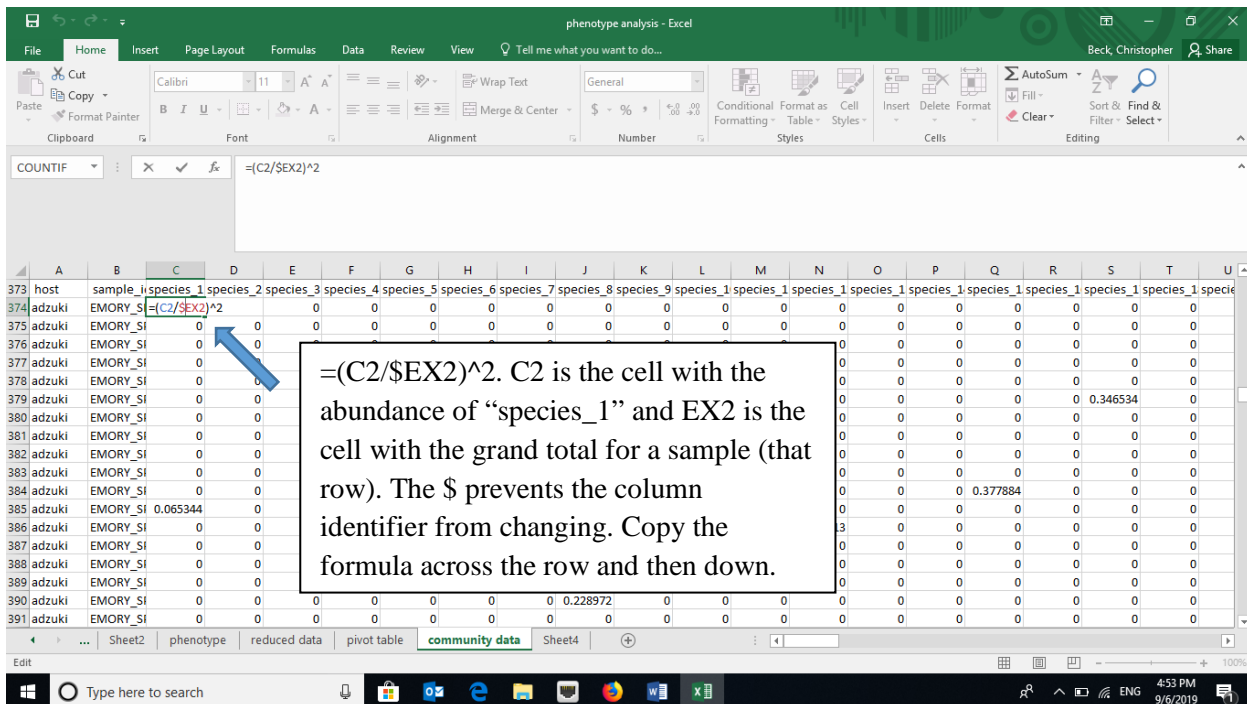
Type here to search

4:49 PM 9/6/2019

2. Simpson Index – the Simpson Index incorporates both species (taxon) richness and species (taxon) evenness.
 - a. $D = \sum (n/N)^2$, where n =number of individuals of a particular species (taxon) and N =total number of individuals in a sample. D increases as diversity decreases, which is counterintuitive. A reciprocal or inverse index would be more intuitive and are easily calculated.
 - b. Reciprocal Simpson = $1/D$ and scales so the maximum value is the species richness of a community.
 - c. Inverse Simpson = $1-D$ and scales to a maximum value of 1.0.
 - d. Create a new data array below the original using the same row labels (treatment variables) and the same column labels (species).



- e. To calculate the proportion squared for each taxa, use the grand totals for each treatment. Using the Excel trick that \$ before a column or row prevents Excel from iterating when copying a formula makes this easy. For example, $= (C2/\$EX2)^2$. Copy the formula across the row and then down.



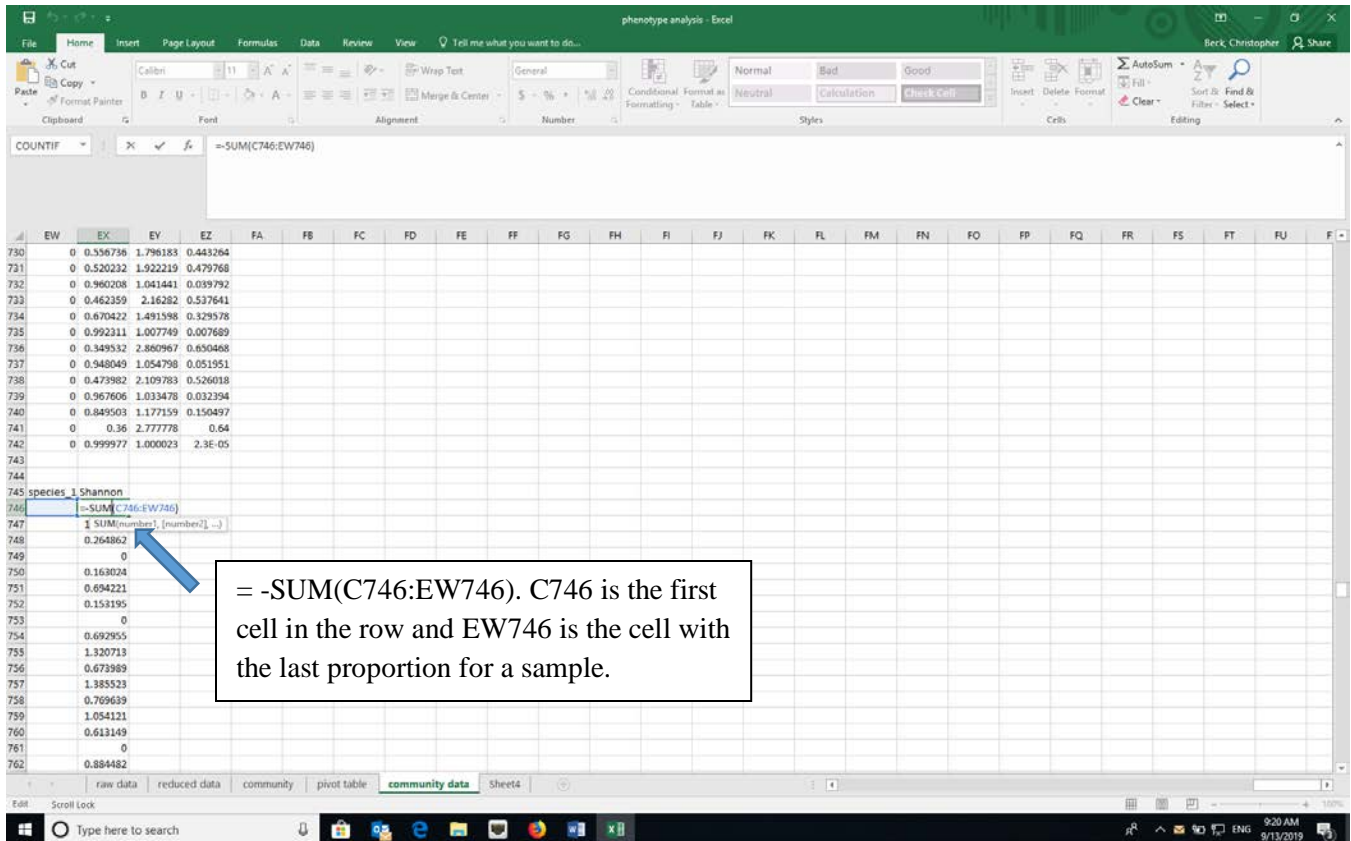
- f. Calculate the sum of the proportions squared ($= \text{SUM}$ in Excel for each row, a different microbial community) to calculate the Simpson Index.

The screenshot shows an Excel spreadsheet titled "phenotype analysis - Excel". The formula bar at the top displays the formula `=SUM(C374:EW374)`. A callout box points to cell EW374, explaining that the formula calculates the sum of the range C374:EW374, where C374 is the first cell in the row and EW374 is the cell with the last proportion for a sample. The spreadsheet contains a data table with columns labeled EK through FD and rows 373 through 391. The table includes columns for Simpson, Reciprocal, and Inverse indices.

g. Calculate the reciprocal (e.g., $=1/\text{EX374}$) and inverse Simpson (e.g., $=1/\text{EX374}$) using formulas in Excel.

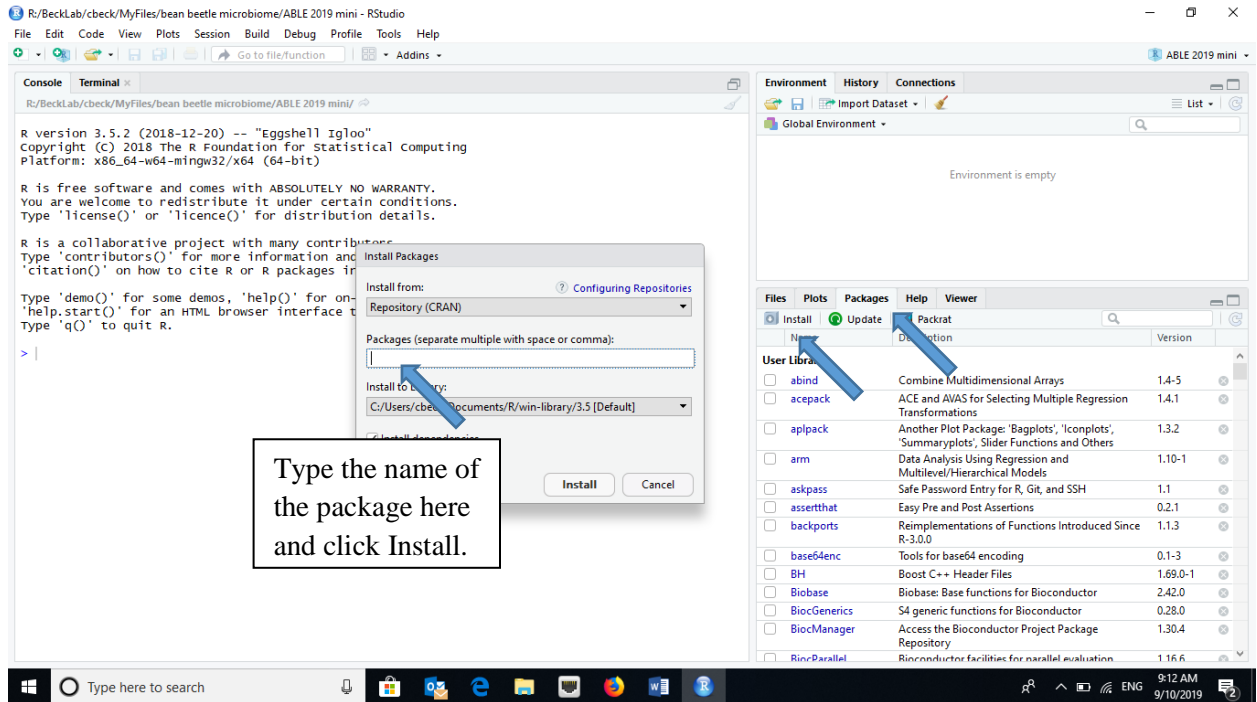
3. Shannon-Weaver (Shannon-Weiner) Index – also incorporates species (taxon) richness and species (taxon) evenness

