# MicrobiomeAnalyst 2.0

Comprehensive statistical, functional and integrative analysis of microbiome data

xialab@mcgill 2023-Mar-03

# Tutorial for Statistical Meta-analysis

# Motivation

- Increasing microbiome studies result in tremendous data designed for understanding different experimental variables, such as diseases and environment pressure, associated with changes in microbial community structure.

- However, it remains a major challenge to achieve reproducible features across different microbiome studies due to the variation in experimental design, analysis methods and quantitative assessment.

- There is an unmet need for analytical tools that provide rigorous statistical analysis dedicated to mine available data against same hypothesis and obtain consistent interpretation across different studies.
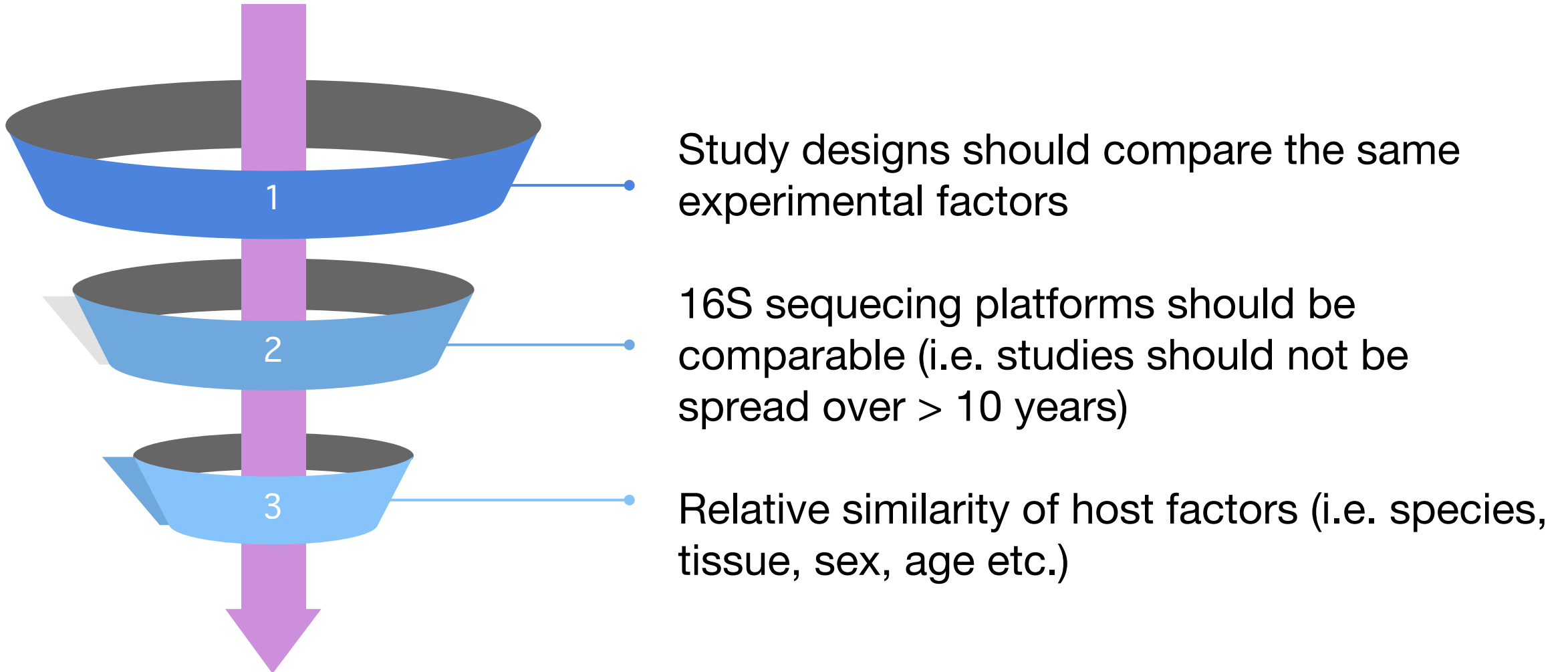
# Overview

Goal: To provide a framework for integrating multiple maker gene studies to help identify robust and reproducible features from multiple microbiome studies.

Strategy and Approach:

- The MMUPhin method is employed to alleviates batch effects in the joint analysis of microbial profiles. It adjust for differences in technical or experimental variation between studies by considering batch/study effects which can significantly increases the comparability of different microbiome studies.

