

PathoStat User Manual

Sandro L. Valenzuela, Eduardo Castro-Nallar, Solaiappan Manimaran

2017-12-07

Package

PathoStat 1.4.1

Contents

1	Introduction	3
2	Installation and Run	3
3	Relative Abundance	3
3.1	Taxonomy level	3
3.2	Heatmap	4
3.3	Summary	5
3.4	RA Table (%)	5
3.5	Count table	6
4	Diversity	6
4.1	Alpha Diversity	6
4.2	Beta Diversity	7
4.3	Exploratory Tree	7
4.4	Biplot	8
4.5	Co-Occurrence	8
5	Differential Expression	9
5.1	Expression Plots	9
5.2	Summary, Table and LIMMA	10
6	Confidence Region	11
7	PCA and PCoA	12
7.1	Explained Variation	12

8 Time Series 13

9 Core OTUs 14

1 Introduction

Welcome! This is the manual for the PathoStat package. PathoStat is a Shiny App interactive package that will let you explore metagenomic datasets, e.g., microbiome abundance tables, for exploratory data analysis, differential abundance hypothesis testing, and more. PathoStat can take metagenomic abundance data produced by any taxonomic profiling pipeline, however, to get the most out of PathoStat we recommend using PathoScope2. PathoStat is not limited to whole-metagenomic shotgun data but also can take metataxonomic data, i.e., 16S rRNA, ITS, etc. Plots can be exported in vector-based file formats (svg, PDF) for sharing or further editing. For detailed installation instructions check the introductory and advanced vignettes.

```
require(PathoStat)
vignette("PathoStatIntro")
vignette("PathoStatUserManual")
vignette("PathoStatAdvanced")
```

2 Installation and Run

While there are a few functions that can work from the console, PathoStat is designed to be run interactively. From the R console simply type:

```
source("http://bioconductor.org/biocLite.R")
biocLite("PathoStat")
```

If all went well you should now be able to load PathoStat:

```
require(PathoStat)
runPathoStat()
```

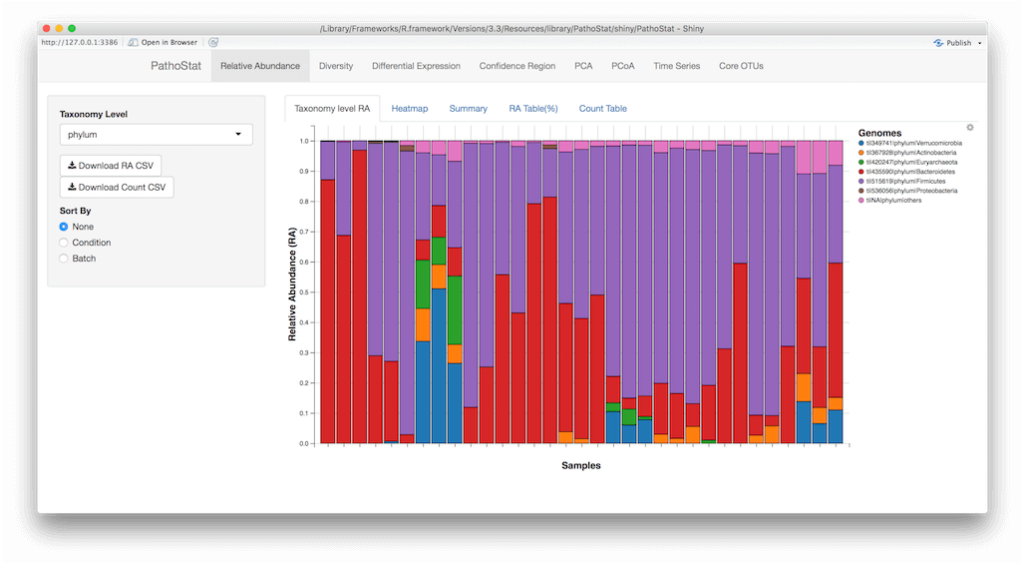
PathoStat functionality is grouped into eight tabs that comprise all analyses and visualizations. Here, we will go over each of those eight tabs and their main features.

3 Relative Abundance

The relative abundance tab is further subdivided into five subtabs where you can explore your data in the form of a stacked bar chart, a heatmap, summary, and searchable relative abundance and count tables. You can download these data both in read counts or in relative proportions.