- $H = -\sum p \ln p$, where p is the proportion of individuals of each species (taxon) in a community (i.e., n/N).
- Create a new data array below the original using the same row labels (treatment variables) and the same column labels (species).
- Using the grand totals for each community, calculate the proportions ($p \ln p$). Using the Excel trick that $\$$ before a column or row prevents Excel from iterating when copying a formula makes this easy.
- Note that $\ln p$ is undefined if $p=0$, so you can use an "IF" statement in Excel. For example, `=IF(C2>0,(C2/$EX2)*LN((C2/$EX2)),"")`

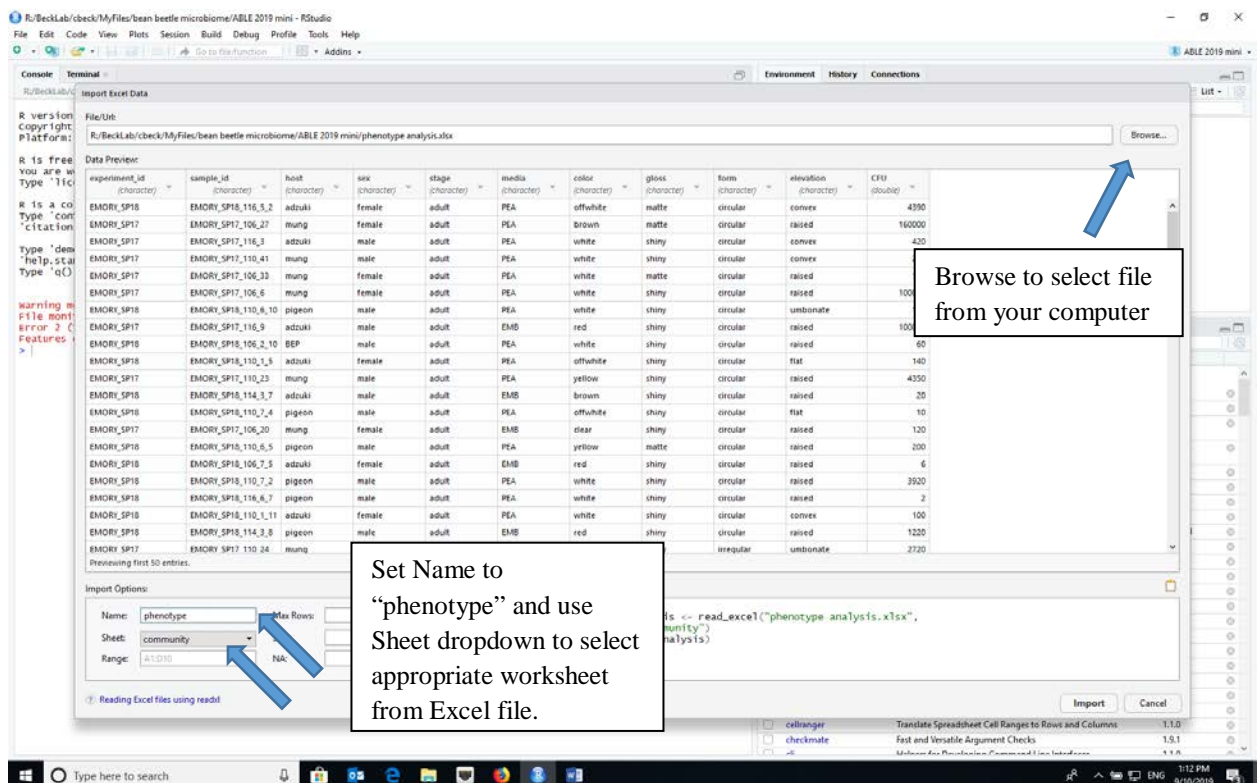
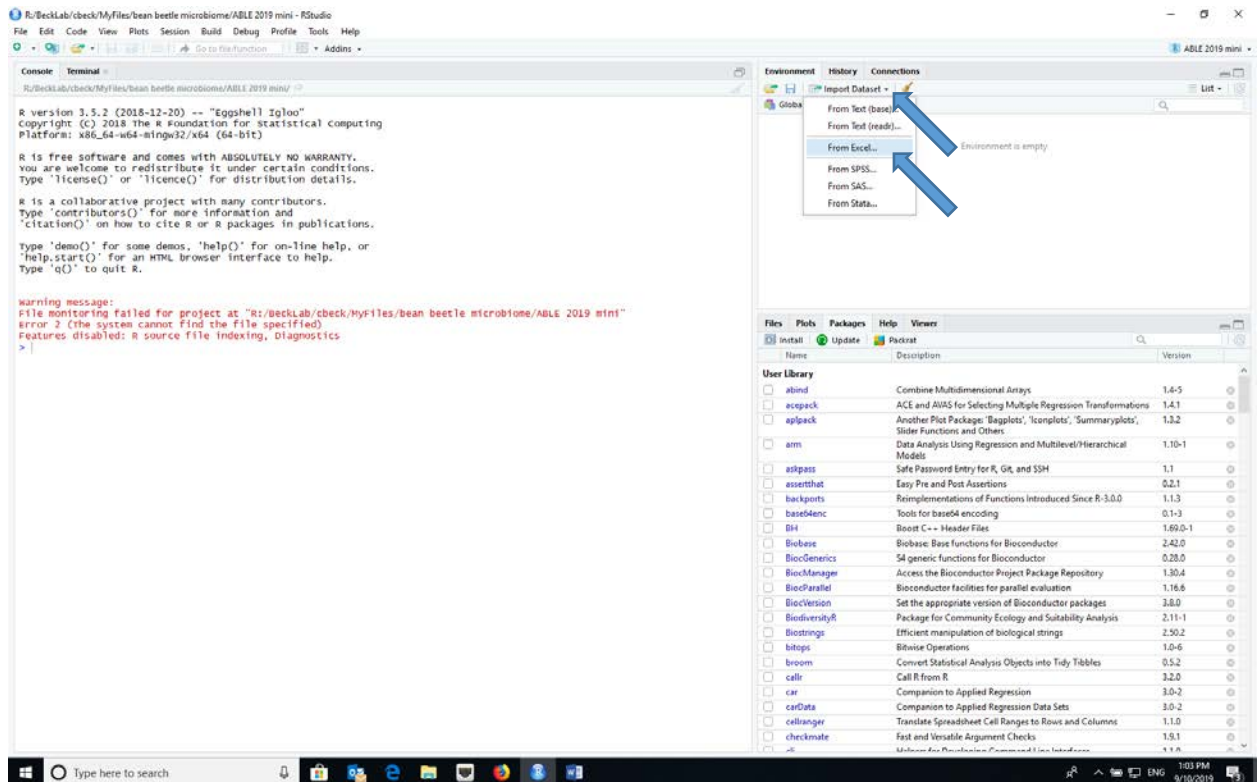


Data Manipulation in R

1. Open RStudio and create a new project using the New Project option under File and select for the new project to be in an existing folder where your data are.
2. Install the following packages either using the Packages tab in RStudio or the command `install.packages("name_of_package")` in the console. Note that BiodiversityR requires QuartzX on a Mac. If you are using a MacOS and don't have QuartzX, install it first and restart your computer before install these packages.
 - a. dplyr
 - b. reshape2
 - c. vegan
 - d. BiodiversityR
 - e. ggplot2



3. Load the packages listed above by clicking the checkboxes for the appropriate packages in the Packages tab or the command `library("name_of_package")` in the console.
4. Import the dataset ("community" that you created above in the Excel section) into RStudio.



5. Attach the imported dataset to the dataframe using the attach command in the console (`attach(phenotype)`)
6. Create a community matrix for each sample. Since each “taxon” is defined by the combination of media type and the four phenotypic characters, we can use the `dcast` function in the `reshape2` library, using the following command.


```
> comm_pheno=dcast(phenotype,sample_id~media+color+gloss+form+elevation,value.var = "CFU",fun.aggregate = sum)
```
7. The first column in the resulting data table is the `sample_id`. The `sample_id` needs to become the row name and then deleted, using the following commands.


```
> row.names(comm_pheno)<-(comm_pheno$sample_id)
> comm_pheno<-comm_pheno[,-1]
```
8. Now, we need to create an environment matrix with the sample metadata. First, we create a data table with the sample names and sample metadata, which are in the second through fourth columns of the “phenotype” dataframe. Second, we remove all of the duplicate values for the `sample_id` using the `distinct` function in the `dplyr` package. Third, we need the `sample_id` to become the row name. (Note that this causes an error message and for some reason if you delete the column with the row name, the row names disappear. In this case, we do not need to delete the `sample_id` column because we can specify the factor of interest.) Last, we need to set host and sex as factors.


```
> pheno_meta<-phenotype[,2:4]
> pheno_meta<-pheno_meta %>% distinct(sample_id, .keep_all = TRUE)
> row.names(pheno_meta)<-(pheno_meta$sample_id)
> pheno_meta$host<-as.factor(pheno_meta$host)
> pheno_meta$sex<-as.factor(pheno_meta$sex)
```

You can ignore the error message about setting row names on a tibble being deprecated.

Species accumulation curves

Species accumulation curves are often used to visualize whether all of the taxa in a community have been sampled. As the number of samples increases, the total number of unique species sampled should increase. However, the relationship between the number of samples and the number of unique species should asymptote. If so, we can say that we have sampled the majority of species in the community. However, if the slope of the relationship is steep, this suggests that the community is incompletely sampled.

For each treatment separately, create the data for the species accumulation curve using:

```
> Accum.label1<-specaccum(comm_pheno,method='exact',permutations = 100,
conditioned=TRUE,gamma='jack1',w=NULL,
subset=pheno_meta$factor_variable=="factor")
```


`factor_variable=="factor"` refers to the factor being evaluated such as host in our dataset, and the “factor” is one state of that variable. For example, `host=="mung"` would do a species accumulation curve for microbiome communities of beetles that emerged from mung beans. Note that two equals signs are necessary

You need to run the command above for each treatment group separately. Change the label and factor terms appropriately for additional analyses.

- a. Plot the first species accumulation curve using:

```
> plot(Accum.label1,col="red")
```

- b. Plot each additional species accumulation curve using:

```
> plot(Accum.label2, add=TRUE, col="blue")
```

If the second curve extends beyond the y-axis, replot the curves in the opposite order (i.e., plot curve 2 first and then curve 1).

Calculating diversity indices

You can calculate the diversity indices described above in the Excel exercise using functions in the [BiodiversityR](#) package.

1. Species Richness

```
> diversityresult(comm_pheno,index="richness",method="each site")
```

2. Simpson

```
> diversityresult(comm_pheno,index="Simpson",method="each site")
```

This calculates the inverse Simpson described above.

```
> diversityresult(comm_pheno,index="inverseSimpson",method="each site")
```

This calculates the reciprocal Simpson described above (confusing that it is called in the inverseSimpson).