- Three analysis tracks are offered for user to explore the consistent pattern and potential biomarkers – visual exploration, diversity meta-analysis, and biomarker meta-analysis.

# Datasets selection



1. Study designs should compare the same experimental factors

2. 16S sequecing platforms should be comparable (i.e. studies should not be spread over > 10 years)

3. Relative similarity of host factors (i.e. species, tissue, sex, age etc.)

# Data format: data table

**Sample names** →

**OTU ids**

| #NAME | Sample1 | Sample2 | Sample3 | Sample4 | Sampl5 | Sampl6 | Sample7 | Sampl8 |
|-------|---------|---------|---------|---------|--------|--------|---------|--------|
| OTU1  | -3.06   | -2.25   | -1.15   | -6.64   | 0.4    | 1.08   | 1.22    | 1.02   |
| OTU2  | -1.36   | -0.67   | -0.17   | -0.97   | -2.32  | -5.06  | 0.28    | 1.32   |
| OTU3  | 1.61    | -0.27   | 0.71    | -0.62   | 0.14   |        | 0.11    | 0.98   |
| OTU4  | 0.93    | 1.29    | -0.23   | -0.74   | -2     | -1.25  | 1.07    | 1.27   |

...

Please take a look at these example data tables:

https://www.microbiomeanalyst.ca/MicrobiomeAnalyst/resources/data/metaanal/data1.csv
https://www.microbiomeanalyst.ca/MicrobiomeAnalyst/resources/data/metaanal/data2.csv
https://www.microbiomeanalyst.ca/MicrobiomeAnalyst/resources/data/metaanal/data3.csv

# Data format: meta-data table

Primary meta-data

| #NAME | study_condition | age |
|-------|-----------------|-----|
| SID31004 | CRC | 64 |
| SID31009 | control | 68 |
| SID31021 | control | 60 |
| SID31071 | control | 68 |
| SID31112 | control | 66 |
| SID31129 | control | 73 |
| SID31159 | CRC | 73 |

Sample names

…

The primary meta-data needs to be consistent across datasets.
Only supports case-control Design (two factors)

https://www.microbiomeanalyst.ca/MicrobiomeAnalyst/resources/data/metaanal/data1_meta.csv
https://www.microbiomeanalyst.ca/MicrobiomeAnalyst/resources/data/metaanal/data2_meta.csv
https://www.microbiomeanalyst.ca/MicrobiomeAnalyst/resources/data/metaanal/data3_meta.csv

# Upload data

The first step is to upload and process all your individual datasets. This repeats the steps of a single marker data profiling for each dataset - for more details on each step, see the corresponding tutorial. It is advised to upload raw counts to access all analysis options.

Upload your dataset 1 by 1, make sure that at least one meta-data group is shared across dataset and consists of two factors (case-control)



For the purpose of this tutorial, try our example data

# Example data

The example datasets come from stool samples of three 16S colorectal cancer studies; the datasets have been trimmed for testing purposes.



Uploaded datasets will be displayed here on the left panel

# Data Summary

This page provides general text summary and library size graphical overview on the uploaded datasets



Available file downloads for each page are displayed here

# Data processing

Data processing page offers the same filtering and normalization options available for single gene marker profiling with the addition of batch effect correction to remove study-specific bias

You can perform filtering and normalizine on all datasets at once or one by one.

On each dataset, we show the progress of data processing. (Incomplete vs Finished)

Navigate to:

**Uploaded datasets**

data1.csv
Feature: 435
Sample: 107
Norm. Input: No
Incomplete

data2.csv
Feature: 196
Sample: 55
Norm. Input: No
Incomplete

data3.csv
Feature: 400
Sample: 104
Norm. Input: No
Incomplete

**Downloads of the page**

**Data processing**

By default, all uploaded datasets are processed using the default parameters (see below). You can use the table be...
Scroll down to see graphical summaries of individual omics datasets.