3.1 Taxonomy level

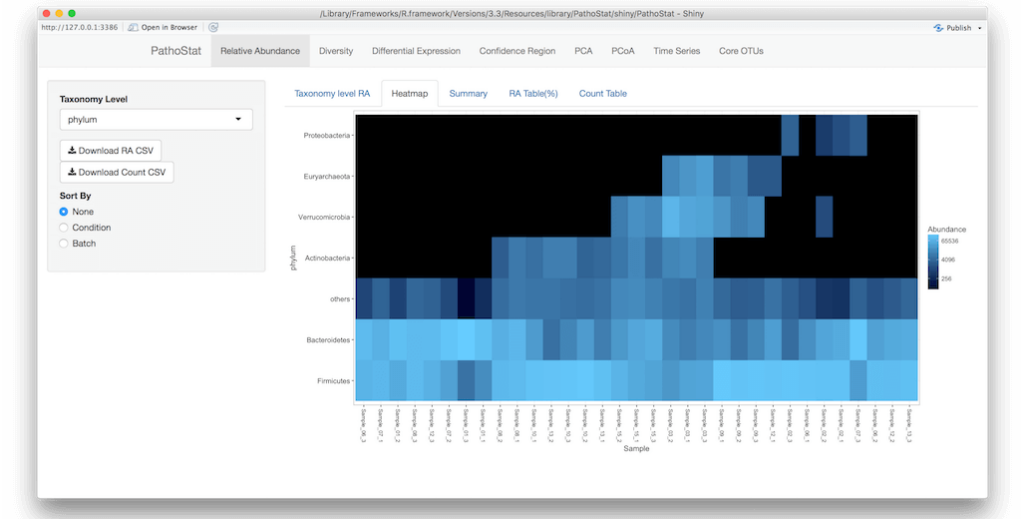
Within the Taxonomy level RA subtab, you can hover over different stacked bars and get the taxonomy membership and relative abundance. In addition you can sort the stacked bars (samples) by different factors in your data and download a svg image for further editing.



The image above shows the user interface of the stacked bar chart. On the left you can select the appropriate taxonomy level to plot, and by clicking on the gear icon on the upper right side of the chart, you can select the file format for your download.

3.2 Heatmap

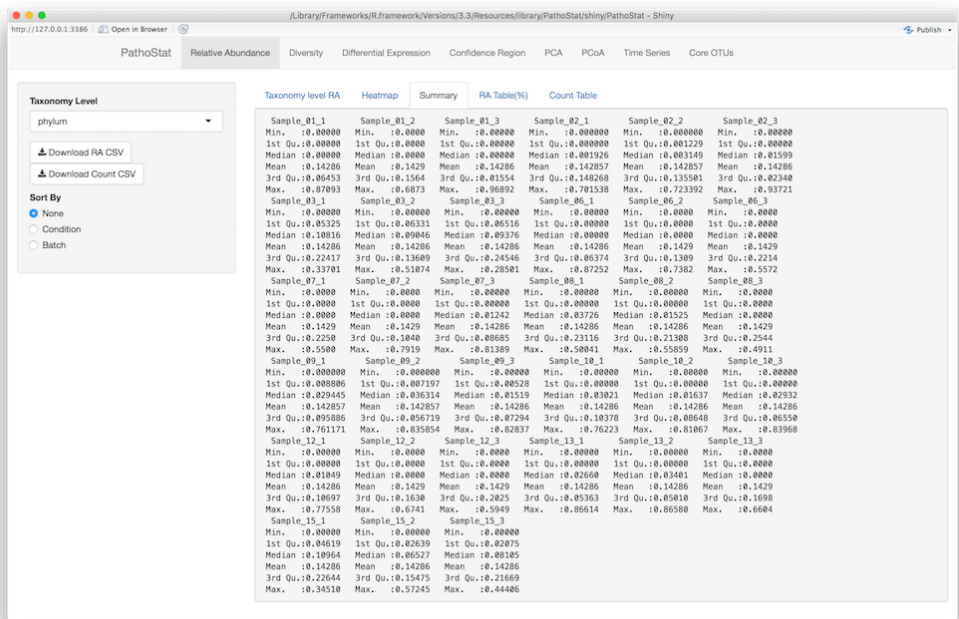
This subtab shows you read count abundance of selected taxonomic levels by sample. As with the previous subtab, you can select what taxonomic level you wish to plot by simply using the dropdown menu at the left of the user interface.



The legend on the right side of the plot shows the scale and the color gradient.

