

# Multi-omic data science with R/Bioconductor

Welcome to Oulu Summer School, June 2022

2022-06-28

# Contents

<b>1</b>	<b>Overview</b>	<b>2</b>
1.1	Contents and learning goals . . . . .	2
1.2	Schedule and organizers . . . . .	2
1.3	How to apply . . . . .	3
1.4	Acknowledgments . . . . .	3
<b>2</b>	<b>Program</b>	<b>5</b>
2.1	Day 1 - Open data science . . . . .	5
2.2	Day 2 - Tabular data . . . . .	5
2.3	Day 3 - Multi-assay data . . . . .	6
2.4	Day 4 - Advanced topics . . . . .	6
<b>3</b>	<b>Project-wide Code of Conduct statement for Bioconductor</b>	<b>7</b>
3.1	Code of Conduct -Version 1.0.2 (July 27, 2021) . . . . .	7
3.2	Enforcement . . . . .	8
3.3	Reporting . . . . .	8
<b>4</b>	<b>Getting started</b>	<b>10</b>
4.1	Checklist (before the course) . . . . .	10
4.2	Support and resources . . . . .	11
4.3	Installing and loading the required R packages . . . . .	11
<b>5</b>	<b>Reproducible reporting with Rmarkdown</b>	<b>13</b>
<b>6</b>	<b>Study material</b>	<b>14</b>
6.1	Online tutorial . . . . .	14
6.2	Lecture slides . . . . .	14
6.3	Tasks . . . . .	14
6.4	Extra material on miaverse and R programming . . . . .	14

# Chapter 1

## Overview

### 1.1 Contents and learning goals

This course will teach the **basics of biomedical data analysis with R/Bioconductor**, a popular open source environment for scientific data analysis. The participants get an overview of the reproducible data analysis workflow in modern multi-omics, with a focus on recent examples from published microbiome studies. After the course you will know how to approach new tasks in biomedical data analysis by utilizing available documentation and R tools.

The teaching will follow open online documentation created by the course teachers, extending the online book *Orchestrating Microbiome Analysis* (<https://microbiome.github.io/OMA>). The openly licensed teaching material will be available online during and after the course, following national recommendations on open education.

The training material walks you through the standard steps of biomedical data analysis covering data access, exploration, analysis, visualization, reproducible reporting, and best practices in open science. We will teach generic data analytical skills that are applicable to common data analysis tasks encountered in modern omics research. The teaching format allows adaptations according to the student's learning speed.

### 1.2 Schedule and organizers

The course will be organized in a live format (Flyer)

**Venue** University of Oulu. June 20-23, 2022.

**Schedule** Contact teaching daily between 9am – 5pm, including lectures, demonstrations, hands-on sessions, and breaks. A detailed schedule is available at the course website: ([https://microbiome.github.io/course\\_2022\\_oulu](https://microbiome.github.io/course_2022_oulu)).

#### Teachers and organizers

Leo Lahti is the main teacher and Associate Professor in Data Science at the University of Turku, with specialization on biomedical data analysis. Course assistants are *Tuomas Borman* (University of Turku) is one of the main developers of the open training material covered by the course, *Jenni Hekkala*, a PhD researcher at the University of Oulu, in the group of the course coordinator Docent *Justus Reunanen*, and *Rajesh Shigdel* who has supported the writing of the course material.

The course is jointly organized by

- Health and Biosciences Doctoral Programme University of Oulu Graduate School
- Cancer & Translational Medicine Research Unit, University of Oulu
- Department of Computing, University of Turku, Finland
- Finnish IT Center for Science (CSC) supports the course with cloud computing services

## 1.3 How to apply

### Target audience

The course is primarily designed for advanced MSc and PhD students, Postdocs, and biomedical researchers who wish to learn and develop new skills in scientific programming and biomedical data analysis. Academic students and researchers from Finland and abroad are welcome and encouraged to apply. The course has limited capacity of max 20 participants, and priority will be given for local students from Oulu.

**Expected background** Some earlier experience with R or another programming language is recommended. However, this can be compensated by familiarizing with the course material in advance, if necessary. The teaching format allows adaptations according to the student's learning speed.

### Application

- Send a brief motivation letter to Jenni Hekkala [first.last@oulu.fi](mailto:first.last@oulu.fi)
- Applications sent before May 20 will be given priority

### Course fee

The course fee covers contact teaching and teaching material.