Currently selected data:  | All Datasets ∨ |   Status: Incomplete

| Processing Step | Parameter Selection | | Action |
|---|---|---|---|
| Filtering ❓ | Variance filter | ○————— [0] | Submit |
| | | Minimum count: ○—— [0] | |
| | Abundance filter | ● Prevalence in samples (%)○——— [10] | |
| | | ○ Mean abundance value | |
| | | ○ Median abundance value | |
| Normalization ❓ | Data rarefying ❓ | Do not rarefy my data ∨ | Submit |
| | Data scaling ❓ | Total sum scaling (TSS) ∨ | |
| | Data transformation ❓ | Do not transform my data ∨ | |

**Adjust study batch effect** ☐   Update

**PCA Overview**  Density Plot

‹‹ Previous    ›› Proceed

# Data processing

# Methods Selection

Note that some methods can only be performed on counts data (i.e. biomarker meta-analysis, alpha diversity, stacked area/taxa abundance bar)

You can choose to exclude some of the datasets before performing analysis

**Uploaded datasets**

data1.csv
Feature: 435
Sample: 107
Norm. Input: No

data2.csv
Feature: 196
Sample: 55
Norm. Input: No

data3.csv
Feature: 400
Sample: 104
Norm. Input: No

**Downloads of the page**

No downloads on this page.

**R Command History**

**Navigate to:**

**Please choose a meta-analysis method to proceed**

**Visual exploration**

Visually explore your data sets through stacked bar/area plot or PCoA plots. It permits both overall patterns as well as sample-level details through zoom and mouse-over interactions

Visualization method: [ Stacked bar/area plot ⌄ ]   Submit   Select projection dataset

**Biomarker meta-analysis**

Identify consistent changes across different data sets. It performs regression analysis in individual studies using MaAsLin2, and then aggregate results with fixed/mixed effect models using MMUPHin.

Differential analysis  [ Linear modeling (LM) ⌄ ]   Submit
Meta-analysis method  [ Random Effect Model ⌄ ]

**...sity meta-analysis**

...ute alpha- and beta- diversity across different datasets, the overall trend, as well as to evaluate the consistency of communities (discrete) or gradients
...nuous structure)

...Diversity option:  [ Alpha Diversity ⌄ ]   Submit

The graphical overview displays a maximum of top significant features. Detailed table contains the results for all features. You can also download the result table in the "Downloads of the page" tab.

# Biomarker meta-analysis

Biomarker meta-analysis

You can visualize overall abundance profiles of individual feature using our Detailed table, under "View" column

# Alpha diversity analysis

This module offers two graphical representation: 1) box plot displays the distribution of diversity metrics; 2) log2 ratio view displays results from comparative analysis between case-control. Detailed table provides more information on the statistical results

# Beta diversity analysis

This module applies PCoA of beta diversity distance matrices along with statistical testing to measure significance on the effect of phenotype on community composition.

The title of each PCoA plot contains result from statistical testing

# Projection to public dataset



The "Projection to public dataset" module has been merged here. To try out this feature, try our second example data that has compatible IDs (i.e. taxonomy id). Our first example data do not have compatible IDs with the collected public datasets.

Only the two methods from "Visual Exploration" are available for this feature.

**Select public dataset**

Project public dataset with uploaded dataset for visual exploration. Make sure to select data with similar conditions. A basic data check will be performed to check the compatibility.

Human Gut | Human Skin | Human Oral | Human Vagina | Mouse | Cow | Environment

| Studies | Target region | Sequence platform | No. of samples | Ref. |
|---------|---------------|-------------------|----------------|------|
| Healthy_whole_body | V2 | 454 GS FLX | 45 | Costello et al. 2009 |
| Dense_timeseries | V4 | Illumina HiSeq 2000 | 467 | Caporaso et al. 2011 |
| HMP_V35 | V3-5 | 454 GS FLX Titanium | 371 | HMP 2012 Consortium |
| HMP_V13 | V1-3 | 454 GS FLX Titanium | 204 | HMP 2012 Consortium |
| Global_gut | V4 | Illumina HiSeq 2000 | 528 | Yatsunenko et al. 2012 |
| Family_study | V2 | Illumina HiSeq 2000 | 169 | Song et al. 2013 |
| Diet_enterotype | V2 | 454 GS FLX Titanium | 85 | Wu et al. 2011 |
| Pregnant_women | V2 | 454 GS FLX and GS FLX Titanium | 667 | Koren et al. 2011 |
| Newborns_and_mothers | V2 | 454 GS FLX | 80 | Dominguez-Bello et al. 2010 |
| US_infant_timeseries | V2 | 454 GS FLX | 61 | Koenig et al. 2011 |
| Obese_twins | V2 | 454 GS FLX | 281 | Turnbaugh et al. 2009 |
| IBD_twins | V2 | 454 GS FLX | 114 | Willing et al. 2010 |

Add

Sample: 73
Norm. Input: No

Projection Data (optional)

🗑 costello_gut

˅ Downloads of the page

No downloads on this page.