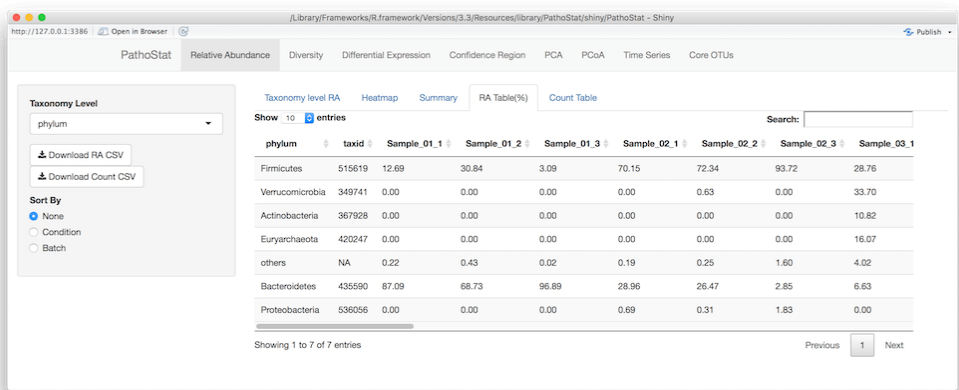
3.3 Summary

The Summary table shows summary statistics that may come in handy for detecting trends in the data and identifying outliers. As with the previous subtabs, Summary offers redundant information that points to the distribution of the taxonomic composition of your samples.



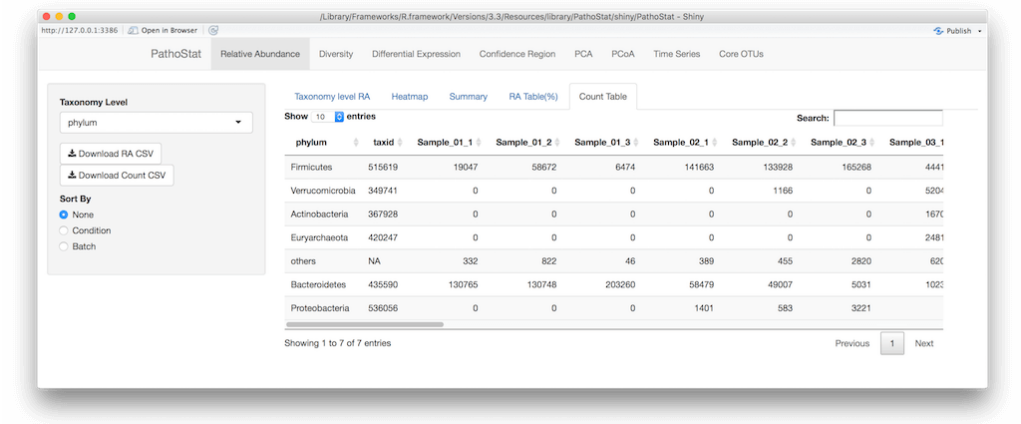
3.4 RA Table (%)

This subtab presents a searchable and sortable table of taxa abundance by sample. At user-selected taxonomy levels, you get information about the taxonomy ID (taxid) of the organism as per NCBI's taxonomy database and their relative abundance. All columns are sortable, and at the top right corner of the table you can use the search case to check for specific taxa.



3.5 Count table

This table is similar to the previous one but it shows abundance data as raw read counts instead of proportions. This is useful when you want to test for differential abundance of taxa between conditions. Statistical models such as those implemented in EdgeR and DESeq2 explicitly require count data to ensure specificity.

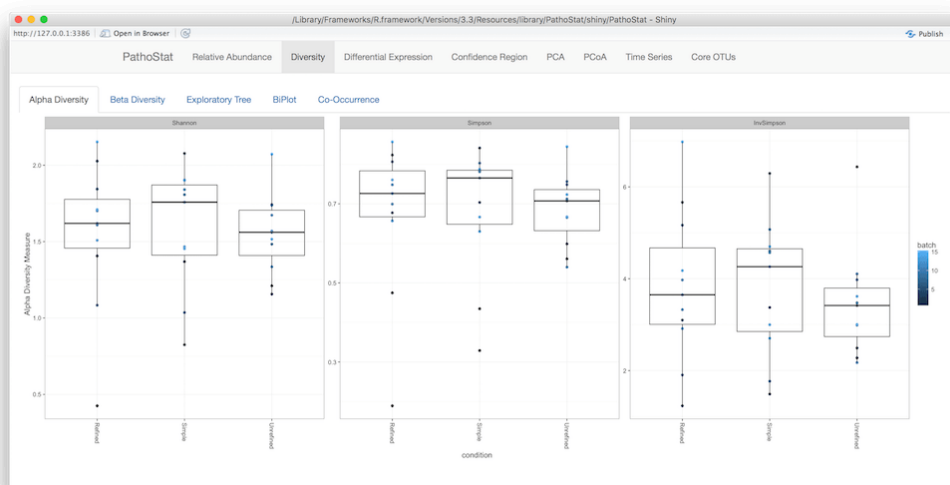


4 Diversity

The second tab in PathoStat is the Diversity tab. Here, users can obtain estimates of alpha and beta diversity, as well as explore a clustering dendrogram decorated with the abundance of certain taxonomic groups by sample. Also, the Biplot subtab illustrates the relationships between user-selected variables using multidimensional scaling. Finally, the Co-Occurrence subtab shows the relationships among taxa at user-defined distances allowing the exploration of potential biological associations.