- 285 euros with registration by May 20, 2022
- 350 euros with registration after May 20, 2022
- Local students are exempted from the fee

### Accommodation

Accommodation and travel costs are not included in the registration fee. For accommodation tips, see <https://visitoulu.fi/en/arrival-overnight/>

## 1.4 Acknowledgments

**Citation** We thank all developers and contributors who have contributed open resources that supported the development of the training material. Kindly cite the course material as Tuomas Borman and Leo Lahti (2022)

**Contact** See <https://microbiome.github.io>

### License and source code

All material is released under the open CC BY-NC-SA 3.0 License and available online during and after the course, following the recommendations on open teaching materials of the national open science coordination in Finland\*\*.

The source code of this repository is reproducible and contains the Rmd files with executable code. All files can be rendered at one go by running the file `main.R`. You can check the file for details on how to clone the repository and convert it into a gitbook, although this is not necessary for the training.

- Source code (github): [miaverse teaching material](#)
- Course page (html): [miaverse teaching material](#)

# Chapter 2

## Program

The course takes place daily from 9am – 5pm (CEST), including coffee and lunch breaks.

We expect that participants will prepare for the course in advance, see section 4. Online support is available.

The material follows open online book created by the course teachers, Orchestrating Microbiome Analysis <https://microbiome.github.io/OMA>. This is R/Bioconductor framework for multi-omic data science.

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.

### 2.1 Day 1 - Open data science

#### Morning session

9-10 Coffee, Welcome & Practicalities

10-11 Lecture: Open & reproducible workflows

11-12 Demo & hands-on: Introduction to CSC RStudio notebook

12-13 Lunch break

#### Afternoon hands-on session

13-15 Demo: Data science framework

15-17 Hands-on: microbiome data summaries & exploration

17-18 Presentations & Discussion

---

### 2.2 Day 2 - Tabular data

#### Morning session

9-10 Lecture: Analysis & visualization of *tabular data*

10-12 Demo & hands-on: Univariate methods

12-13 Lunch break

**Afternoon hands-on session**

13-14 Demo: Multivariate data analysis & visualization

14-17 Hands-on: Multivariate data analysis & visualization

17-18 Presentations & Discussion

---

## 2.3 Day 3 - Multi-assay data

**Morning session**

9-10 Lecture: multi-omic data integration

10-12 Demo & hands-on: multi-assay data container

12-13 Lunch break

**Afternoon hands-on session**

13-15: Demo & hands-on: association analysis

13-17: Demo & hands-on: machine learning

17-18 Presentations & Discussion

---

## 2.4 Day 4 - Advanced topics

**Morning session**

9-10 Summary of the learning material

10-12 Demo & hands-on: custom data & advanced tools

12-13 Q & A session

**Afternoon session**

13-14 Lunch

14-16 Wrap-up

## Chapter 3

# Project-wide Code of Conduct statement for Bioconductor

[link to code of conduct](#)

(Adapted from the BioC 2020 Code of Conduct)

### 3.1 Code of Conduct -Version 1.0.2 (July 27, 2021)

The Bioconductor community values an open approach to science that promotes the

- sharing of ideas, code, software and expertise
- collaboration
- diversity and inclusivity
- a kind and welcoming environment
- community contributions

In line with these values, Bioconductor is dedicated to providing a welcoming, supportive, collegial experience free of harassment, intimidation, and bullying regardless of:

- identity: gender, gender identity and expression, sexual orientation, disability, physical appearance, ethnicity, body size, race, age, religion, language etc.
- intellectual position: approaches to data analysis, software preferences, coding style, scientific perspective, stage of career, etc.

By participating in this community, you agree not to engage in behavior contrary to these values at any Bioconductor-sponsored event (in person or virtual, including but not limited to talks, workshops, poster sessions, social activities) or electronic communication channel (including but not limited to community-bioc Slack, the support site, online forums, package review site and social media communications). Furthermore, we require all participants to have identifiable accounts in Bioconductor online forums. Accounts that do not adhere to this after request to de-anonymise may be deleted.

We do not tolerate harassment, intimidation, or bullying of community members. Sexual language and imagery are not appropriate in presentations, communications or in online venues, including chats.



Any person/s violating the Code of Conduct may be sanctioned or expelled temporarily or permanently from an electronic platform or event at the discretion of the Code of Conduct committee.

