

Causal Discovery on Gut Microbial Data for Disease Risk Prediction

Mariana Paco Mendivil¹ Candus Shi² Nicole Zhang³ Mentor: Biwei Huang⁴ Mentor: Jelena Bradic⁵

¹mpacomendivil@ucsd.edu ²c6shi@ucsd.edu ³nwzhang@ucsd.edu ⁴biho07@ucsd.edu ⁵jbradic@ucsd.edu

Background

- **Causal Discovery & Causal Inference:** These are a set of methods and models that attempt to causally answer scientific questions using observed data rather than RCTs.
- **Gut Microbiome:** The gut microbiome is an important indicator of human health, and extensive research is ongoing to explore its impact on human health and disease.
- **Causality in the Gut Microbiome:** Most studies report associations, but these are often insufficient to answer the scientific question of interest.

Research Questions

1. **Microbe-Microbe:** How do the microbe-microbe interaction networks between healthy and diseased participants differ?
2. **Microbe-Disease:** Which microbes have a causal relationship to disease status, and how is it quantified?
3. **Prediction:** Is it possible to predict disease status with the current composition of the dataset given causal representation learning techniques? How do they differ with the microbes learned in question 2?

Data

- **T2D:** NIH Human Microbiome Project (HMP2) dataset, filtered to healthy visits with 16S sequencing. Includes 153 insulin-sensitive (IS) and 178 insulin-resistant (IR) samples.
- **PCOS:** Meta analysis dataset from 14 different clinical studies across Asia and Europe. Includes 435 healthy controls (HC) and 513 PCOS patients.

Causal Discovery

Causal discovery attempts to recover the true causal structure of a system given observed data. One way to model this causal structure is through a directed graphical model. A widely-used general-purpose causal discovery algorithm is the Peter-Clark (PC) algorithm. It follows these key steps:

1. Start with a **complete undirected graph** (each node connected to all other nodes).
2. **Remove edges** based on statistical independence and conditional independence tests.
3. **Identify v-structures** (patterns like $X \rightarrow Y \leftarrow Z$) to infer causal directions.
4. **Apply Meek's rules** to orient additional edges while preserving v-structures.

The result is a **CPDAG (Completed Partially Directed Acyclic Graph)**, which represents a set of causal structures consistent with the observed data, also known as the Markov Equivalence Class (MEC).

Methods

1. **Filter out rare OTUs.**
2. **Feature selection and sure screening.** SparCC and graphical lasso to reduce the number of edges between pairs of microbes; logistic lasso regression to reduce the number of features that are not helpful in predicting disease.
3. **Causal discovery algorithms.** PC-stable with a max depth of 2 for microbe-microbe; CD-NOD (a variant of PC) for microbe-outcome.
4. **Causal inference.** do-calculus and logistic regression for causal effects; compare with Bayesian Inferential Regression for Differential Microbiome Analysis (BIRDMAN).
5. **Variational autoencoder.** Incorporated a latent space decomposition to distinguish between condition-dependent features. Compared baseline classification models using both the original data and VAE-reconstructed outputs.

T2D

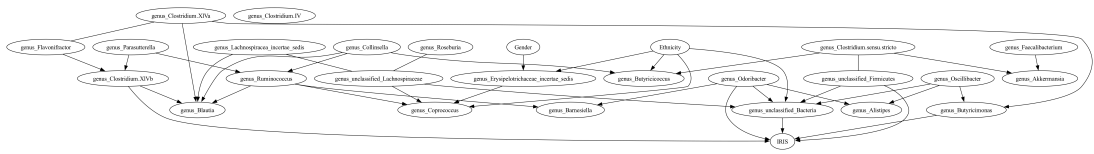


Figure 1. Microbe-Disease Network for T2D.

Model 1: $\text{logit}(\text{disease status}) \sim \text{microbes directly linked}$

Model 2: $\text{logit}(\text{disease status}) \sim \text{microbe} + \text{neighbors}(\text{microbe}) \text{ or mediators}$

BIRDMAN: Bayesian inference with $NegBinomial(\mu, \phi)$ where μ is the mean count and ϕ is the dispersion

Genus	Model 1	Model 2	BIRDMAN	Literature Agreement
<i>Butyriricomonas</i>	-2.0070*	-2.26645*	-5.19385*	Yes
<i>Clostridium XIVb</i>	1.54212*	1.80822*	2.15788*	Inconclusive
<i>Odoribacter</i>	-1.46989*	-3.055047*	-2.43796*	Yes
<i>unclassified Bacteria</i>	-0.12991*	-0.12284*	0.12409	N/A
<i>unclassified Firmicutes</i>	-0.69477*	-0.933718*	-1.47437*	N/A

Table 1. Do-Calculus Results for T2D.

Table 1. shows the log odds ratio for Model 1 and Model 2, and the mean CLR from BIRDMAN. Significant values are denoted with an asterisk (*).