4.1 Alpha Diversity

Alpha diversity considers presence and absence of taxonomit units as well as their homogeneity or evenness. In this subtab, you obtain estimates for alpha diversity using the Shannon, Simpson, and Inverse Simpson metrics as a function of a factor.



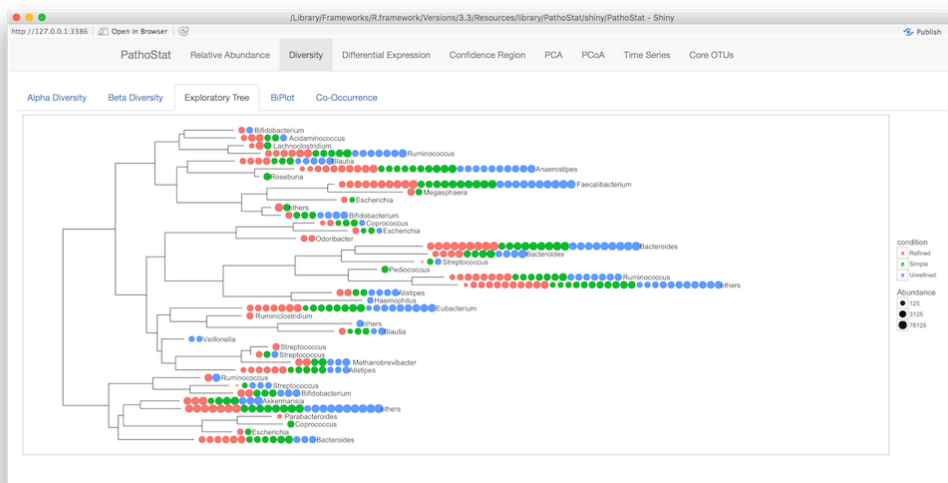
4.2 Beta Diversity

The beta diversity subtab provides the user with a heatmap of sample to sample variation using the Bray-Curtis distance metric. The color key indicates a normalized score (the row z-score) and the dendrogram is constructed using hierarchical clustering.



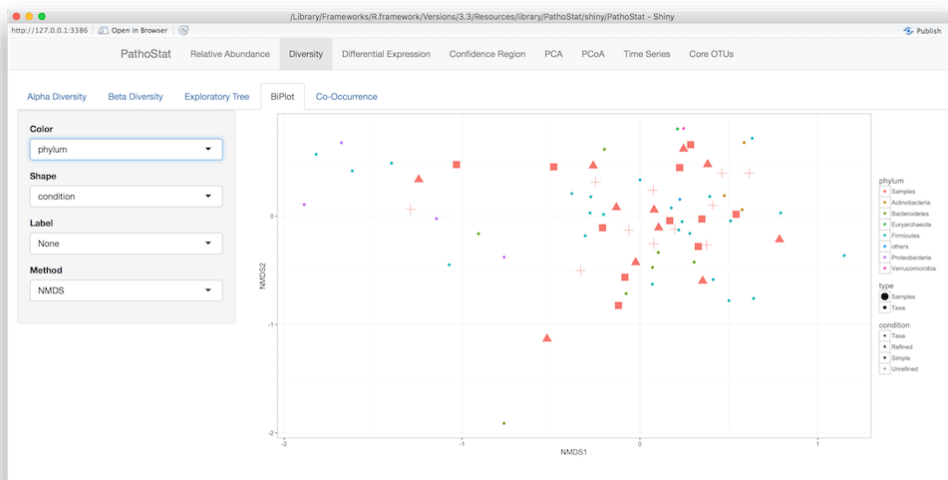
4.3 Exploratory Tree

This subtab shows a dendrogram (tree) among all the samples in your dataset. In the figure below, the tree is decorated with the most abundant taxon and their relative abundance by condition. This tree can reveal not only similarity among samples but also what taxon dominates each sample.