3. Shannon

```
> diversityresult(comm_pheno,index="Shannon",method="each site")
```

4. To calculate all of the biodiversity indices and merge them with the metadata for future analysis. You can then use this dataset for analysis of differences between treatments with t-tests, ANOVAs, or their non-parametric equivalents.

```
> pheno_diversity<-diversityvariables(comm_pheno,pheno_meta)
```

Because the `diversityvariables` function calculates a range of diversity indices, sometimes an error occurs because a particular index cannot be calculated with the dataset (e.g., requires taking log of zero). If this is the case, you can use the `diversityresult` function above and place the results in a dataframe. Then, you can combine the dataframes with the metadata using the `cbind` function.

```
> Simpson<-diversityresult(comm_pheno,index="Simpson",method="each site")
> Shannon<-diversityresult(comm_pheno,index="Shannon",method="each site")
> pheno_diversity<-cbind(pheno_meta,Simpson,Shannon)
```

Calculating community similarity (distance)

Sometimes we are interested in how similar (or different) two communities are based on what species (taxa) are present and the relative abundance of those species (taxa) in the two communities. One of the most common measures of distance is the Bray Curtis Dissimilarity. Similarity can be measured as 1-BC.

$$BC_{ij} = 1 - \frac{2C_{ij}}{S_i + S_j}$$

Where:

- i & j are the two samples,
- S_i is the total number of specimens counted in sample i,
- S_j is the total number of specimens counted in sample j,
- C_{ij} is the sum of only the lesser counts for each taxa found in both sites.

Although Bray-Curtis Dissimilarity is often used in community ecology, it is not robust to incomplete sampling of the community (all taxa are not sampled) or unbalanced sampling (all treatments are not equally sampled). An alternative is the Morista-Horn Index of Dissimilarity (1- C_H). Morista-Horn Index of Similarity is

$$C_H = \frac{2 \sum_{i=1}^{D_{12}} \frac{X_i}{n} \frac{Y_i}{m}}{\sum_{i=1}^{D_1} \left(\frac{X_i}{n}\right)^2 + \sum_{i=1}^{D_2} \left(\frac{Y_i}{m}\right)^2}$$

Where

- D_1 =number of taxa in sample 1
- D_2 =number of taxa in sample 2
- D_{12} =number of taxa in shared in both communities
- X_i =number of individuals of taxon i in sample 1
- Y_i =number of individuals of taxon i in sample 2
- n =total number of individuals in sample 1
- m =total number of individuals in sample 2

So that X_i/n and Y_i/m are proportion of individuals of taxon i in each of the samples (communities).

To produce a matrix of all of the pair-wise distances between samples using the Bray Curtis index of distance, use the following command.

```
> vegdist(comm_pheno, method="bray", binary=FALSE, diag=FALSE, upper=FALSE)
```

To produce a matrix of all of the pair-wise distances between samples using the Morista-Horn index of distance.

```
> vegdist(comm_pheno, method="horn", binary=FALSE, diag=FALSE, upper=FALSE)
```

How different (similar) are the communities?

Adonis is an approach to testing whether communities differ based on a treatment. It is the equivalent of an analysis of variance, but comparing distance matrices. “community_adonis” stores the results of the analysis, “comm_pheno” is the community matrix, “factor_variable” is the treatment (e.g., host), and “pheno_meta” is the name of the dataset that has the treatment data for each community. In the code below, we use Morista-Horn to estimate distance between communities.

```
> community_adonis<-adonis2(comm_pheno ~ factor_variable, data = pheno_meta, method="horn")  
> community_adonis
```

Cited References

Christian N, Whitaker BK, Clay K. 2015. Microbiomes: unifying animal and plant systems through the lens of community ecology theory. *Front. Microbiol.* 6:1–15.

Cole MF, Acevedo-Gonzalez T, Gerardo NM, Harris EV, Beck CW. 2018. Effect of diet on bean beetle microbial communities. Article 3 In: McMahon K, editor. *Tested studies for laboratory teaching*.

Volume 39. Proceedings of the 39th Conference of the Association for Biology Laboratory Education (ABLE).

Costello EK, Stagaman K, Dethlefsen L, Bohannan BJM, Relman DA. 2012. The application of ecological theory toward an understanding of the human microbiome. *Science*. 336:1255–1262

Engel P, Moran NA. 2013. The gut microbiota of insects – diversity in structure and function. *FEMS Microbiol. Rev.* 37:699-735.

Krebs CJ. 1999. *Ecological Methodology*, 2nd edition. New York: Benjamin Cummings.

McFall-Ngai M, Hadfield MG, Bosch TCG, Carey HV, Domazet-Lošo T, Douglas AE, Dubilier N, Eberl G, Fukami T, Gilbert SF, Hentschel U, King N, Kjelleberg S, Knoll AH, Kremer N, Mazmanian SK, Metcalf JL, Nealson K, Pierce NE, Rawls JF, Reid A, Ruby EG, Rumpho M, Sanders JG, Tautz D, Wernegreen JJ. 2013. Animals in a bacterial world, a new imperative for the life sciences. *Proc. Natl. Acad. Sci. U.S.A.* 110:3229–3236.

The Human Microbiome Consortium. 2012. Structure, function and diversity of the healthy human microbiome. *Nature* 486: 207-214.

Young E. 2016. *I Contain Multitudes: The Microbes Within Us and a Grander View of Life*. New York: HarperCollins Publishers.

This study is based on Blumer LS, Beck CW 2020. **Introducing community ecology and data skills with the bean beetle microbiome project**. *Advances in Biology Laboratory Education* 41.

Microbial Community Analysis Using Colony-based Sequencing Database

Student Handout

Objectives

- Manipulate large datasets to conduct community-level ecological analyses
- Use community-level data to address questions about insect microbiomes
- Use Excel or RStudio programs to calculate community ecology variables
- Compare microbial community using community ecology variables

Introduction

Microbiomes are the communities of microbes (bacteria, viruses, fungi and archaea) living symbiotically with all metazoans. In the past decade, both interest and research on microbiomes, including their implications for human health, have increased dramatically (Christian *et al.* 2015, Costello *et al.* 2012, McFall-Ngai *et al.* 2013, The Human Microbiome Consortium 2012, Young 2016). Insects have been used as model species to study the importance of microbiomes, because of their ease of use and the fact that microbial communities play diverse roles in insects (Engel and Moran 2013).

The data that are collected in any microbiome study consists of lists of the taxonomic units identified and their abundance. The same types of data are evaluated in an ecological community analysis, but now the communities are the collections of microbes that constitute different microbiomes. The community variables, “species” richness and relative abundance, are the same and the statistical methods used to compare communities, diversity and difference indices, also are the same. Perhaps the simplest measure of community structure used by ecologists is “species” or taxon richness, a count of the number of unique taxa in a sample. However, species richness does not consider the relative abundance of species in a community. Imagine two communities with five different species. In one community, all of the species have the same relative abundance. In the other community, one species dominates comprising 95% of individuals in the community. The other four species are very rare. Based on species richness as a measure of community structure, these two communities are the same, although they are clearly very different. As a result, ecologists use other species diversity indices that consider both the number of species and the relative abundance of species in a community. Two common indices are the Simpson Index and the Shannon-Weaver Index. Communities with greater numbers of species and higher evenness (i.e., similar relative abundance of species within a community) are considered more diverse. Finally, measures of species richness and species diversity do not consider the identity of species in a community. So, communities could have the same level of species diversity, but have completely different species. Measure of community similarity, such as the Bray-Curtis Index, compare the similarity (or dissimilarity) between two communities based on the identity of species in the communities, as well as their relative abundances. For more information on indices of species diversity and measures of community similarity, see Krebs (1999).

In this study, bean beetle gut microbiome data were collected by undergraduate students using the protocols developed by Cole *et al.* (2018). Three types of data were collected: colony phenotypes from cultured bacteria, 16s rRNA gene sequencing of specific bacterial colonies, and whole community 16s rRNA gene sequencing, but we will limit our analyses to the colony phenotype and colony-based 16s data.

Questions

Using data from the colony-based sequencing database and the analyses described below, answer the following questions.

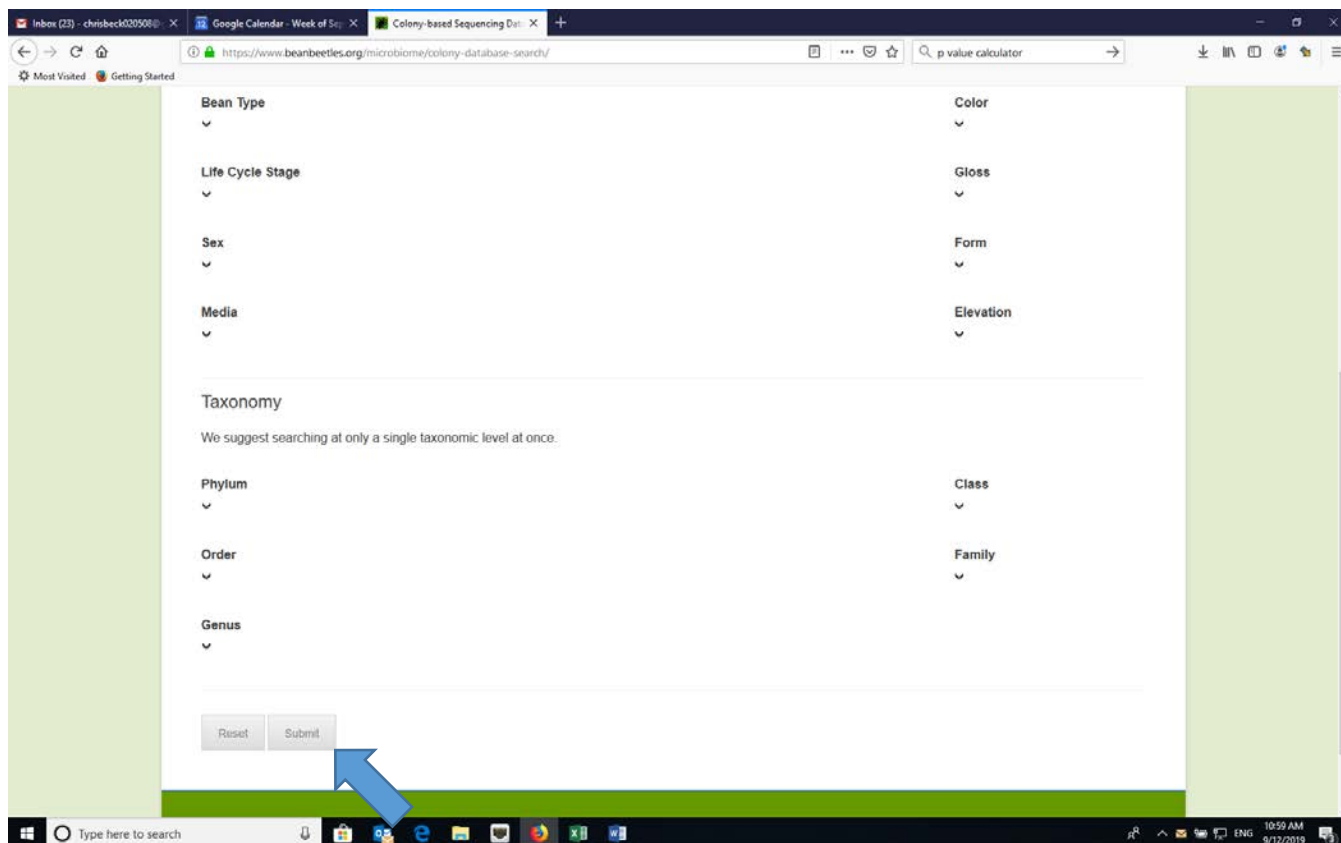
1. Which taxa are most prevalent in the bacterial communities in bean beetles?
2. Do the most prevalent taxa vary based on host bean type?
3. Based on the diversity indices that you calculated, which treatment had the highest (lowest) diversity?
4. Does the answer depend on the measure of species (taxon) diversity that you use?
5. Is there a relationship between number of samples and taxonomic diversity? If so, what might explain this?
6. Which communities are most similar (different)?
7. Do your answers to the questions above depend on the taxonomic level of analysis?

