# University of Padua

---

Department of mathematics Tullio-Levi Civita

*Master Degree in Data Science*

# Metagenomics data analysis: development and application of ai methods for gut microbiome community characterization in ibd

*Supervisor*
prof. Francesco Rinaldi
University of Padua

*Co-supervisor*
Loris Bertoldi, PhD
BMR Genomics

*Master Candidate*
Veronica Pederiva
2010641

*Academic Year*

2021-2022

To all my beloved ones.

# Abstract

The gut microbiome plays a crucial role in maintaining the host homeostasis. On the contrary, a dysregulation in the gut microbial composition can seriously affect the host health leading to a condition known as Inflammatory bowel disease (IBD), one of the most prevalent forms of dysbiosis. The serious impact the microbiome regulation has on the human health has led researchers to investigate which microbial and metabolite elements constitute a healthy core microbiome. In particular, understanding which species must be protected from pathogens proliferation in order to ensure a healthy functional environment may help with the definition of targeted therapies that can be either prebiotic or probiotic-based. Differentially abundant analysis (DAA) is usually applied to investigate species and metabolic pathways that are enriched or depleted in the dysbiotic condition compared to the healthy one. However, since interactions play a major role in the microbiome regulation, an innovative approach based on community detection was proposed in this thesis to identify communities characterizing a healthy or IBD-affected microbiota. Results of this latter approach were compared to the DAA outcomes and interestingly the IBD enriched *Phascolarctobacterium succinatutens* species emerged as an IBD community leading bacteria, too. Consequently, this succinate-consumer bacteria might be studied as a potential target of new therapies. Notwithstanding, marker-based approaches as DAA are still valid to identify features that can be used for the definition of machine learning models. Indeed, the integration of data-driven models in the medical practice might provide a reliable evaluation of the IBD risk avoiding invasive procedures. In this work, a Random Forest classifier was successfully designed and trained to discriminate between healthy and IBD samples.

# Contents

# Listing of figures

# Listing of tables

# Listing of acronyms

**BH** . . . . . . . . . . . Benjamini-Hochberg correction procedure

**BMI** . . . . . . . . . . . Body Mass Index

**BMR-ITA** . . . . . . Italian Cohort Dataset identifier

**CD** . . . . . . . . . . . . Chron's Disease

**CRC** . . . . . . . . . . Colorectal Cancer

**DA** . . . . . . . . . . . . Differentially Abundant

**DAA** . . . . . . . . . . Differential Abundance Analysis

**DT** . . . . . . . . . . . . Decision Tree

**EDA** . . . . . . . . . . . Exploratory Data Analysis

**FMT** . . . . . . . . . . Fecal Microbiota Transplantation

**GN** . . . . . . . . . . . . Girvan-Newman algorithm

**HMP2** . . . . . . . . . Integrative Human Microbiome Project

**IBD** . . . . . . . . . . . Inflammatory Bowel Disease

**IBS** . . . . . . . . . . . . Irritable Bowel Syndrome

**MetaHIT** . . . . . . METAgenomics of the Human Intestinal Tract

**PCoA** . . . . . . . . . . Principal Coordinates Analysis

**RF** . . . . . . . . . . . . Random Forest

**SCFA** . . . . . . . . . . Short-Chain Fatty Acid

**SRA** . . . . . . . . . . . Sequence Read Archive

**Sar** . . . . . . . . . . . . Adjusted Rand similarity coefficient

**Sfm** . . . . . . . . . . . Fowlkes-Mallows similarity coefficient

**Sg** . . . . . . . . . . . . . Gamma similarity coefficient

**Sj** . . . . . . . . . . . . . Jaccard similarity coefficient

**Sm** . . . . . . . . . . . Minkowski similarity coefficient

**Sr** . . . . . . . . . . . . Rand similarity coefficient

**T2D** . . . . . . . . . . Type 2 Diabetes

**UC** . . . . . . . . . . . Ulcerative Colitis