### Advances in Microbiome Data Science with R/Bioconductor

Workshop, Tirana May 24, 2022



**Associate Prof. Leo Lahti** | datascience.utu.fi Department of Computing, University of Turku, Finland

# Microbiome bioinformatics

# Data science workflows

Open science & reproducible research

The demise of alchemy provides further evidence, if further evidence were needed, that what marks out modern science is not the conduct of experiments (alchemists conducted plenty of experiments), but the formation of a *critical community capable of assessing discoveries and replicating results*. Alchemy, as a clandestine enterprise, could never develop a community of the right sort. Popper was right to think that science can flourish only in an open society.

The Invention of Science: A New History of the Scientific Revolution, by David Wootton

Open reporting and communication were part of academic culture since the early days



A family of alchemists at work, an engraving by Philip Galle, after a painting by Pieter Bruegel the Elder, published by Hieronymus Cock, *c*.1558.

# Initial sequencing and analysis of the human genome $\sim 2001$

International Human Genome Sequencing Consortium

\* A partial list of authors appears on the opposite page. Affiliations are listed at the end of the paper.

The human genome holds an extraordinary trove of information about human development, physiology, medicine and evolution. Here we report the results of an international collaboration to produce and make freely available a draft sequence of the human genome. We also present an initial analysis of the data, describing some of the insights that can be gleaned from the sequence.





# FINRISK cohort (2002)

20 year follow-up

### Population register data







### Taxonomic signatures of cause-specific mortality risk in human gut microbiome

Aaro Salosensaari, Ville Laitinen, Aki S. Havulinna, Guillaume Meric, Susan Cheng, Markus Perola, Liisa Valsta, Georg Alfthan, Michael Inouye, Jeramie D. Watrous, Tao Long, Rodolfo A. Salido, Karenina Sanders, Caitriona Brennan, Gregory C. Humphrey, Jon G. Sanders, Mohit Jain, Pekka Jousilahti, Veikko Salomaa, Rob Knight, Leo Lahti⊠ & Teemu Niiranen ⊠

Nature Communications 12, Article number: 2671 (2021) Cite this article

## Gut microbiome and mortality risk



![](_page_5_Picture_2.jpeg)

#### There's A New Self-Administered Medical Craze For 2018 And It's Much More Horrifying Than Tide Pods

![](_page_6_Picture_1.jpeg)

![](_page_6_Picture_2.jpeg)

People are attempting DIY fecal transplants without medical supervision and putting themselves at risk of deadly infections, warn experts

### Gut microbiome similarity among Finnish adults

![](_page_7_Figure_1.jpeg)

- How many clusters?
- Which clusters?
- Where are the boundaries
- Who belongs to which cluster?

<u>How to choose a correct model?</u>  $\rightarrow$  a community typing example

![](_page_8_Figure_1.jpeg)

![](_page_9_Figure_0.jpeg)

#### Reproducible Research: Enterotype Example

Susan Holmes and Joey McMurdie

http://statweb.stanford.edu/~susan/papers/EnterotypeRR.html

#### Taxonomic Signatures of Long-Term Mortality Risk in Human Gut Microbiota

- 💿 Aaro Salosensaari, 💿 Ville Laitinen, 💿 Aki Havulinna, Guillaume Meric, 💿 Susan Cheng,
- 💿 Markus Perola, Liisa Valsta, 💿 Georg Alfthan, 💿 Michael Inouye, Jeramie D. Watrous, Tao Long,

O Comment on this paper

- D Rodolfo Salido, Karenina Sanders, Caitriona Brennan, Gregory C. Humphrey, Jon G. Sanders,
- 10 Mohit Jain, Pekka Jousilahti, 0 Veikko Salomaa, 0 Rob Knight, 0 Leo Lahti, 0 Teemu Niiranen doi: https://doi.org/10.1101/2019.12.30.19015842