Database description

This database contains data for the microbial community of bean beetles based on 16s rRNA sequencing of individual bacterial colonies cultured from bean beetle homogenates plated on different media. Since only a small number of colonies are sequenced from each plate, the data do not represent the entire microbial community for a particular sample. However, qualitative comparisons can be made based on bean host species, sex of beetle, and other variables.

Access the database at <https://www.beanbeetles.org/microbiome/colony-database-search/>.

The database allows you to limit your search by bean host type, sex, life cycle stage, media on which bacteria were grown, colony phenotype, and bacterial taxonomy. Since we are interested in making comparisons between bacterial communities based on host species and sex, we want to download the entire database. Clicking “Submit” without limiting the search will allow you to view all of the data.



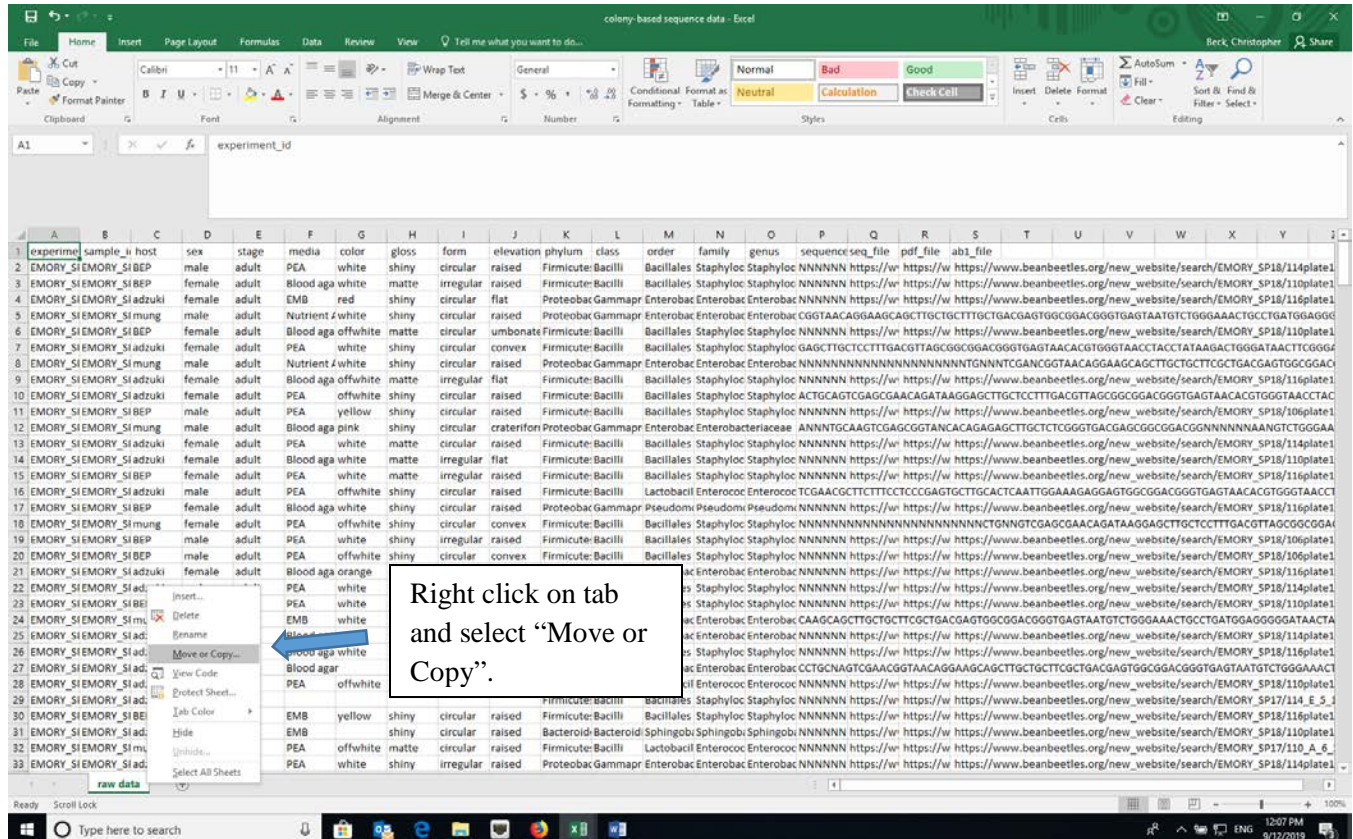
Downloading Data

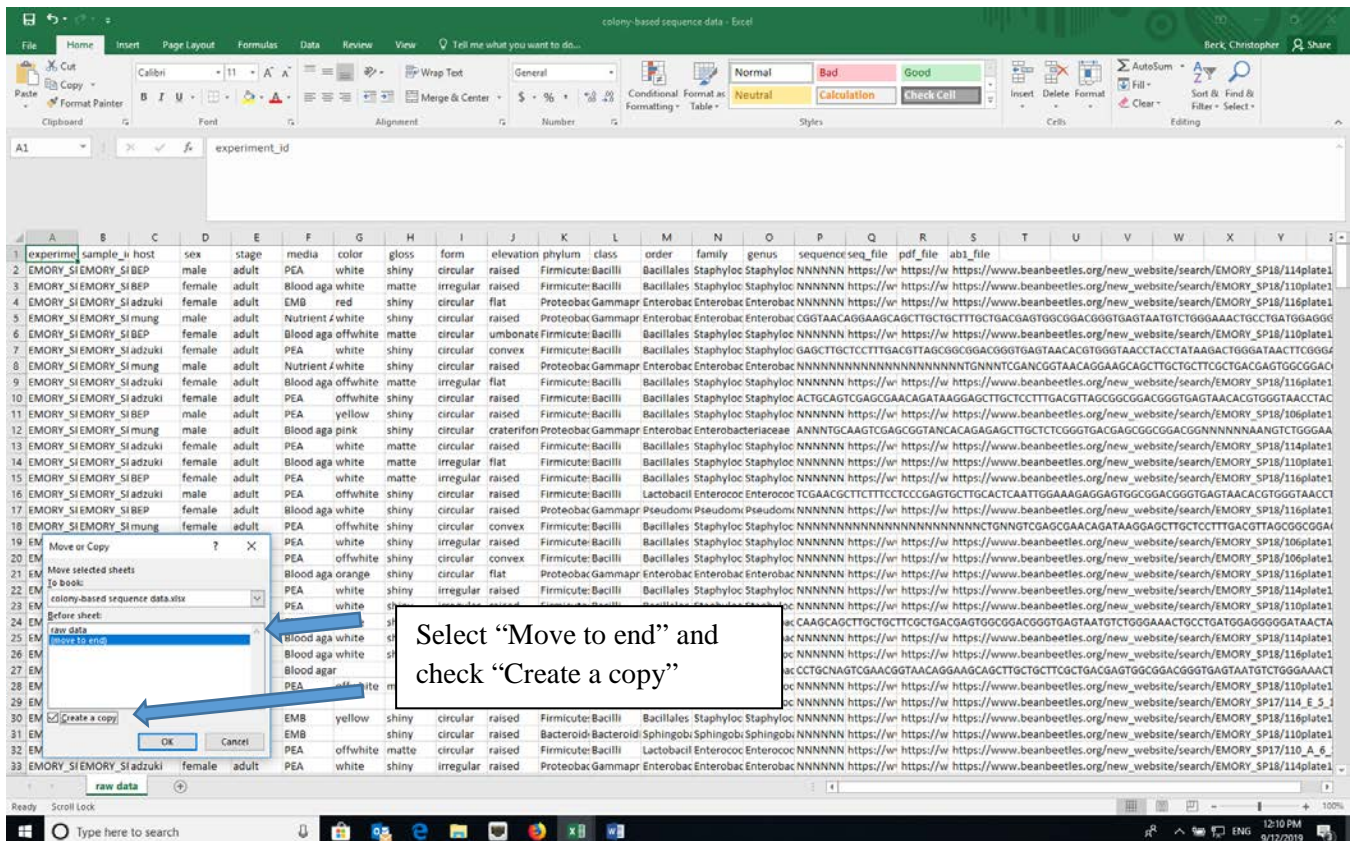
While we can view the data on the website, we want to download the data to manipulate. Click the download link to download a csv file with the data. Then, save the file as an Excel file (name the file “colony-based sequence data”) and rename the tab “raw data.”

Double click tab and rename as "raw data"

Data Reduction

1. Make a copy of the raw data in a new sheet using the sheet copy function in Excel (right click on the tab and select "Move or copy" and rename the tab ("reduced raw data").





2. In the “reduced raw data” sheet, delete any columns that we don’t need, such as the colony phenotype (color, gloss, form, elevation) and sequence data columns. The “reduced raw data” sheet is the data source if you choose to analyze these data in RStudio. Additional data manipulation and formatting (below) is required if you choose to analyze these data in Excel.

Data manipulation

1. We need to consolidate the data for each host species, each sex, or the combination of the two by the bacterial taxa. The easiest way to do this is with the Pivot Table function in Excel.
2. When clicked on a cell within the data, create a Pivot Table (Insert -> Pivot Table OR Data -> Summarize with Pivot Table). Make sure that the data source includes the top row (the cell range should include “\$A\$1”), which has the column headings. Click OK to create the pivot table in a new worksheet and label the tab “pivot table”.

colony-based sequence data - Excel

Beck, Christopher Share

The Excel ribbon is displayed with the following tabs and icons:

