

Serial Group Comparisons Script

Task

Calculate non-parametric ANOVA (Kruskal-Wallis Rank Sum Test) and Fisher tests across the input numerical variables over a selected categorical variable.

Background

A common objective of microbial profiles analysis is the comparison of variables among groups of samples sharing a certain characteristic or treatment in order to detect differences in composition and abundances. This can be determined by performing an Analysis of Variance (ANOVA) type of test to establish, based on the values seen across the groups, how likely it is for the values in those samples to originate from different distributions. As a parametric test, classical ANOVA assumes normality of distribution. Since this is rarely the case for OTU data, we use the non-parametric Kruskal-Wallis Rank Sum Test in Rhea [1]. When more than two groups are compared, pairwise tests are needed to determine which of the groups are significantly different. Again, we use a non-parametric test (Mann-Whitney Test [1]) therefore. The obtained pairwise test significance values are corrected for multiple testing using the Benjamini-Hochberg method [2] and are reported together with the original values. Rhea was designed to perform a systematic testing of all available OTUs or taxonomies in a given experiment. This results in many tests and per extension in a high trade-off for p-values correction. Hence, a reduction of tests can be applied by removing unnecessary tests using e.g. prevalence cutoffs, since it has been shown that pre-filtering datasets increase the power of analysis [3].

References

1. Myles Hollander and Douglas A. Wolfe (1973), Nonparametric Statistical Methods. New York: John Wiley & Sons. Pages 115–120
2. Benjamini, Y., and Hochberg, Y. (1995). Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society Series B* 57, 289–300
3. Bourgon, R., Gentleman, R., & Huber, W. (2010). Independent filtering increases detection power for high-throughput experiments. *Proceedings of the National Academy of Sciences*, 107(21), 9546-9551.

Input

The input of the script is a tab-delimited text file following the format shown in the picture below. This table is the result of combining the OTU and Taxonomic binning relative abundance tables with the *alpha*-diversity and others meta-variables. The script create_input_table.R can assist in the preparation of the final input table by combining the specified files. Since... something missing here?!?!?!?!?

input_filename - The name of the table file used as input

independent_variable_name - Name of the column with the categorical variable (groups) used for the comparison to detect differences among the dependent numerical variables.

dependant_variables_start - The column where the dependant numerical variables start. The first column containing sample names do not count (see picture below).

taxonomic_variables_start - The column number where the taxonomic variables start. The first column containing sample names do not count (see picture below).

The diagram illustrates the structure of the input table with 10 columns. Brackets above the columns group them into three categories: 'Independent Categorical Variables' (columns 1-2), 'Meta Variables' (columns 3-7), and 'Taxonomic Variables' (columns 8-10). A green arrow labeled 'Position' points to column 10. A green arrow labeled 'Samples Names' points to column 1. A bracket below columns 3-7 is labeled 'Dependent Numerical Variables'. A red arrow points down from column 3 with the text 'Numeric Variables starts at 3'. Another red arrow points down from column 8 with the text 'Taxonomic Variables starts at 8'.

	1	2	3	4	5	6	7	8	9	10
#SampleID	Facility	Diet	Cholic_acid	Muricholic_acid	Deoxycholic_acid	Richness	Shannon.effective	OTU_1	OTU_3	OTU_12
11.HFW.HFD	HFW	HFD	1028.571004	3329.622948	4355.408455	109	47.13	8.69233	1.752794	0.14328
12.HFW.CD	HFW	CD	773.1665828	2945.613748	1810.557744	92	35.5	0.231237	13.24084	1.543256
13.SPF.CD	SPF	CD	342.5568553	1026.617105	2560.313083	101	30.45	5.925529	5.422457	4.222506
14.SPF.CD	SPF	CD	227.4375746	1436.135551	2874.451237	90	24.02	2.237588	11.0823	1.094785
20.SPF.HFD	SPF	HFD	134.5116952	2038.89837	4556.363049	104	33.06	12.96095	9.707729	1.747014
21.SPF.HFD	SPF	HFD	258.5543231	1316.12489	3578.413007	94	31.59	18.8355	1.837128	0.069139
3.HFW.CD	HFW	CD	31.1397895	1789.611976	1926.024534	93	23.55	0.02854	2.711342	8.493635
4.HFW.HFD	HFW	HFD	96.57664502	1523.306444	3984.614981	113	47.07	9.297885	4.837446	0.259499