#### The influence of hidden researcher decisions in applied microeconomics

Nick Huntington-Klein ➡, Andreu Arenas, Emily Beam, Marco Bertoni, Jeffrey R. Bloem, Pralhad Burli, Naibin Chen, Paul Grieco, Godwin Ekpe, Todd Pugatch, Martin Saavedra, Yaniv Stopnitzky

First published: 22 March 2021 https://doi.org/10.1111/ecin.12992

![](_page_10_Figure_3.jpeg)

Researchers make hundreds of decisions about data collection, preparation, and analysis in their research. We use a many-analysts approach to measure the extent and impact of these decisions. Two published causal empirical results are replicated by seven replicators each. We find large differences in data preparation and analysis decisions, many of which would not likely be reported in a publication. No two replicators reported the same sample size. Statistical significance varied across replications, and for one of the studies the effect's sign varied as well. The standard deviation of estimates across replications was 3-4 times the mean reported standard error.

P < 0.04P < 0.06P < 0.05Effect? No effect?

![](_page_11_Figure_1.jpeg)

(barely) not statistically significant (p=0.052) a barely detectable statistically significant difference (p=0.073) a borderline significant trend (p=0.09) a certain trend toward significance (p=0.08) a clear tendency to significance (p=0.052) a clear trend (p<0.09) a clear, strong trend (p=0.09) a considerable trend toward significance (p=0.069) a decreasing trend (p=0.09) a definite trend (p=0.08) a distinct trend toward significance (p=0.07) a favorable trend (p=0.09) a favourable statistical trend (p=0.09) a little significant (p < 0.1) a margin at the edge of significance (p=0.0608) a marginal trend (p=0.09) a marginal trend toward significance (p=0.052) a marked trend (p=0.07) a mild trend (p<0.09) a moderate trend toward significance (p=0.068) a near-significant trend (p=0.07) a negative trend (p=0.09) a nonsignificant trend (p<0.1) a nonsignificant trend toward significance (p=0.1)

![](_page_12_Figure_1.jpeg)

![](_page_13_Figure_0.jpeg)

![](_page_13_Picture_1.jpeg)

RESEARCH PRIORITIES Shining Light into Black Boxes

A. Morin<sup>1</sup>, J. Urban<sup>2</sup>, P. D. Adams<sup>3</sup>, I. Foster<sup>4</sup>, A. Sali<sup>5</sup>, D. Baker<sup>6</sup>, P. Sliz<sup>1,\*</sup>

![](_page_13_Picture_4.jpeg)

# How we choose which model to apply?

![](_page_14_Figure_1.jpeg)

I DON'T TRUST LINEAR REGRESSIONS WHEN IT'S HARDER TO GUESS THE DIRECTION OF THE CORRELATION FROM THE SCATTER PLOT THAN TO FIND NEW CONSTELLATIONS ON IT. Data science workflow reproducibility & automation sharing & collaboration transparency

![](_page_15_Figure_1.jpeg)

Program

REVISED Bioconductor Workflow for Microbiome Data Analysis: from raw reads to community analyses [version 2; peer review: 3 approved]

Ben J. Callahan<sup>1</sup>, Kris Sankaran<sup>1</sup>, Julia A. Fukuyama<sup>1</sup>, Paul J. McMurdie<sup>2</sup>, 🔀 Susan P. Holmes

![](_page_15_Picture_5.jpeg)

### open data science ecosystems

![](_page_16_Picture_1.jpeg)

![](_page_16_Picture_2.jpeg)

#### PeerJ >

# Anvi'o: an advanced analysis and visualization platform for 'omics data

Research article Bioinformatics Biotechnology Computational Biology Genomics Microbiology

A. Murat Eren<sup>≤1,2</sup>, Özcan C. Esen<sup>1</sup>, Christopher Quince<sup>3</sup>, Joseph H. Vineis<sup>1</sup>, Hilary G. Morrison<sup>1</sup>, Mitchell L. Sogin<sup>1</sup>, Tom O. Delmont<sup>1</sup>

![](_page_16_Picture_7.jpeg)

#### Anvi'o in a nutshell

![](_page_16_Figure_9.jpeg)

Anvi'o is an open-source, community-driven analysis and visualization platform for 'omics data.