- File**: Save, Open, Recent, Print, etc.
- Home**: Font, Paragraph, Styles, etc.
- Insert**: Tables, Charts, Links, Text, etc.
- Page Layout**: Themes, Background Images, Page Setup, etc.
- Formulas**: Calculation Groups, Formula Auditing, etc.
- Data**: Data Tools, Data Validation, etc.
- Review**: Proofing, Changes, Comments, etc.
- View**: Views, Windows, etc.
- Help**: Tell me what you want to do...
- Store**: My Add-ins, Bing Maps, People Graph, etc.
- Recommended Charts**: Charts, Sparklines, etc.
- PivotChart**: PivotCharts, etc.
- 3D Map**: 3D Maps, etc.
- Line**: Line, Column, Win/Loss, etc.
- Slicer**: Slicers, etc.
- Timeline**: Timelines, etc.
- Hyperlink**: Hyperlinks, etc.
- Text Box**: Text Boxes, etc.
- Header & Footer**: Headers, Footers, etc.
- WordArt**: WordArt, etc.
- Signature**: Signatures, etc.
- Object**: Objects, etc.
- Equation**: Equations, etc.
- Symbol**: Symbols, etc.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z
1	experim	sample_id	host	sex	stage	media	phylum	class	order	family	genus															
2	EMORY_SI	EMORY_SI	BEP	female	adult	PEA	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
3	EMORY_SI	EMORY_SI	BEP	female	adult	Blood aga	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
4	EMORY_SI	EMORY_SI	adzuki	female	adult	EMB	Proteobac	Gammapr	Enterobac	Enterobac	Enterobacter															
5	EMORY_SI	EMORY_SI	mung	male	adult	Nutrient	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
6	EMORY_SI	EMORY_SI	BEP	female	adult	Blood aga	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
7	EMORY_SI	EMORY_SI	adzuki	female	adult	PEA	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
8	EMORY_SI	EMORY_SI	mung	male	adult	Nutrient	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
9	EMORY_SI	EMORY_SI	adzuki	female	adult	Blood aga	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
10	EMORY_SI	EMORY_SI	adzuki	female	adult	PEA	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
11	EMORY_SI	EMORY_SI	BEP	male	adult	PEA	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
12	EMORY_SI	EMORY_SI	mung	male	adult	Blood aga	Proteobac	Gammapr	Enterobac	Enterobacteriaceae																
13	EMORY_SI	EMORY_SI	adzuki	female	adult	PEA	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
14	EMORY_SI	EMORY_SI	adzuki	female	adult	Blood aga	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
15	EMORY_SI	EMORY_SI	BEP	female	adult	PEA	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
16	EMORY_SI	EMORY_SI	adzuki	male	adult	PEA	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
17	EMORY_SI	EMORY_SI	BEP	female	adult	Blood aga	Proteobac	Gammapr	Pseudomi	Pseudomi	Pseudomonas															
18	EMORY_SI	EMORY_SI	mung	female	adult	PEA	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
19	EMORY_SI	EMORY_SI	BEP	male	adult	PEA	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
20	EMORY_SI	EMORY_SI	BEP	male	adult	PEA	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
21	EMORY_SI	EMORY_SI	adzuki	female	adult	Blood aga	Proteobac	Gammapr	Enterobac	Enterobac	Enterobacter															
22	EMORY_SI	EMORY_SI	adzuki	male	adult	PEA	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
23	EMORY_SI	EMORY_SI	BEP	male	adult	PEA	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
24	EMORY_SI	EMORY_SI	mung	male	adult	EMB	Proteobac	Gammapr	Enterobac	Enterobac	Enterobacter															
25	EMORY_SI	EMORY_SI	adzuki	male	adult	Blood aga	Proteobac	Gammapr	Enterobac	Enterobac	Enterobacter															
26	EMORY_SI	EMORY_SI	adzuki	female	adult	Blood aga	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
27	EMORY_SI	EMORY_SI	adzuki	male	adult	Blood aga	Proteobac	Gammapr	Enterobac	Enterobac	Enterobacter															
28	EMORY_SI	EMORY_SI	adzuki	male	adult	PEA	Firmicute: Bacilli	Lactobacil	Enterococ	Enterococcus																
29	EMORY_SI	EMORY_SI	adzuki	female	adult	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																	
30	EMORY_SI	EMORY_SI	BEP	female	adult	EMB	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
31	EMORY_SI	EMORY_SI	adzuki	female	adult	EMB	Bacteroid: Bacteroidi	Sphingobi: Sphingobi	Sphingobacterium																	
32	EMORY_SI	EMORY_SI	mung	male	adult	PEA	Firmicute: Bacilli	Lactobacil	Enterococ	Enterococcus																
33	EMORY_SI	EMORY_SI	adzuki	female	adult	PEA	Proteobac	Gammapr	Enterobac	Enterobac	Enterobacter															

Ready Scroll Lock Type here to search

12:18 PM 9/12/2019

colony-based sequence data - Excel

File Home Insert Page Layout Formulas Data Review View Tell me what you want to do... Beck, Christopher Share

PivotTable Recommended Table PivotTables Tables Pictures Online Shapes SmartArt Screenshot My Add-ins Bing Maps People Recommended Charts Charts PivotChart 3D Map Tours Sparklines Line Column Win/Loss Slicer Timeline Hyperlink Text Box Header & Footer WordArt Signature Line Object Equation Symbol

B1 sample_id

Create PivotTable

Choose the data that you want to analyze

☒ Select a table or range

Table/Range: reduced raw data: \$A\$1:\$X\$303

☐ Use an external data source

Connection name:

☐ Use this workbook's Data Model

Choose where you want the PivotTable report to be placed

☒ New Worksheet

☐ Existing Worksheet

Location:

Choose whether you want to analyze multiple tables

☐ Add this data to the Data Model

OK Cancel

Make sure that the cell range includes \$A\$1

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z
2	EMORY_SI	EMORY_SI	BEP	male	adult	PEA	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
3	EMORY_SI	EMORY_SI	BEP	female	adult	Blood aga	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
4	EMORY_SI	EMORY_SI	adzuki	female	adult	EMB	Proteobac	Gammapr	Enterobac	Enterobac	Enterobacter															
5	EMORY_SI	EMORY_SI	mung	male	adult	Nutrient	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
6	EMORY_SI	EMORY_SI	BEP	female	adult	Blood aga	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
7	EMORY_SI	EMORY_SI	BEP	female	adult	PEA	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
8	EMORY_SI	EMORY_SI	mung	male	adult	Nutrient	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
9	EMORY_SI	EMORY_SI	adzuki	female	adult	Blood aga	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
10	EMORY_SI	EMORY_SI	adzuki	female	adult	PEA	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
11	EMORY_SI	EMORY_SI	BEP	male	adult	PEA	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
12	EMORY_SI	EMORY_SI	mung	male	adult	Blood aga	Proteobac	Gammapr	Enterobac	Enterobacteriaceae																
13	EMORY_SI	EMORY_SI	adzuki	female	adult	PEA	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
14	EMORY_SI	EMORY_SI	adzuki	female	adult	Blood aga	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
15	EMORY_SI	EMORY_SI	BEP	female	adult	PEA	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
16	EMORY_SI	EMORY_SI	adzuki	male	adult	PEA	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
17	EMORY_SI	EMORY_SI	BEP	female	adult	Blood aga	Proteobac	Gammapr	Pseudomi	Pseudomi	Pseudomonas															
18	EMORY_SI	EMORY_SI	mung	female	adult	PEA	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
19	EMORY_SI	EMORY_SI	BEP	male	adult	PEA	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
20	EMORY_SI	EMORY_SI	BEP	male	adult	PEA	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
21	EMORY_SI	EMORY_SI	adzuki	female	adult	Blood aga	Proteobac	Gammapr	Enterobac	Enterobac	Enterobacter															
22	EMORY_SI	EMORY_SI	adzuki	male	adult	PEA	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
23	EMORY_SI	EMORY_SI	BEP	male	adult	PEA	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
24	EMORY_SI	EMORY_SI	mung	male	adult	EMB	Proteobac	Gammapr	Enterobac	Enterobac	Enterobacter															
25	EMORY_SI	EMORY_SI	adzuki	male	adult	Blood aga	Proteobac	Gammapr	Enterobac	Enterobac	Enterobacter															
26	EMORY_SI	EMORY_SI	adzuki	female	adult	Blood aga	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
27	EMORY_SI	EMORY_SI	adzuki	male	adult	Blood aga	Proteobac	Gammapr	Enterobac	Enterobac	Enterobacter															
28	EMORY_SI	EMORY_SI	adzuki	male	adult	PEA	Firmicute: Bacilli	Lactobacil	Enterococ	Enterococcus																
29	EMORY_SI	EMORY_SI	adzuki	female	adult	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																	
30	EMORY_SI	EMORY_SI	BEP	female	adult	EMB	Firmicute: Bacilli	Bacillales	Staphyloc	Staphylococcus																
31	EMORY_SI	EMORY_SI	adzuki	female	adult	EMB	Bacteroid: Bacteroidi	Sphingobi: Sphingobi	Sphingobacterium																	
32	EMORY_SI	EMORY_SI	mung	male	adult	PEA	Firmicute: Bacilli	Lactobacil	Enterococ	Enterococcus																
33	EMORY_SI	EMORY_SI	adzuki	female	adult	PEA	Proteobac	Gammapr	Enterobac	Enterobac	Enterobacter															
34	EMORY_SI	EMORY_SI	adzuki	female	adult	PEA	Proteobac	Gammapr	Enterobac	Enterobacteriaceae																

raw data reduced raw data

Enter Scroll Lock Type here to search 12:20 PM 9/12/2019

- Set the treatment(s) that you are interested (for example, host species) in as the rows and the bacterial taxonomic level you are interested in as the columns. The Values should be a COUNT of the sample_id, as each row in the dataset represents a single sample.

Count of sample_id

Count of sample_id	Column Labels														
Row Labels	Acinetobacter	Bacillus	Burkholderia	Caballeronia	Paraburkholderia	Corynebacterium	1	Enterobacter	Enterococcus	Escherichia-Shigella	Klebsiella	Paenibacillus	Pseudomonas	Ralstonia	Sphingomonas
adzuki	1	1				3	55	20			1		4	1	
BEP					1	1	27	1		1				4	
mung							56	10						4	
pigeon						1	9				1	1		5	
Grand Total	1	1			1	7	147	31		1	2	1	17	1	

Drag treatment of interest (e.g., host) to ROWS, taxonomic level (e.g., genus) to COLUMNS, and sample_id to VALUES

- You can add zeros to all of the empty cells in the Pivot Table using the Options menu and remove the Grand totals for Columns using the Options menu or the Design tab depending of your version of Excel. (You want to keep the Grand totals for rows to calculate diversity indices.)

PivotTable Name: PivotTable1

Count of sample_id

Row Labels: Acinetobacter

Column Labels: Bacillus, Burkholderia-Caballeronia-Paraburkholderia, Corynebacterium 1

Grand Total: 1, 1, 1, 7

PivotTable Options

Layout & Format

Layout

☐ Merge and center cells with labels

When in compact form indent row labels: 1 (character(s))

Display fields in report filter area: Down, Then Over

Report filter fields per column: 0

Format

☐ For error values show: 0

☐ For empty cells show: 0

☒ Autofill column widths on update

☒ Preserve cell formatting on update

OK Cancel

PivotTable Name: PivotTable1

Count of sample_id

Row Labels: Acinetobacter

Column Labels: Bacillus, Burkholderia-Caballeronia-Paraburkholderia, Corynebacterium 1

Grand Total: 1, 1, 1, 7

PivotTable Options

Layout & Format

Grand Totals

☒ Show grand totals for rows

☐ Show grand totals for columns

Filters

☐ Subtotal filtered page

☐ Allow multiple filters per field

Sorting

☒ Use Custom Lists when sorting

OK Cancel

Unselect checkbox for "Show grand totals for columns"

- You can remove the “blanks” column using the Column labels dropdown (located at upper left of the sheet) and unselecting “blank”.

The screenshot shows an Excel spreadsheet with a PivotTable. The PivotTable is named 'Count of sample_id' and is located in the range A3:L33. The PivotTable Fields task pane is open on the right, showing the 'Columns' section with 'genus' selected. The 'Rows' section is empty. The 'Values' section shows 'Count of sam...'. The 'Filters' section is empty. The PivotTable data is as follows:

Count of sample_id	Burkholderia	Caballeronia	Paraburkholderia	Corynebacterium	1 Enterobacter	Enterococcus	Escherichia-Shigella	Klebsiella	Paenibacillus	Pseudomonas	Ralstonia	Sph...
0	3	55	20	1	1	0	1	0	0	4	1	
1	1	27	1	1	1	0	1	0	0	4	0	
0	0	56	10	0	1	1	1	4	0			
0	1	9	0	0	0	0	0	5	0			

The PivotTable Fields task pane shows the following fields:

- Columns: genus
- Rows: host
- Values: Count of sam...
- Filters: (empty)

The 'Column Labels' dropdown in the PivotTable is set to 'genus'. The '(blank)' checkbox in the task pane is unchecked. A text box with the text 'Unselect checkbox for "blanks"' is overlaid on the task pane.