4.4 Biplot

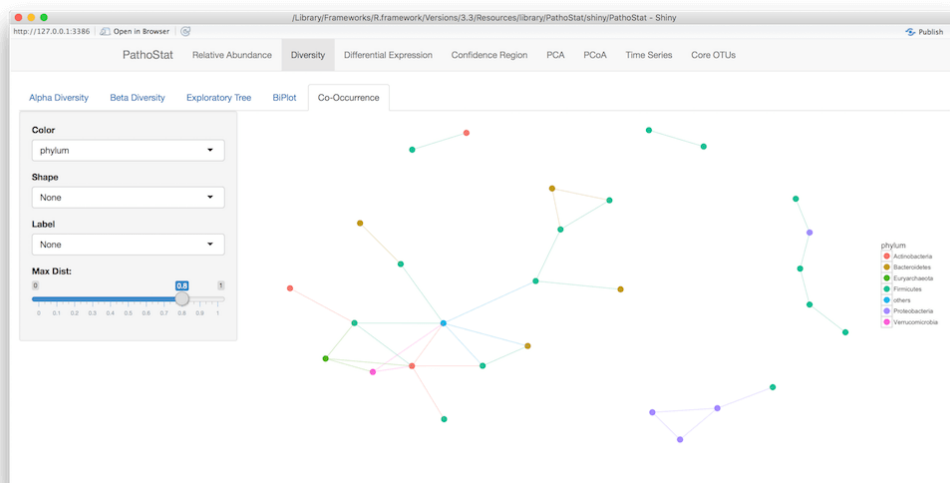
Another useful feature of PathoStat is the biplot. Here, the user can explore potential associations among taxa, and among taxa and factors in the dataset, e.g., condition, treatment, etc. The menu on the left side of the panel lets you select the variables to explore as well as the method to estimate the distance among variables.



4.5 Co-Occurrence

One recurrent question in metagenomics is whether two or more taxa tend to co-occur among a set of samples. In general, one might suspect that if two or more taxa co-occur then there might be a functional relationship among them. In the Co-Occurrence subtab, the user selects

relevant taxonomic levels to be compared and the maximum distance that will connect two taxa together. The result is a network graph where edges are proportional to the strength of co-occurrence.

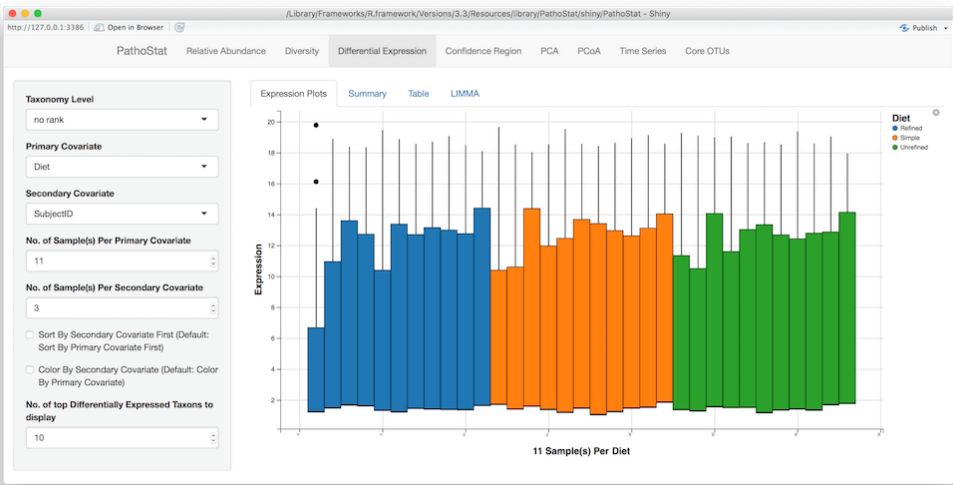


5 Differential Expression

The third tab in PathoStat allows you to test for differential abundance of taxa between conditions. This test is analogous to differential expression of genes in transcriptomic experiments and uses similar statistical models.

5.1 Expression Plots

In the Expression Plot subtab, you first select the taxonomy level of choice for the comparison and then you select the primary and secondary covariates. The results are shown as BoxPlots colored by covariate.