# Varying cultures of open collaboration

![](_page_17_Figure_1.jpeg)

![](_page_18_Picture_0.jpeg)

Started 2001 as a platform for analysis & understanding of microarray data

More than 1,600 packages. Domains of expertise:

- Sequencing (RNASeq, ChIPSeq, single-cell, called variants, ...)
- Microarrays (methylation, expression, copy number, ...)
- Flow cytometry
- Proteomics
- Multi-Omics data integration

#### Important themes

- Reproducible research
- Interoperability between packages & workflows ... even from different authors

- Usability

![](_page_18_Picture_12.jpeg)

![](_page_18_Picture_13.jpeg)

![](_page_18_Figure_14.jpeg)

### "Omics" data taxonomic abundance table

*Omics* in Oxford English Dictionary: *in cellular and molecular biology*, forming nouns with the sense "all constituents considered collectively"

Gut microbiota: 1000 western adults (Lahti *et al.* Nature Comm. 2014)

Features x samples

![](_page_19_Figure_4.jpeg)

Genomics Epigenomics Microbiomics Lipidomics Proteomics Glycomics Foodomics Transcriptomics Metabolomics Culturomics

![](_page_20_Figure_0.jpeg)

![](_page_21_Figure_0.jpeg)

![](_page_22_Figure_0.jpeg)

### "Data container" (*TreeSummarizedExperiment*)

![](_page_23_Figure_1.jpeg)

### Reduce overlapping efforts, improve interoperability, ensure sustainability.

![](_page_24_Figure_1.jpeg)

# Package ecosystem

```
pheatmap(mat, annotation_row = taxa_clusters,
```

annotation\_col = sample\_data,

breaks = breaks,

color = colors)

![](_page_25_Figure_4.jpeg)

Figure 7.1: Prevalence of top phyla as judged by prevalence

plotRowTree(x[rowData(x)\$Phylum %in% top\_phyla\_mean,], edge\_colour\_by = "Phylum", tip\_colour\_by = "prevalence", node\_colour\_by = "prevalence")

![](_page_25_Figure_7.jpeg)

# Open & reproducible workflow

#### **Import Data** The TreeSE object The tse object is uniquely positioned to This workflow starts with either raw data directly from support the next generation of microbiome relative abundance estimation or taxonomic classification data manipulation and visualization. OR pre-existing data objects from widely used software. ples **RAW DATA** . 8 samples metadata taxa tree Taxonomic Tree Row Row Samples (rowTree) Link Data (Columns) colData(tse) counts . 0 taxc Taxa Rows) **EXISTING DATA** ological **bservation** metadata(tse) atry rowLinks(tse) rowData(tse) rowTree(tse) aggValue(tse) phyloseq assays(tse) subsetByNode(tse)

#### The mia Pipeline

Accessing Taxonomic Info.

![](_page_26_Figure_4.jpeg)

[1] mia::addTaxonomyTree(tse)

[2] TreeSE::aggValue(tse)

#### **Quality Control**

![](_page_26_Picture_7.jpeg)

[3] scatter::addPerCellQC(tse)

#### Visualizing with miaViz

![](_page_26_Figure_10.jpeg)

[4] miaViz::plotRowTree(tse)

![](_page_27_Picture_0.jpeg)

This image was created by Scriberia for The Turing Way community DOI: 10.5281/zenodo.3 332807. Licensed with Creative Commons Attribution 4.0 International license.