› R Command History

xa Abundance › Downloads

**meta-analysis method to proceed**

Exploration

Visually explore your data sets through stacked bar/area plot or PCoA plots. It permits both overall patterns as well as sample-level details throu mouse-over interactions

Visualization method: [Stacked bar/area plot ˅]   Submit   Select projection dataset

**Biomarker meta-analysis**

Identify consistent changes across different data sets. It performs regression analysis in individual studies using MaAsLin2, and then aggregate results with fixed/mixed effect models using MMUPHin.

Differential analysis [Linear modeling (LM) ˅]   Submit

Meta-analysis method [Random Effect Model ˅]

**Diversity meta-analysis**

Compute alpha- and beta- diversity across different datasets, the overall trend, as well as to evaluate the consistency of communities (discrete) or gradients (continuous structure)
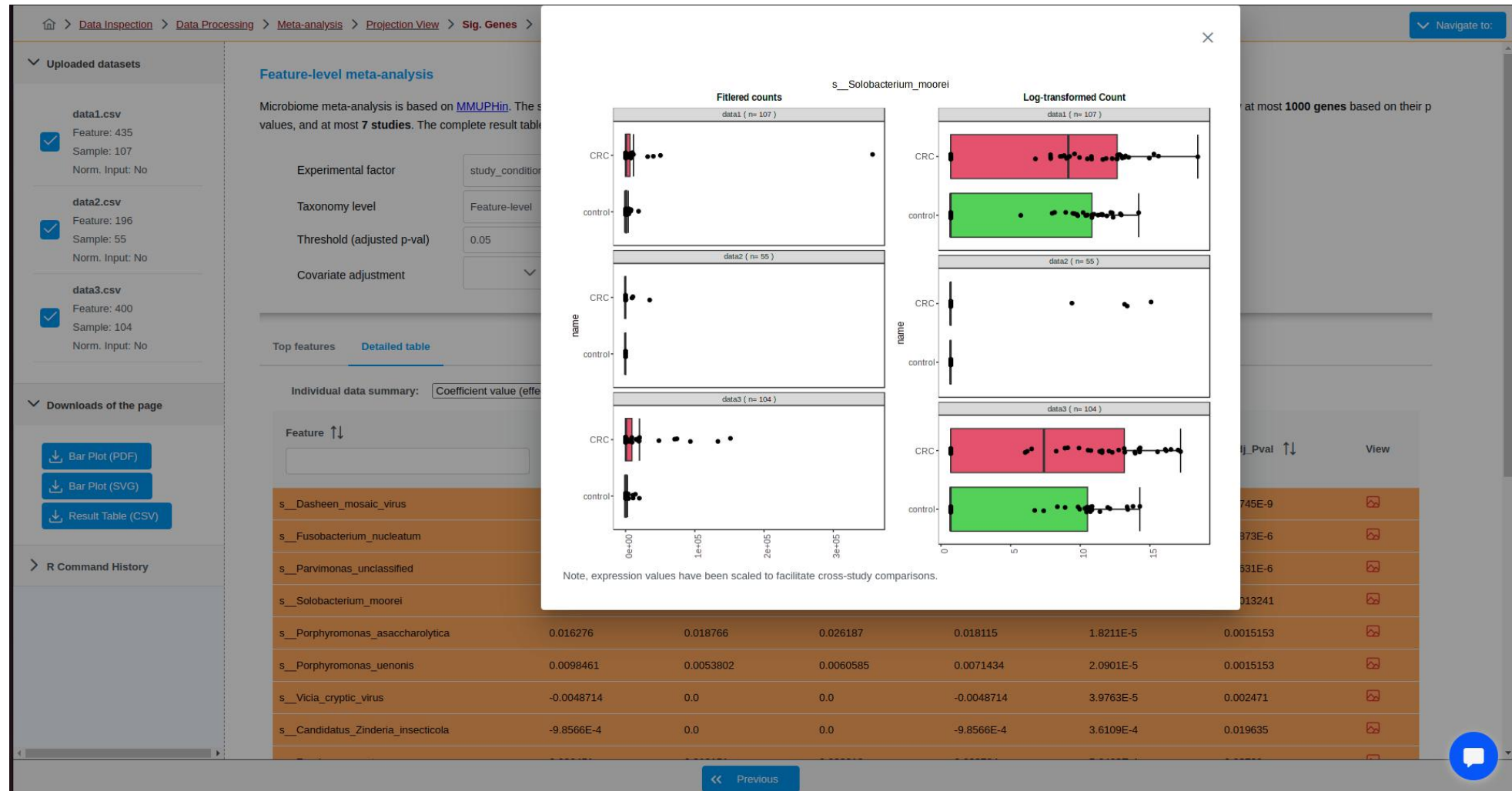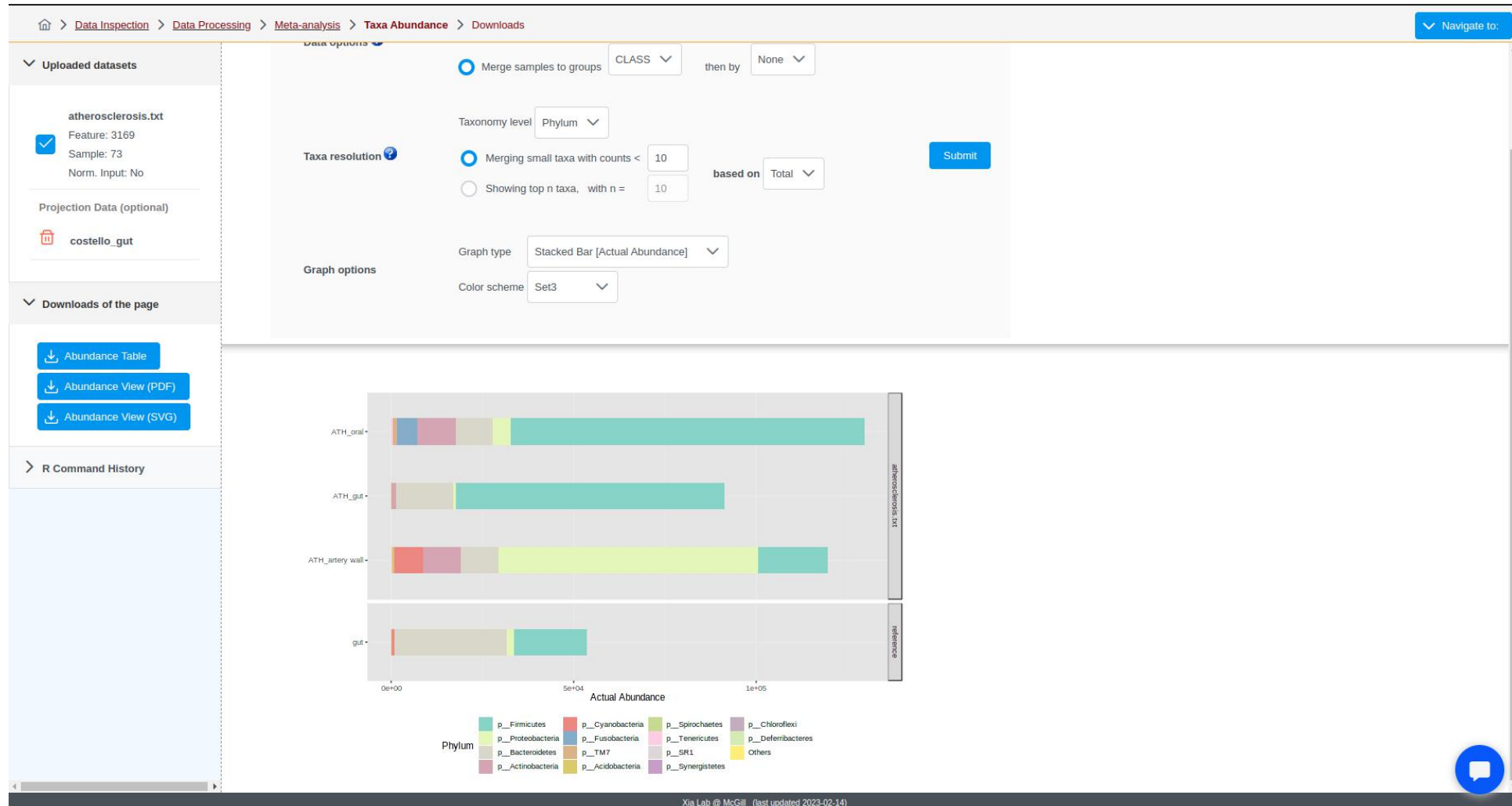