Options

Several options are available to increase flexibility of the script. These options are set to default values that are most often used in common analysis, but can be changed if required. In more detail:

abundance_cutoff - The minimum relative abundance for taxonomic variables to be considered as effectively present. Variables less than this cut-off are zeroed. This masking of borderline abundances helps focusing the comparisons on samples where the taxonomic variables are important components of the communities. This cut-off should be adjusted to the environment and experiment of interest. The default setting (0.5) is a proposal for explorative studies pertaining to the mouse and human gut microbiota.

prevalence_cutoff - The minimum prevalence (number of samples positive for the given variable) in at least one group in the study. If the variable does not appear in more than the cut-off in any sample group, it is not considered for statistical testing. The default value is set to 0.3 (30%) prevalence in at least one group.

max_median_cutoff - If the median relative abundance of the taxonomic variable across all groups is less than this cut-off, it is not tested. This filter removes variables that have generally very low abundances across all samples. The default value is 1%.

ReplaceZero - This option determines the treatment of zero abundances in taxonomic variables. Replacing zeros with NA ("YES") will remove them from the statistical calculations, while selecting "NO" will treat them as true values in all calculations.

PlotOption - This option controls the graphical output of the variables showing significant differences across groups. There are three possible choices (1, 2 and 3). Selecting 1 will plot boxplots and violin plots without showing individual data points. Selecting 2 will produce boxplots and violin plots showing individual data points. Selecting 3 will add the samples names over each individual data point.

Output

Each time the script is executed, the output is placed in a new folder names according to the variable/group used for comparisons and a date-time stamp. The files included in the output folder are the following:

plot_box.pdf - Boxplots of all significant comparisons (before correction).

plot_point.pdf - Dot-plots of all significant comparisons (before correction).

plot_violin.pdf - Violin plots of all significant comparisons (before correction).

my_analysis_log.txt - A text file capturing all the options used in the analysis for future reference

The remaining files have generic names deriving from the name of the input file and the categorical variable used for analysis. If for example the input file name was "OTUsCombined" and the grouping variable was "Diet", then the outputs are:

OTUsCombined-Diet-FisherTestAll.tab - A tabular file with the calculated p-value for the Fisher test for all the variables. A column with adjusted p-values for multiple testing is also calculated.

OTUsCombined-Diet-FisherTestPairWise.tab - A tabular file with the calculated p-value for the Fisher test for all the pairs of groups variables. A column with adjusted p-values for multiple testing is also calculated.

OTUsCombined-Diet-modified.txt - The modified input table after all filters and transformations were applied.

OTUsCombined-Diet-pvalues.tab - A tabular file with the calculated p-value for the Kruskal-Wallis Rank Sum test for all the variables. A column with adjusted p-values for multiple testing is also calculated.

OTUsCombined-Diet-sign_pairs.tab - A tabular file with the calculated p-value for the Mann-Whitney test for all the pairs of groups variables. A column with adjusted p-values for multiple testing is also calculated.

Important Notes

Please pay attention to the formatting of your input table and the ordering of the variables (categorical, numerical-nonSeq, numerical-sequence). To determine the position where dependant variable (numerical) start, do not consider the sample names in the first column (1 is the first column after sample names and so on). The decision to use or remove zero and near-to-zero values has important consequences for the results and is left to the discretion of users. Always remember that only variables with significant Fisher or Kruskal-Wallis test are plotted. If no graphical output is produced, it means that not a single variable delivered significant results, which can be also checked by opening the tabular files.

Common problems

- The path to the script is not set correctly
- The input file names are incorrect
- The column name selected for grouping does not exist or contains typos.
- Only one group or too few samples are available for statistics