# Number of open analysis tools has grown exponentially

![](_page_28_Figure_1.jpeg)

![](_page_28_Picture_2.jpeg)

![](_page_28_Picture_3.jpeg)

![](_page_28_Picture_4.jpeg)

![](_page_28_Picture_5.jpeg)

1. Ampvis2 Tools for visualising amplicon sequencing data

- 2. CCREPE Compositionality Corrected by PErmutation and REnormalization
- 3. DADA2 Divisive Amplicon Denoising Algorithm
- 4. DESeq2 Differential expression analysis for sequence count data
- 5. edgeR empirical analysis of DGE in R
- 6. mare Microbiota Analysis in R Easily
- 7. Metacoder An R package for visualization and manipulation of community taxonomic diversity data
- 8. metagenomeSeq Differential abundance analysis for microbial marker-gene surveys
- 9. microbiome R package Tools for microbiome analysis in R
- 10. MINT Multivariate INTegrative method
- 11. mixDIABLO Data Integration Analysis for Biomarker discovery using Latent variable approaches for 'Omics studies
- 12. mixMC Multivariate Statistical Framework to Gain Insight into Microbial Communities
- 13. MMinte Methodology for the large-scale assessment of microbial metabolic interactions (MMinte) from 16S rDNA data
- 14. pathostat Statistical Microbiome Analysis on metagenomics results from sequencing data samples
- 15. phylofactor Phylogenetic factorization of compositional data
- 16. phylogeo Geographic analysis and visualization of microbiome data
- 17. Phyloseq Import, share, and analyze microbiome census data using R
- 18. qiimer R tools compliment qiime
- 19. RAM R for Amplicon-Sequencing-Based Microbial-Ecology
- 20. ShinyPhyloseq Web-tool with user interface for Phyloseq
- 21. SigTree Identify and Visualize Significantly Responsive Branches in a Phylogenetic Tree
- 22. SPIEC-EASI Sparse and Compositionally Robust Inference of Microbial Ecological Networks
- 23. structSSI Simultaneous and Selective Inference for Grouped or Hierarchically Structured Data
- 24. Tax4Fun Predicting functional profiles from metagenomic 16S rRNA gene data
- 25. taxize Taxonomic Information from Around the Web
- 26. labdsv Ordination and Multivariate Analysis for Ecology
- 27. Vegan R package for community ecologists
- 28. igraph Network Analysis and Visualization in R
- 29. MicrobiomeHD A standardized database of human gut microbiome studies in health and disease Case-Control
- 30. Rhea A pipeline with modular R scripts
- 31. microbiomeutilities Extending and supporting package based on microbiome and phyloseq R package
- 32. breakaway Species Richness Estimation and Modeling

### <u>A survey for 16S</u> Github.com/microsud/ Tools-Microbiome-Analysis

![](_page_29_Picture_33.jpeg)

![](_page_29_Picture_34.jpeg)

Journal of Biosciences October 2019, 44:115 | <u>Cite as</u>

### Microbiome data science

Authors

Authors and affiliations

Sudarshan A Shetty, Leo Lahti 🖂

### **Orchestrating Microbiome Analysis**

Authors: Leo Lahti [aut], Sudarshan Shetty [aut], Felix GM Ernst [aut, cre] Version: 0.98.9 Modified: 2021-04-10 Compiled: 2021-07-29 Environment: R version 4.1.0 (2021-05-18), Bioconductor 3.14 License: CC BY-NC-SA 3.0 US Copyright:

Source: https://github.com/microbiome/OMA

![](_page_30_Picture_3.jpeg)

Figure source: Moreno-Indias et al. (2021) Statistical and Machine Learning Techniques in Human Microbiome Studies: Contemporary Challenges and Solutions. Frontiers in Microbiology 12:11.

# Online tutorial (beta)

microbiome.github.io

Forthcoming courses & events: June 20-23, Oulu, Finland July 11-15, Nijmegen, The Netherlands August 8-12, Pune, India (online / hybrid?) Sep 26-30, CSC, Finland (online) Oct 5-7, Barcelona, Spain (ML4microbiome) Nov, Turku, Finland?

![](_page_31_Picture_0.jpeg)

![](_page_31_Figure_1.jpeg)