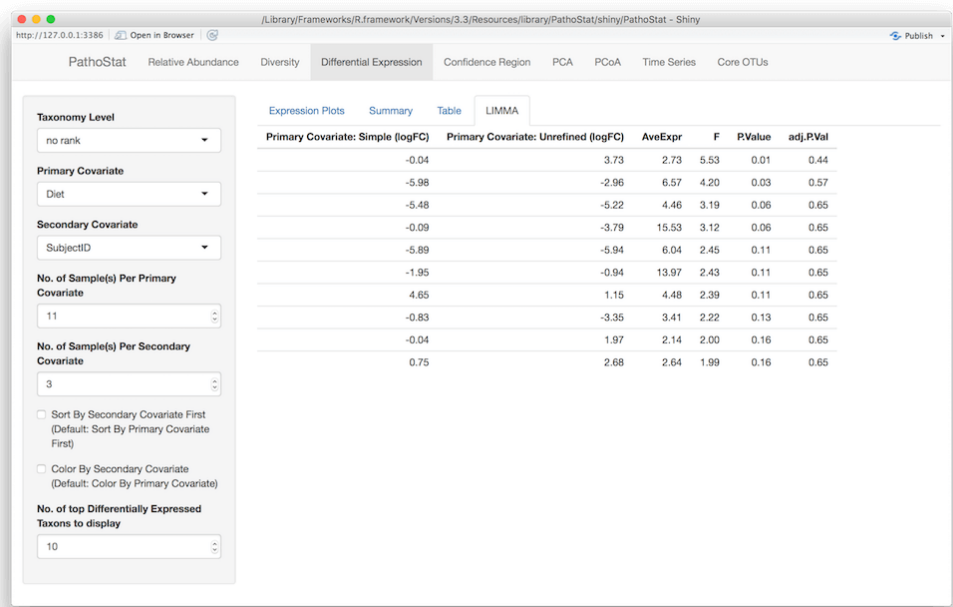


5.2 Summary, Table and LIMMA

The Summary subtab shows text summary of the data used for the BoxPlots in the Expression Plots subtab. The next two subtabs summarize the results as a table of raw values and as a Fold Change table, respectively.

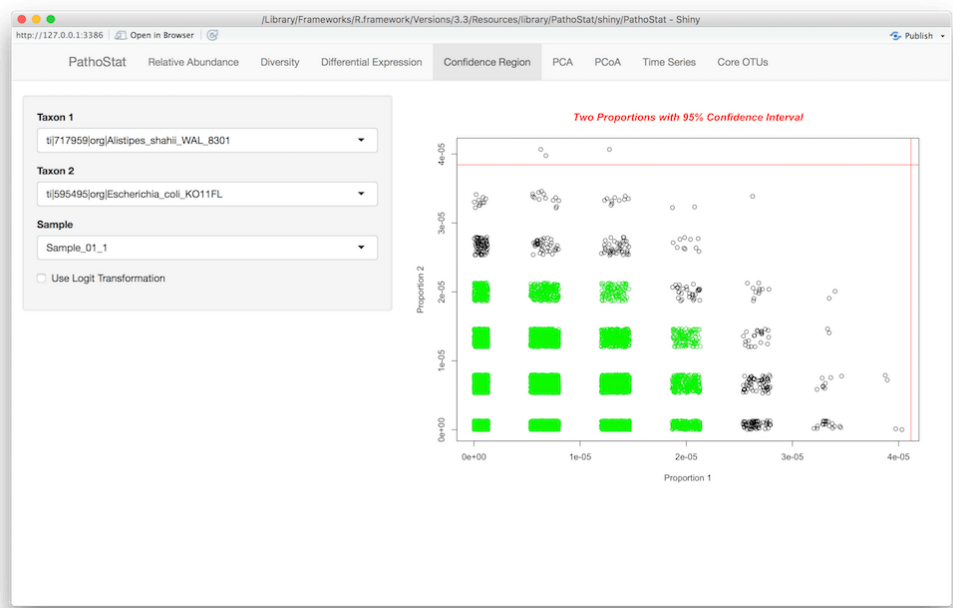
The screenshot shows the PathoStat Shiny app interface with the 'Table' subtab active. The table displays raw values for gene expression levels across 11 samples per diet, grouped by Diet (Refined, Simple, Unrefined). The table has 11 columns for each diet and 11 rows for each sample size. The values are displayed in a grid format.