- If you selected more than one treatment for the rows, you can get the treatment data to repeat for each sample. In the Design tab, select “Report Layout” and choose “Show in Tabular form” and “Repeat All Item Labels”.

The screenshot shows an Excel PivotTable with the following data:

host	Acinetobacter	Bacillus	Burkholderia-Caballeronia-Paraburkholderia	Corynebacterium 1	Enterobacter
adzuki	1	0	0	0	5
BEP	0	0	1	1	1
mung	0	0	0	0	0
pigeon	0	0	0	0	1

The PivotTable Fields task pane on the right shows the following configuration:

- Choose fields to add to report:** experiment_id, sample_id, host, sex, stage, media, phylum, class, order, family, genus.
- Drag Fields between areas below:**
 - FILTERS:** (empty)
 - COLUMNS:** genus
 - ROWS:** host
 - VALUES:** Count of sam...

- Copy and paste (as values) the pivot table to a new worksheet and remove any extra rows at the top. The top row should now have the taxa names. Name this tab "analysis". Conduct the community ecology analyses that follow in Excel on the "analysis" sheet.

colony-based sequence data - Excel

Beck, Christopher

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z
1	host	Acinetob	Bacillus	Burkholder	Coryneb	Enterobac	Enterococ	Escherich	Klebsiella	Paenibaci	Pseudom	Kalstonia	Sphingobi	Staphyloc	Stenotrop	Grand Total										
2	adzuki	1	1	0	5	55	20	0	1	0	4	1	1	64	1	154										
3	BEP	0	0	1	1	27	1	1	0	0	4	0	0	50	2	87										
4	mung	0	0	0	0	56	10	0	1	1	4	0	0	22	0	94										
5	pigeon	0	0	0	1	9	0	0	0	0	5	0	0	7	0	22										
6																										
7																										
8																										
9																										
10																										
11																										
12																										
13																										
14																										
15																										
16																										
17																										
18																										
19																										
20																										
21																										
22																										
23																										
24																										
25																										
26																										
27																										
28																										
29																										
30																										
31																										
32																										
33																										

raw data | Sheet2 | analysis | reduced raw data

Select destination and press ENTER or choose Paste

Type here to search

2:00 PM

9/12/2019

Calculating diversity indices

1. Species (taxon) richness – the number of unique species (taxa) in a sample
 - a. Although you could manually count the number of cells with values greater than zero for each treatment, using the COUNTIF formula in Excel is easier (e.g., =COUNTIF(range,">0")). Where "range" is the cell range in the datasheet, for example "C2:M2", a single row or treatment.

host	Acinetob	Bacillus	Burkholder	Coryneb	Enterobac	Enterococ	Escherich	Klebsiella	Paenibac	Pseudom	Raistonia	Sphingob	Staphyloc	Stenotroph	Grand Tot.	Richness
adzuki	1	1	0	5	55	20	0	1	0	4	1	1	64	1	154	=COUNTIF(B2:O2,>0)
BEP	0	0	1	1	27	1	1	0	0	4	0	0	50	2	87	COUNTIF(range, criteria)
mung	0	0	0	0	56	10	0	1	1	4	0	0	22	0	94	
pigeon	0	0	0	1	9	0	0	0	0	5	0	0	7	0	22	

2. Simpson Index – the Simpson Index incorporates both species (taxon) richness and species (taxon) evenness.

- $D = \sum (n/N)^2$, where n =number of individuals of a particular species (taxon) and N =total number of individuals in a sample. D increases as diversity decreases, which is counterintuitive. A reciprocal or inverse index would be more intuitive and are easily calculated.
- Reciprocal Simpson = $1/D$ and scales so the maximum value is the species richness of a community.
- Inverse Simpson = $1-D$ and scales to a maximum value of 1.0.
- Create a new data array below the original using the same row labels (treatment variables) and the same column labels (bacterial taxa).

colony-based sequence data - Excel

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z
1	host	Acinetoba	Bacillus	Burkholde	Coryneb	Enterobac	Enterococ	Escherich	Klebsiella	Paenibaci	Pseudom	Ralstonia	Sphingobi	Staphyloc	Stenotroph	Grand Tot	Richness									
2	adzuki	1	1	0	5	55	20	0	1	0	4	1	1	64	1	154	11									
3	BEP	0	0	1	1	27	1	1	0	0	4	0	0	50	2	87										
4	mung	0	0	0	0	56	10	0	1	1	4	0	0	22	0	94										
5	pigeon	0	0	0	1	9	0	0	0	0	5	0	0	7	0	22										
6																										
7	host	Acinetoba	Bacillus	Burkholde	Coryneb	Enterobac	Enterococ	Escherich	Klebsiella	Paenibaci	Pseudom	Ralstonia	Sphingobi	Staphyloc	Stenotroph	Grand Tot	Richness									
8	adzuki																									
9	BEP																									
10	mung																									
11	pigeon																									
12																										
13																										
14																										
15																										
16																										
17																										
18																										
19																										
20																										
21																										
22																										
23																										
24																										
25																										
26																										
27																										
28																										
29																										
30																										
31																										
32																										
33																										

raw data | Sheet2 | analysis | reduced raw data

Ready | Scroll Lock | Type here to search | 2:04 PM | 9/12/2019

- e. To calculate the proportion squared for each taxa, use the grand totals for each treatment. Using the Excel trick that \$ before a column or row prevents Excel from iterating when copying a formula makes this easy. For example, $= (C2/\$P2)^2$. Copy the formula across the row and then down.

colony-based sequence data - Excel

File Home Insert Page Layout Formulas Data Review View Tell me what you want to do...

Clipboard Font Alignment Number Styles Editing

COUNTIF $= (B2/SP2)^2$

	host	Acinetobacte	Bacillus	Burkholder	Corynebaci	Enterobac	Enterococ	Escherichi	Klebsiella	Paenibaci	Pseudom	Ralstonia	Sphingobi	Staphyloc	Stenotroph	Grand Tot	Richness
1	adzuki	1	1	0	5	55	20	0	1	0	4	1	1	64	1	154	11
2	BEP	0	0	1	1	27	1	1	0	0	4	0	0	30	2	87	8
3	mung	0	0	0	0	56	10	0	1	1	4	0	0	22	0	94	6
4	pigeon	0	0	0	1	9	0	0	0	0	5	0	0	7	0	22	4

8 host Acinetobacte Bacillus Burkholder Corynebaci Enterobac Enterococ Escherichi Klebsiella Paenibaci Pseudom Ralstonia Sphingobi Staphyloc Stenotrophomonas

9 adzuki $= (B2/SP2)^2$ 4.22E-05 0.001054 0.127551 0.016866 0 4.22E-05 0 0.000675 4.22E-05 4.22E-05 0.17271 4.22E-05

10 BEP

11 mung

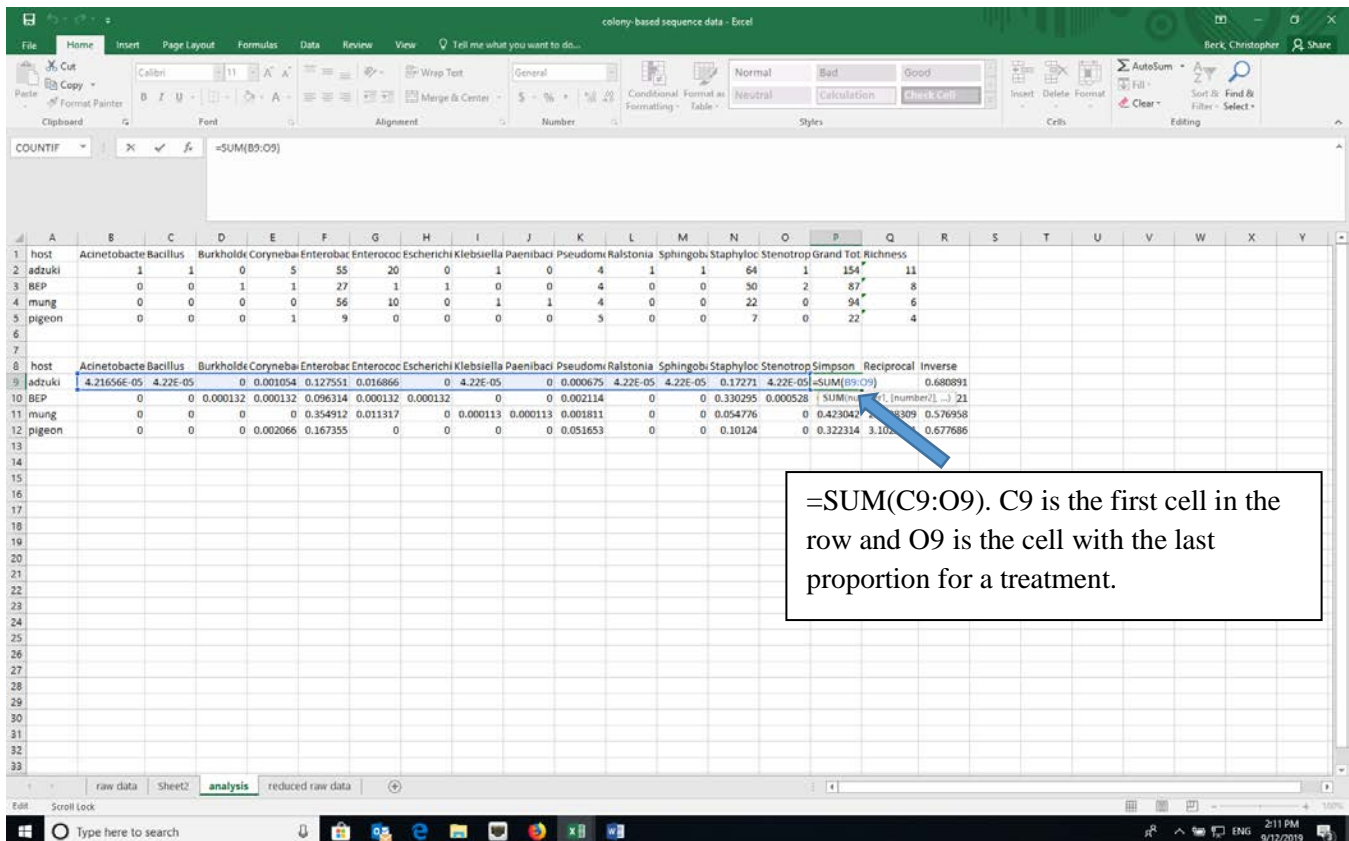
12 pigeon

$= (B2/SP2)^2$. B2 is the cell with the abundance of the first taxon and P2 is the cell with the grand total for a treatment. The \$ prevents the column identifier from changing. Copy the formula across the row and then down.

raw data Sheet2 analysis reduced raw data

Type here to search 2:07 PM 9/12/2019

- f. Calculate the sum of the proportions squared ($=SUM$ in Excel) to calculate the Simpson Index.



g. Calculate the reciprocal (e.g., $=1/P9$) and inverse Simpson (e.g., $=1-P9$) using formulas in Excel.

3. Shannon-Weaver (Shannon-Weiner) Index – also incorporates species (taxon) richness and species (taxon) evenness

- $H = -\sum p \ln p$, where p is the proportion of individuals of each bacterial taxon in a community (i.e., n/N).
- Create a new data array below the original using the same row labels (treatment variables) and the same column labels (species).
- Using the grand totals for each treatment, calculate the proportions ($p \ln p$). Using the Excel trick that \$ before a column or row prevents Excel from iterating when copying a formula makes this easy.
- Note that $\ln p$ is undefined if $p=0$, so you can use an "IF" statement in Excel. For example, $=IF(B2>0,(B2/$P2)*LN((B2/$P2)),"")$

colony-based sequence data v2 - Excel

File Home Insert Page Layout Formulas Data Review View Tell me what you want to do...