Diversity option: [Alpha Diversity ˅]   Submit
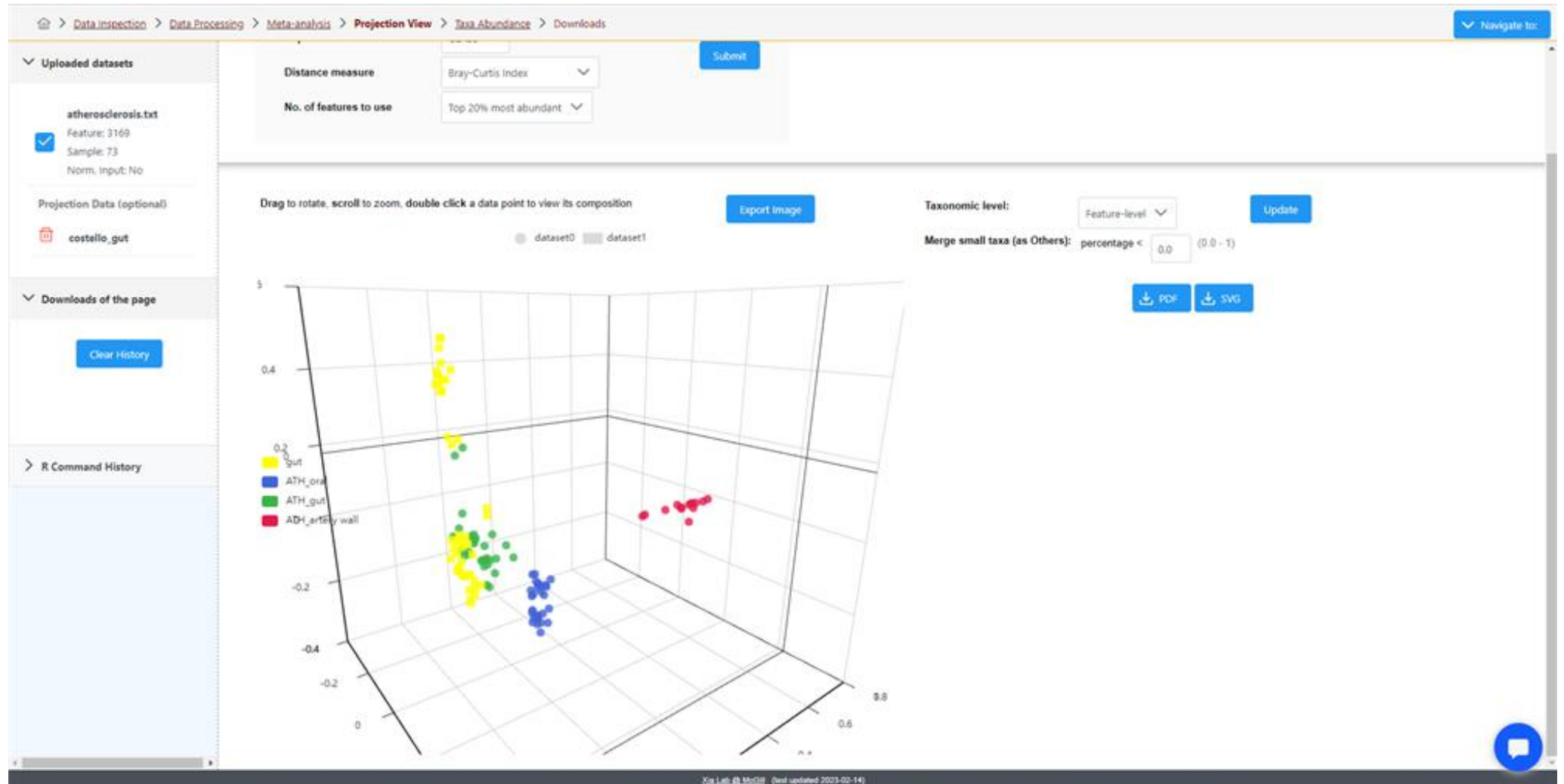
# Stacked bar plot

You can compare taxa abundance of your uplodaded data with reference data using "Stacked bar/area plot".

PCoA projection

Similarly, you can visualize beta-diversity community composition with reference dataset in PCoA space.

# The End

For more information, visit Tutorials, Resources and Contact pages on www.microbiomeanalyst.ca
Also visit our forum for FAQs on www.omicsforum.ca