	3	6	7	11	14	18	20	22	27	30	31	1	5	9	12	15
1.25	1.50	1.70	15.98	15.94	13.77	1.47	14.12	14.74	13.54	1.66	1.74	1.43	1.62	17.19	15.84	
1.25	12.26	1.70	1.64	1.35	1.25	1.47	1.43	1.41	1.39	1.66	1.74	1.43	1.62	1.41	12.60	
1.25	1.50	1.70	1.64	1.35	1.25	1.47	1.43	1.41	1.39	1.66	1.74	1.43	1.62	1.41	1.22	
8.57	1.50	1.70	1.64	1.35	1.25	1.47	1.43	1.41	1.39	1.66	1.74	1.43	1.62	1.41	1.22	
1.25	18.87	1.70	18.32	16.92	17.93	18.15	1.43	17.09	16.42	1.66	1.74	18.51	1.62	16.29	16.68	
1.25	1.50	18.36	1.64	1.35	1.25	15.88	1.43	1.41	1.39	17.07	1.74	12.62	18.01	1.41	1.22	
1.25	1.50	16.72	1.64	1.35	1.25	1.47	1.43	1.41	1.39	16.50	1.74	1.43	15.31	1.41	1.22	
6.68	1.50	1.70	1.64	1.35	1.25	1.47	1.43	1.41	1.39	1.66	11.20	1.43	1.62	1.41	1.22	
1.25	1.50	1.70	1.64	1.35	1.25	1.47	1.43	1.41	1.39	16.88	1.74	1.43	1.62	1.41	1.22	
1.25	10.95	13.60	12.73	1.35	1.25	14.13	1.43	1.41	14.52	15.25	1.74	13.80	14.78	11.96	1.22	
1.25	1.50	17.29	1.64	1.35	1.25	15.69	1.43	1.41	1.39	1.66	1.74	1.43	17.79	1.41	1.22	
1.25	12.93	1.70	1.64	1.35	1.25	1.47	1.43	1.41	1.39	1.66	1.74	1.43	1.62	1.41	1.22	
1.25	14.60	1.70	1.64	1.35	1.25	1.47	1.43	1.41	1.39	1.66	1.74	1.43	1.62	1.41	1.22	
1.25	1.50	1.70	1.64	1.35	1.25	1.47	1.43	1.41	1.39	1.66	1.74	1.43	1.62	1.41	1.22	
1.25	1.50	1.70	1.64	1.35	1.25	1.47	1.43	1.41	1.39	1.66	1.74	1.43	1.62	1.41	1.22	
1.25	1.50	1.70	1.64	1.35	1.25	1.47	1.43	1.41	1.39	1.66	1.74	1.43	1.62	1.41	1.22	
1.25	1.50	14.53	16.90	13.42	13.97	13.51	1.43	1.41	16.94	15.32	1.74	11.04	1.62	16.47	13.84	
1.25	1.50	1.70	1.64	1.35	1.25	1.47	1.43	1.41	1.39	1.66	1.74	1.43	1.62	1.41	1.22	
1.25	1.50	1.70	1.64	14.27	1.25	1.47	1.43	13.99	1.39	1.66	14.24	1.43	1.62	1.41	13.75	



6 Confidence Region

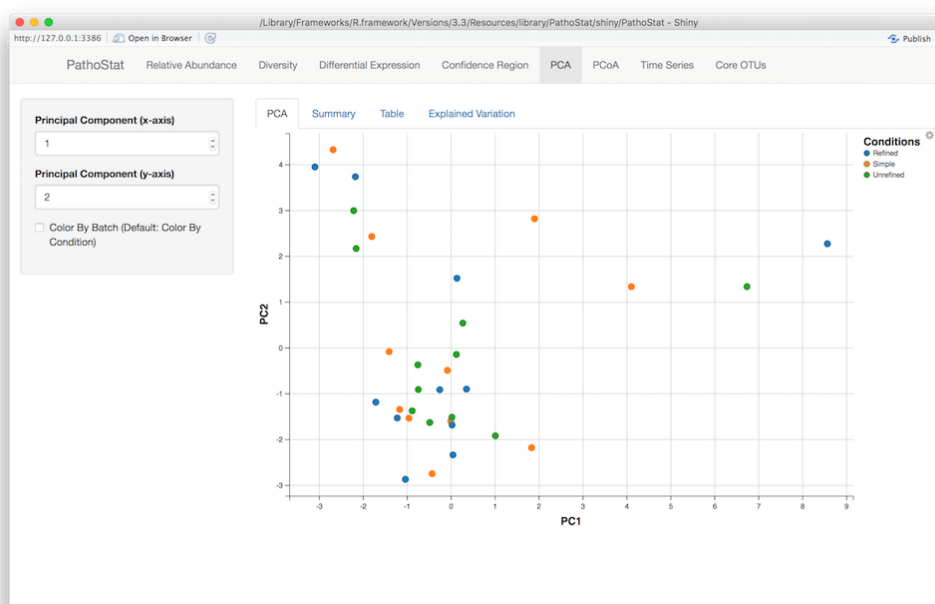
Many times researchers are interested in the accuracy of taxon abundance estimates. In this subtab, we provide a way to compare within-sample taxa in terms of their abundance estimate and 95% confidence interval.



On the left-side menu, you can select both the sample and the taxa to be used in the comparison. The results are plotted as a jitter plot that indicates the 95% confidence interval between the selected taxa.