Clipboard Font Alignment Number Styles Editing

COUNTIF =IF(B2>0,(B2/\$P2)*LN((B2/\$P2)), "")

host	Acinetobacte	Bacillus	Burkholder	Coryneb	Enterobac	Enterococ	Escherich	Klebsiella	Paenibaci	Pseudom	Ralstonia	Sphingobi	Staphyloc	Stenotro	Simpson	Grand Tot	Richness
adzuki	1	1	0	5	55	20	0	1	0	4	1	1	64	1	154	11	
BEP	0	0	1	1	27	1	1	0	0	4	0	0	30	2	87	8	
mung	0	0	0	0	56	10	0	1	1	4	0	0	22	0	94	6	
pigeon	0	0	0	1	9	0	0	0	0	5	0	0	7	0	22	4	

host Acinetobacte Bacillus Burkholder Coryneb Enterobac Enterococ Escherich Klebsiella Paenibaci Pseudom Ralstonia Sphingobi Staphyloc Stenotro Simpson Reciprocal Inverse

adzuki 4.21656E-05 4.22E-05 0 0.001054 0.127551 0.016866 0 4.22E-05 0 0.000675 4.22E-05 0.17271 4.22E-05 0.319109 3.1337209 0.680891

BEP 0 0 0.000132 0.000132 0.096314 0.000132 0.000132 0 0 0.002114 0 0 0.330295 0.000528 0.429779 2.3267753 0.570221

mung 0 0 0 0 0.354912 0.011317 0 0.000113 0.001811 0 0 0.054776 0 0.423042 2.3638309 0.576958

pigeon 0 0 0 0.002066 0.167355 0 0 0 0 0.051653 0 0 0.10124 0 0.322314 3.1025641 0.677686

host Acinetobacte Bacillus Burkholder Coryneb Enterobac Enterococ Escherich Klebsiella Paenibaci Pseudom Ralstonia Sphingobi Staphyloc Stenotro Shannon

adzuki =IF(B2>0,(B2/\$P2)*LN((B2/\$P2)), "") -0.36772 -0.26509 -0.03271 -0.09482 -0.03271 -0.03271 -0.36491 -0.03271 1.400077

BEP IF(logical_test, [value_if_true], [value_if_false]) 13 -0.36313 -0.05133 -0.05133 -0.14159 -0.31832 -0.08673 1.115101

mung -0.30856 -0.23837 -0.04833 -0.13434 -0.33989 1.11783

pigeon -0.1405 -0.36565 -0.04833 -0.13434 -0.33989 1.207243

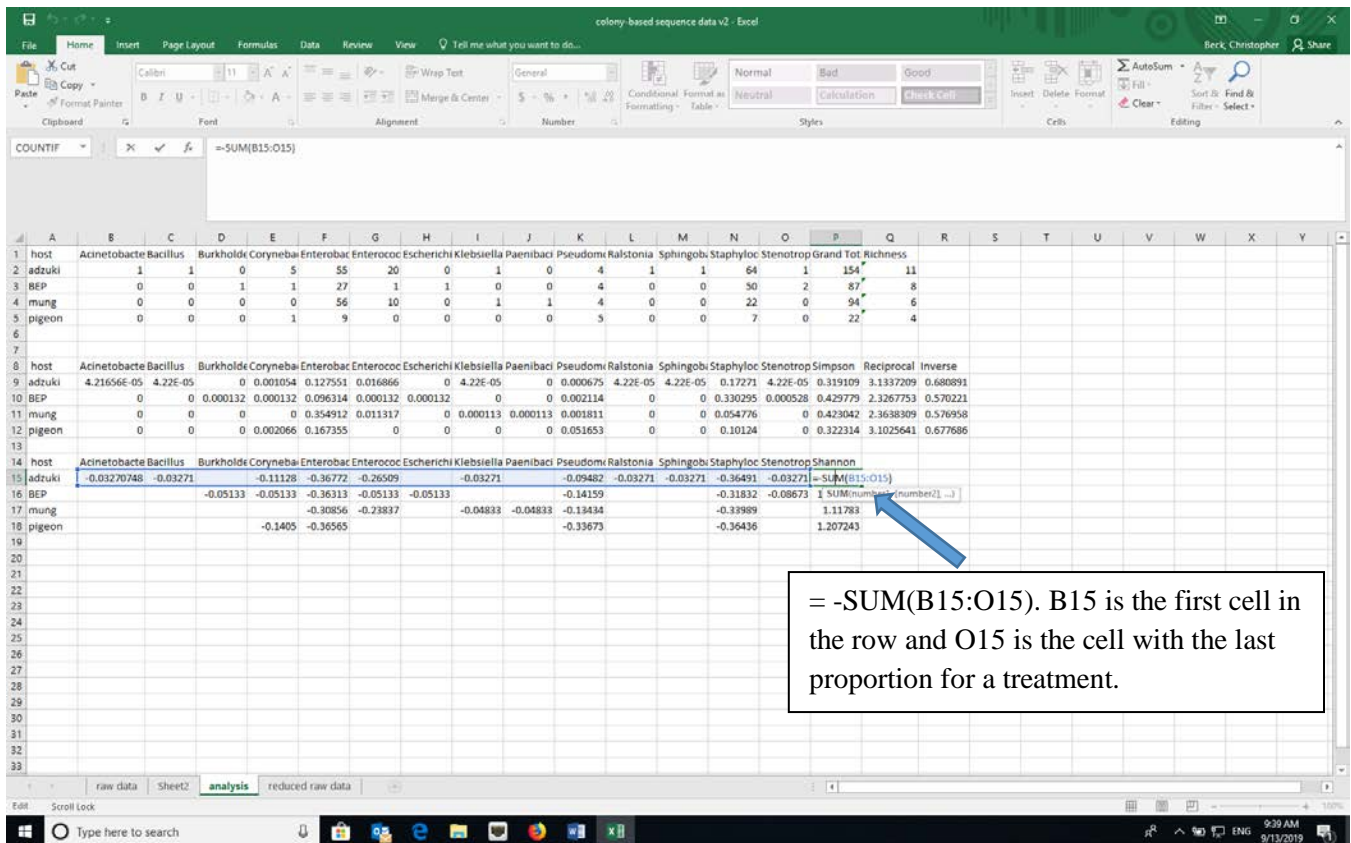
raw data Sheet2 analysis

Type here to search

9:32 AM 9/13/2019

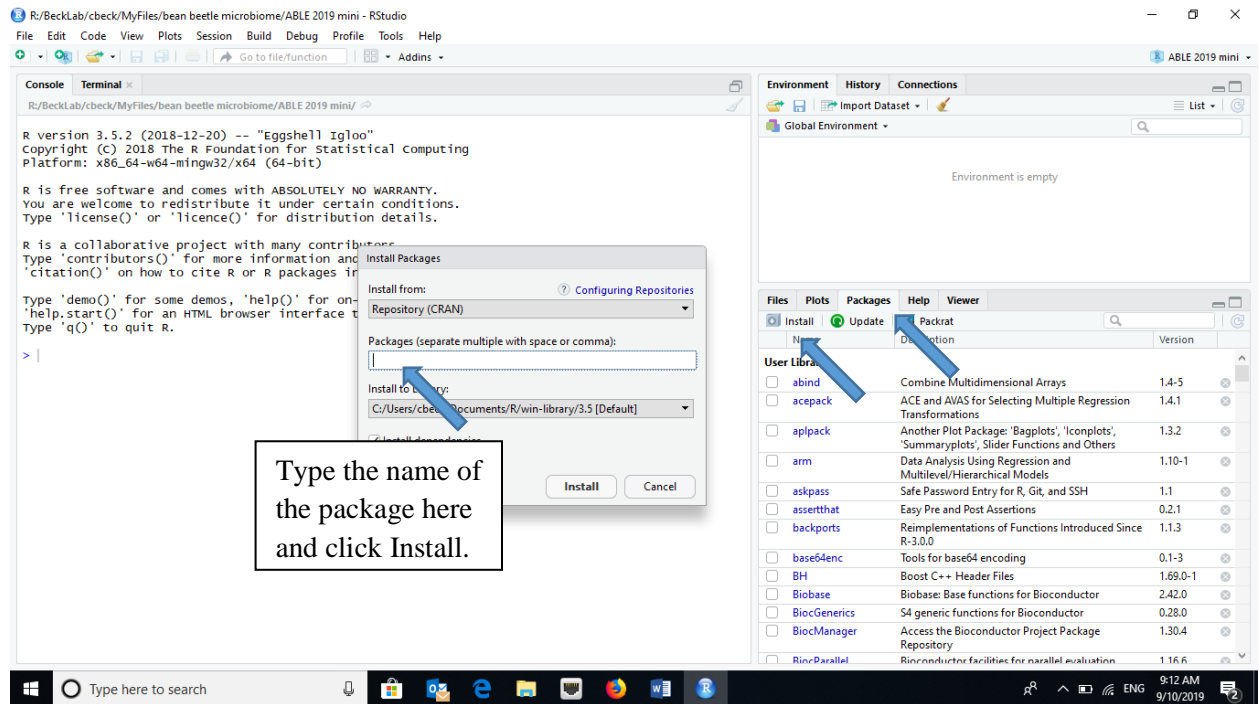
=IF(B2>0,(B2/\$P2)*LN((B2/\$P2)), ""). B2 is the cell with the abundance of the first taxon and P2 is the cell with the grand total for a treatment. The \$ prevents the column identifier from changing. Copy the formula across the row and then down.

- e. Calculate the negative sum of the proportions ($p \ln p$) (=SUM in Excel for each row, a different microbial community) to calculate the Shannon-Weaver Index.