### **Examples of unacceptable harassment, intimidation, and bullying behavior**

Harassment includes, but is not limited to:

- Making comments in chats, to an audience or personally, that belittle or demean another person
- Sharing sexual images online
- Harassing photography or recording
- Sustained disruption of talks or other events
- Unwelcome sexual attention
- Advocating for, or encouraging, any of the above behavior

Intimidation and bullying include, but are not limited to:

- Aggressive or browbeating behavior
- Mocking or insulting another person's intellect, work, perspective, or question/comment
- Making reference to someone's gender, gender identity and expression, sexual orientation, disability, physical appearance, body size, race, age, religion, or other personal attributes in the context of a scientific discussion
- Deliberately making someone feel unwelcome
- Trolling behaviour (deliberately inflammatory or offensive messages)
- Sustained off-topic posts

## **3.2 Enforcement**

Anyone asked to stop harassing or intimidating behavior are expected to comply immediately.

If a person/s contravene the Code of Conduct the Code of Conduct committee retains the right to take any action that ensures a welcoming environment for all community members. This includes warning the alleged offender or temporary/permanent expulsion from the event and/or electronic platforms under Bioconductor's control.

The Code of Conduct committee may take action to redress anything designed to, or with the clear impact of, disrupting an event or electronic communication platform or making the environment hostile for any community member.

We expect everyone in the Bioconductor community to comply with the Code of Conduct when participating in Bioconductor events and online communication platforms.

## **3.3 Reporting**

If someone makes you or anyone else feel unsafe or unwelcome, please report it as soon as possible. You can make a report either anonymously or personally. All reports will be reviewed by the Code of Conduct Committee and will be kept confidential.

*Electronically* You can make an anonymous or non-anonymous report via the following link: <https://forms.gle/gEWHBWnXvZbEdFsq5>. It is a free-form text box that will be forwarded to the Code of Conduct Committee. Alternatively you can email the Code of Conduct Committee ([code-of-conduct@bioconductor.org](mailto:code-of-conduct@bioconductor.org)). If

you are uncomfortable reporting to the Code of Conduct committee as a group, you can contact any individual committee member via email or a direct message on the community-bioc Slack channel. Please include screenshots/copies of all relevant electronic conversations whenever possible (you don't need to compromise your anonymity!).

We can't follow up an anonymous report with you directly, but we will fully investigate it and take whatever action is necessary to prevent a recurrence.

**Personal Report (for any Bioconductor events: in-person or virtual)**

You can make a personal report to any member of the event Code of Conduct committee present at an event.

When taking a personal report, we will ensure you are safe and cannot be overheard. We may involve other event staff to ensure your report is managed properly. Once safe, we'll ask you to tell us about what happened. This can be upsetting, but we'll handle it as respectfully as possible, and you can bring someone to support you. You won't be asked to confront anyone, and we won't tell anyone who you are.

Our team will be happy to help you get the relevant support (e.g. help contacting hotel/venue security, local law enforcement, local support services, provide escorts, or otherwise assist you to feel safe for the duration of the event).

We value your attendance and participation at Bioconductor events and in our community.

# Chapter 4

## Getting started

### 4.1 Checklist (before the course)

#### 4.1.1 CSC Notebook

We will provide a temporary access to a cloud computing environment that readily contains the available software packages. Instructions to access the environment will be sent to the registered participants.

1. Read the instructions
2. Go to the CSC notebook frontpage
3. Login
  - a. Haka login
    - If you have a Finnish university account, you should be able to login with Haka
    - 1. Press **Login** button from the frontpage
    - 2. Press **Haka** button
    - 3. Select right organization
    - 4. Enter login information
  - b. CSC login
    - You can create a CSC account by following the instructions
    - 1. Press **Login** button from the frontpage
    - 2. Press **CSC** button
    - 3. Enter login information
  - c. Special login
    - For those who cannot login with Haka or CSC account
    - 1. Contact Tuomas by email (first.v.last@utu.fi) if you are not able to login
    - 2. We give you a guest account
    - 3. Press **Special Login** button from the frontpage (below the **Login** button)
    - 4. Enter login information (username goes to **email** slot)
4. Join workspace
  - a. Press **Join workspace** button (Top right corner)
  - b. Enter the **Join Code** (Check your email)

5. Start session
  - a. Press **ON** button
6. You can save files to **my-work** directory. They are kept stored even when the session is closed. **shared** folder is shared with all participants.

### 4.1.2 (Your own computer)

Setting up the system on your own computer is not required for the course but it can be useful for later use. The required software:

- R (version >4.1.0)
- RStudio; choose “Rstudio Desktop” to download the latest version. Optional but preferred. For further details, check the Rstudio home page.
- Install and load the required R packages (see Section 4.3)
- After a successful installation you can start with the case study examples in this training material

## 4.2 Support and resources

- We recommend to have a look at the additional reading tips and try out online material listed in Section 6.