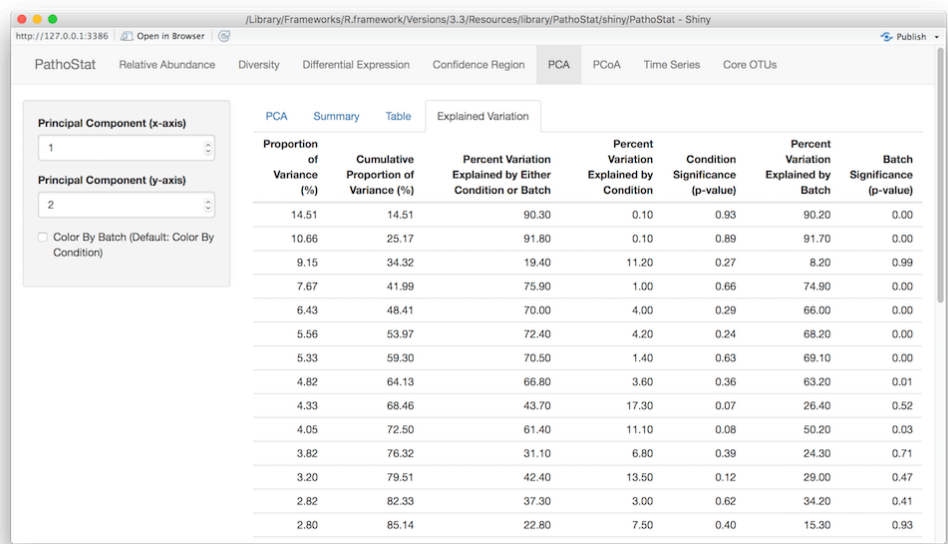
7 PCA and PCoA

The next two tabs in PathoStat calculate multidimensional scaling using Principal Components and Principal Coordinates Analysis. These tools can help you identify overall trends in the data. For instance, you can explore whether your samples are related by some biological condition of interest or by technical batch, in which case you would need to denoise the data first using methods such as surrogate variable analysis.



7.1 Explained Variation

An important aspect of PCA analysis is to understand to what extent the resulting vectors can explain the observed multidimensional variation. In this subtab, you can get an idea of that by exploring the different columns in terms of Proportion of Variance, Cumulative Variance and percent variation, among others.



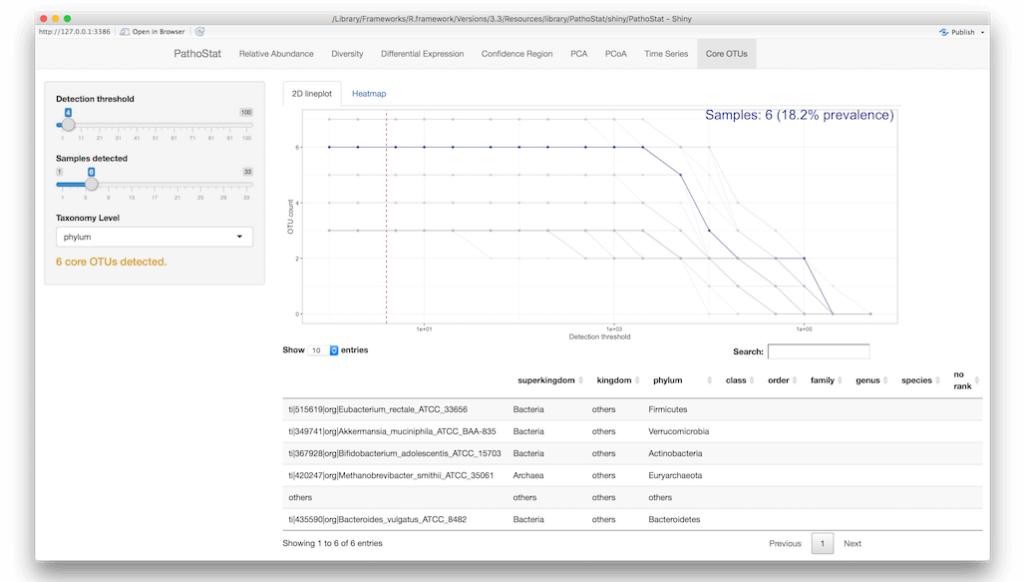
8 Time Series

The Time Series tab allows you to understand variation in relative abundance as a function of time. Strictly speaking, you could use any discrete numerical variable from your dataset for visualization. In the example below, we use SubjectID (not a numerical variable though) and agglomerate the data by Phylum.



9 Core OTUs

Many times we would like to understand what is particular and what is general about a set of metagenomic samples. Analogous to pangenome analysis in comparative genomics, the Core OTU tab allows you to identify shared taxa or core taxa among samples.



You first set up a detection threshold, meaning a minimum abundance proportion a OTU must have in order to be considered in the analysis, and the number of samples selected from your dataset. You can also select the taxonomic level of interest at which you would like to perform the analysis.