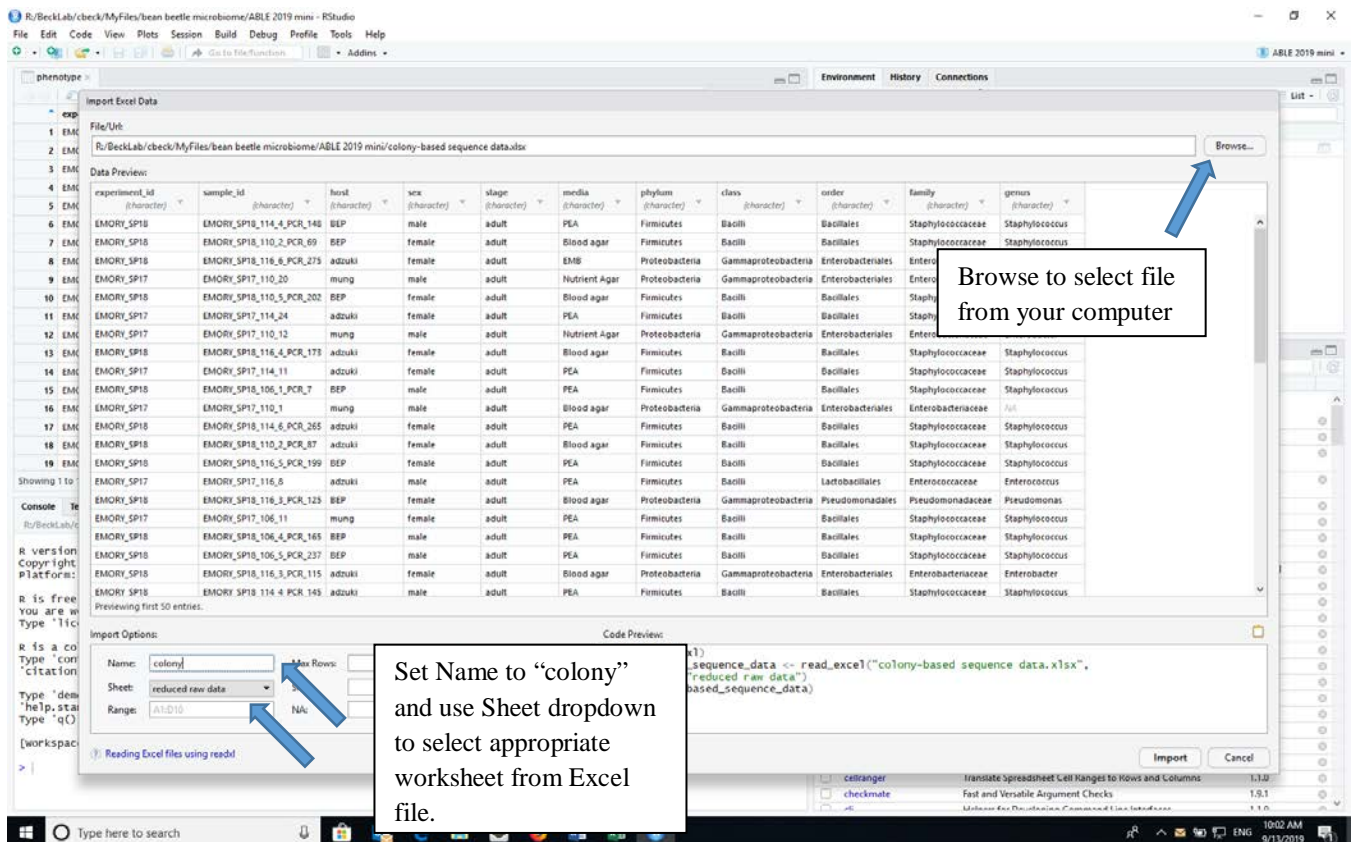
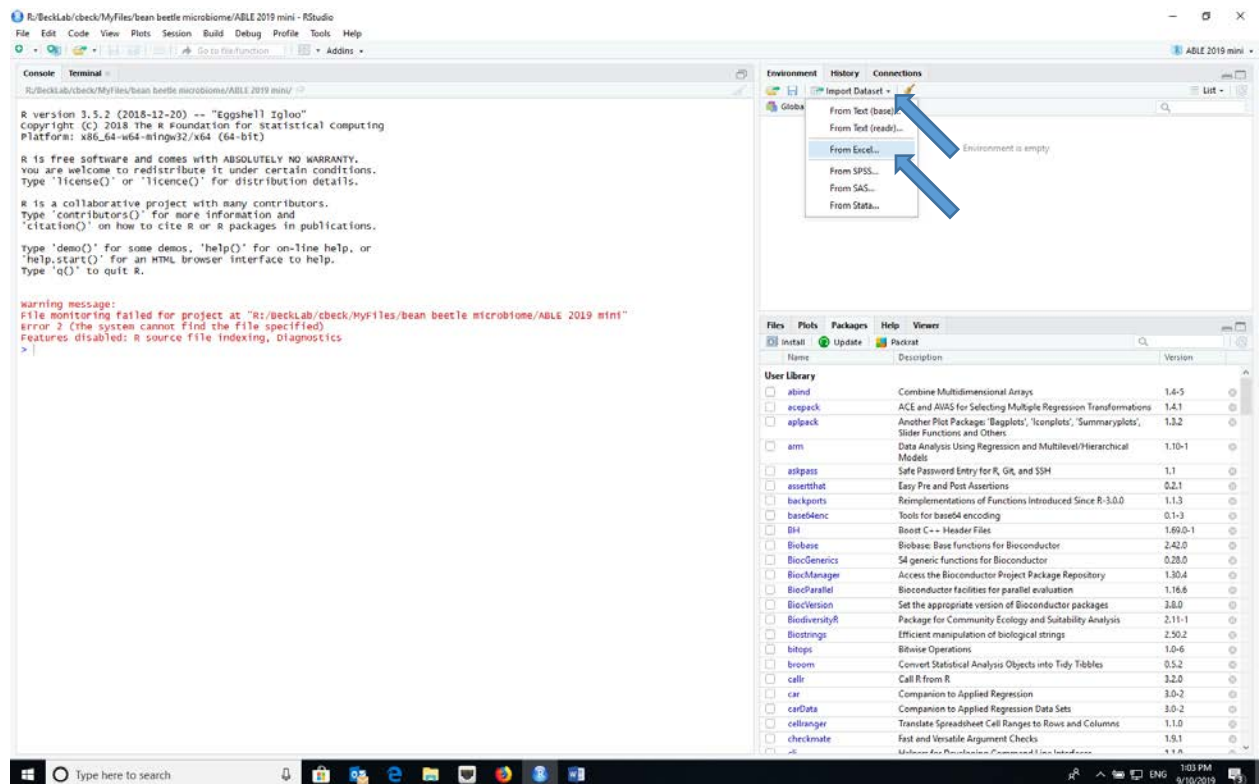


Data Manipulation in R

1. Open RStudio and create a new project using the New Project option under File and select for the new project to be in an existing folder where your data are.
2. Install the following packages either using the Packages tab in RStudio or the command `install.packages("name_of_package")` in the console. Note that BiodiversityR requires QuartzX on a Mac. If you are using a MacOS and don't have QuartzX, install it first and restart your computer before install these packages.
 - a. dplyr
 - b. reshape2
 - c. vegan
 - d. BiodiversityR
 - e. ggplot2



3. Load the packages listed above by clicking the checkboxes for the appropriate packages in the Packages tab or the command `library("name_of_package")` in the console.
4. Import the dataset "reduced raw data" (dataset without the extra metadata that you created in the Excel section) into RStudio.



5. Attach the imported dataset ("colony") to the dataframe using the attach command in the console (`attach(colony)`)
6. Create a community matrix (named "community" in this example) for a particular treatment. This example assumes that you are doing the analysis at the genus level. This can be changed to other taxonomic levels using the appropriate variable name

```
>community<-table(host,genus)
```

7. If you want to look at two factors at the same time, creating the community matrix is a little more complicated. The first command calculates the count of each genus by each sex and host combination and drops any missing values. The second command creates a community matrix.

```
> community_2 <- colony %>% count(sex,host,genus) %>% drop_na()
```

```
> comm2<-dcast(community_2, sex+host~genus, value.var = "n", fun.aggregate = sum)
```

"genus" in both command lines may be whatever taxon level you wish to evaluate in the dataset. For example, it could be changed to "family" or "order".

Calculating diversity indices

Note: "community" is the name of the community matrix

1. Species Richness

```
> diversityresult(community,index="richness",method="each site")
```

2. Simpson

```
> diversityresult(community,index="Simpson",method="each site")
```

This calculates the inverse Simpson described above

```
> diversityresult(community,index="inverseSimpson",method="each site")
```

This calculates the reciprocal Simpson described above. (confusing that it is called in the inverseSimpson)

3. Shannon

```
> diversityresult(community,index="Shannon",method="each site")
```

Calculating community similarity (distance)

Sometimes we are interested in how similar (or different) two communities are based on what species (taxa) are present and the relative abundance of those species (taxa) in the two communities. One of the most common measures of distance is the Bray Curtis Dissimilarity. Similarity can be measured as $1-BC$.

$$BC_{ij} = 1 - \frac{2C_{ij}}{S_i + S_j}$$

Where:

- i & j are the two samples,
- S_i is the total number of specimens counted in sample i ,
- S_j is the total number of specimens counted in sample j ,
- C_{ij} is the sum of only the lesser counts for each taxa found in both sites.

Although Bray-Curtis Dissimilarity is often used in community ecology, it is not robust to incomplete sampling of the community (all taxa are not sampled) or unbalanced sampling (all treatments are not equally sampled). An alternative is the Morista-Horn Index of Dissimilarity ($1-C_H$). Morista-Horn Index of Similarity is

$$C_H = \frac{2 \sum_{i=1}^{D_{12}} \frac{X_i}{n} \frac{Y_i}{m}}{\sum_{i=1}^{D_1} \left(\frac{X_i}{n}\right)^2 + \sum_{i=1}^{D_2} \left(\frac{Y_i}{m}\right)^2}$$

Where

- D_1 =number of taxa in sample 1
- D_2 =number of taxa in sample 2
- D_{12} =number of taxa in shared in both communities
- X_i =number of individuals of taxon i in sample 1
- Y_i =number of individuals of taxon i in sample 2
- n =total number of individuals in sample 1
- m =total number of individuals in sample 2

So that X_i/n and Y_i/m are proportion of individuals of taxon i in each of the samples (communities).

To produce a matrix of all of the pair-wise distances between samples using the Bray Curtis index of distance, use the following command.

```
> vegdist(community, method="bray", binary=FALSE, diag=FALSE, upper=FALSE)
```

To produce a matrix of all of the pair-wise distances between samples using the Morista-Horn index of distance.

```
> vegdist(community, method="horn", binary=FALSE, diag=FALSE, upper=FALSE)
```

Cited References

- Christian N, Whitaker BK, Clay K. 2015. Microbiomes: unifying animal and plant systems through the lens of community ecology theory. *Front. Microbiol.* 6:1–15.
- Cole MF, Acevedo-Gonzalez T, Gerardo NM, Harris EV, Beck CW. 2018. Effect of diet on bean beetle microbial communities. Article 3 In: McMahon K, editor. *Tested studies for laboratory teaching*. Volume 39. Proceedings of the 39th Conference of the Association for Biology Laboratory Education (ABLE).
- Costello EK, Stagaman K, Dethlefsen L, Bohannan BJM, Relman DA. 2012. The application of ecological theory toward an understanding of the human microbiome. *Science*. 336:1255–1262
- Engel P, Moran NA. 2013. The gut microbiota of insects – diversity in structure and function. *FEMS Microbiol. Rev.* 37:699-735.
- Krebs CJ. 1999. *Ecological Methodology*, 2nd edition. New York: Benjamin Cummings.
- McFall-Ngai M, Hadfield MG, Bosch TCG, Carey HV, Domazet-Lošo T, Douglas AE, Dubilier N, Eberl G, Fukami T, Gilbert SF, Hentschel U, King N, Kjelleberg S, Knoll AH, Kremer N, Mazmanian SK, Metcalf JL, Neelson K, Pierce NE, Rawls JF, Reid A, Ruby EG, Rumpho M, Sanders JG, Tautz D, Wernegreen JJ. 2013. Animals in a bacterial world, a new imperative for the life sciences. *Proc. Natl. Acad. Sci. U.S.A.* 110:3229–3236.
- The Human Microbiome Consortium. 2012. Structure, function and diversity of the healthy human microbiome. *Nature* 486: 207-214.
- Young E. 2016. *I Contain Multitudes: The Microbes Within Us and a Grander View of Life*. New York: HarperCollins Publishers.

This study is based on Blumer LS, Beck CW 2020. **Introducing community ecology and data skills with the bean beetle microbiome project**. *Advances in Biology Laboratory Education* 41.