**You can run the workflows by simply copy-pasting the examples.** For further, advanced material, you can test and modify further examples from the online book, and apply these techniques to your own data.

- Online support on installation and other matters, join us at Gitter

## 4.3 Installing and loading the required R packages

Note that the CSC/RStudio environment has readily installed setup. You may need the examples from this subsection if you are installing the environment on your own computer. If you need to add new packages, you can modify the examples below.

This section shows how to install and load all required packages into the R session, if needed. Only uninstalled packages are installed.

```
# List of packages that we need from cran and bioc
cran_pkg <- c("BiocManager", "bookdown", "dplyr", "ecodist", "ggplot2",
             "gridExtra", "kableExtra", "knitr", "scales", "vegan", "matrixStats")
bioc_pkg <- c("yulab.utils", "ggtree", "ANCOMBC", "ape", "DESeq2", "DirichletMultinomial", "mia", "miaViz")

# Get those packages that are already installed
cran_pkg_already_installed <- cran_pkg[cran_pkg %in% installed.packages()]
bioc_pkg_already_installed <- bioc_pkg[bioc_pkg %in% installed.packages()]
```

```

# Get those packages that need to be installed
cran_pkg_to_be_installed <- setdiff(cran_pkg, cran_pkg_already_installed)
bioc_pkg_to_be_installed <- setdiff(bioc_pkg, bioc_pkg_already_installed)

# Reorders bioc packages, so that mia and miaViz are first
bioc_pkg <- c(bioc_pkg[ bioc_pkg %in% c("mia", "miaViz") ],
              bioc_pkg[ !bioc_pkg %in% c("mia", "miaViz") ] )

# Combine to one vector
packages <- c(bioc_pkg, cran_pkg)
packages_to_install <- c( bioc_pkg_to_be_installed, cran_pkg_to_be_installed )

# If there are packages that need to be installed, install them
if( length(packages_to_install) ) {
  BiocManager::install(packages_to_install)
}

```

Now all required packages are installed, so let's load them into the session. Some function names occur in multiple packages. That is why *miaverse*'s packages *mia* and *miaViz* are prioritized. Packages that are loaded first have higher priority.

```

# Loading all packages into session. Returns true if package was successfully loaded.
loaded <- sapply(packages, require, character.only = TRUE)
as.data.frame(loaded)

```

##	loaded
## mia	TRUE
## miaViz	TRUE
## yulab.utils	TRUE
## ggtree	TRUE
## ANCOMBC	TRUE
## ape	TRUE
## DESeq2	TRUE
## DirichletMultinomial	TRUE
## BiocManager	TRUE
## bookdown	TRUE
## dplyr	TRUE
## ecodist	TRUE
## ggplot2	TRUE
## gridExtra	TRUE
## kableExtra	TRUE
## knitr	TRUE
## scales	TRUE
## vegan	TRUE
## matrixStats	TRUE

## Chapter 5

# Reproducible reporting with Rmarkdown

Reproducible reporting is the starting point for robust interactive data science. Perform the following tasks:

- If you are entirely new to Markdown, take this 10 minute tutorial to get introduced to the most important functions within Markdown. Then experiment with different options with Rmarkdown
- Create a Rmarkdown template in RStudio, and render it into a document (markdown, PDF, docx or other format). In case you are new to Rmarkdown Rstudio provides resources to learn about the use cases and the basics of Rmarkdown.
- Further examples and tips for Rmarkdown are available in the online tutorial to reproducible reporting by Dr. C Titus Brown.

# Chapter 6

## Study material

### 6.1 Online tutorial

The course will utilize material from the online book (beta version) Orchestrating Microbiome Analysis with R/Bioconductor (OMA).

We encourage to familiarize with the material and test examples already before the course.

### 6.2 Lecture slides

Slides (will be added).

### 6.3 Tasks

Exercises

### 6.4 Extra material on miaverse and R programming

Resources

In this course, we will analyze HintikkaXOData. In this rat study, it was analyzed whether fats and prebiotics affects the microbiome.

The data consist of 4 groups:

- High-fat diet without prebiotics
- High-fat diet with prebiotics
- Low-fat diet without prebiotics
- Low-fat diet with prebiotics

You can find the data from [here](#).

# Bibliography

Tuomas Borman and Leo Lahti (2022). *Multi-omic data science with R/Bioconductor*.