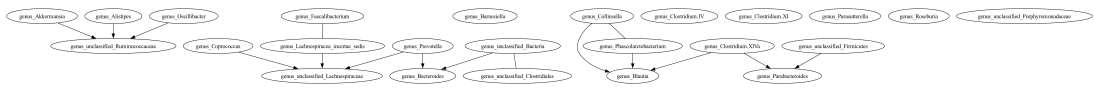


Figure 2. Microbe-Microbe Network for T2D, IS cohort.



Figure 3. Microbe-Microbe Network for T2D IR cohort.

- IR and IS cohorts share 21 microbes (Figures 2 and 3).
- IS have an additional three: *Clostridium.XI*, *Clostridium.XIVa*, and *Parasutterella*.
- IR have an additional four: *Dorea*, *Ruminococcus*, *Veillonella*, and *unclassified Ersipelotrichaceae*.

VAE

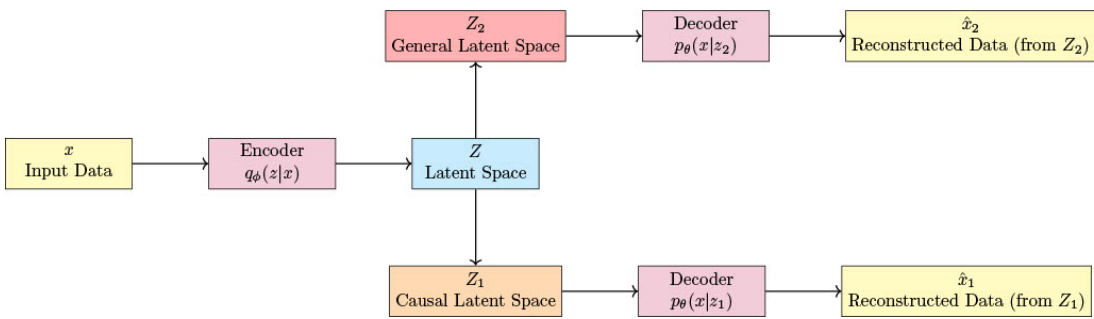


Figure 4. VAE Overview

PCOS

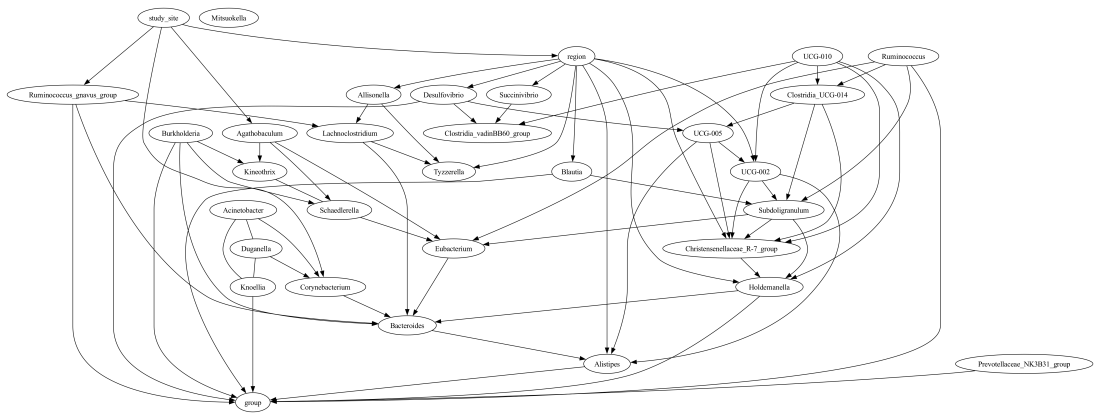


Figure 5. Microbe-Disease Network for PCOS.

Genus	Model 1	Model 2	BIRDMAN	Literature Agreement
<i>Alistipes</i>	0.13272*	0.15346*	1.28613*	Inconclusive
<i>Blautia</i>	0.07461*	0.07008*	0.82554*	No
<i>Burkholderia</i>	-7.60599	-0.48578	-10.95696*	Inconclusive
<i>Desulfovibrio</i>	-0.79283*	-1.14492*	-0.17153	No
<i>Holdemanella</i>	-0.22801*	-0.17267*	-0.13299	Yes
<i>Knoellia</i>	592.26751	1.40864	5.57650*	Inconclusive
<i>Prevotellaceae NK3B31 group</i>	-0.42407	-0.47231*	-1.76743*	Inconclusive
<i>Ruminococcus</i>	-0.14137*	-0.13490*	-0.12796	Inconclusive
<i>Ruminococcus gnavus group</i>	0.24152*	0.18259*	2.01842	Yes

Table 2. Do-Calculus Results for PCOS.

For the PCOS microbe-microbe networks, please see our website for the graphs and the specific differing microbes.

- HC and PCOS cohorts share xx microbes.
- HC have an additional xxx. Includes (bring up a microbe-disease microbe).
- PCOS have an additional xxx. Includes (bring up a microbe-disease microbe).

Conclusion & Future Work

1. **Microbe-Microbe:** Healthy and diseased participants share certain microbes, but also differ on which microbes are present and how they interact with other microbes. It is important to consider these microbes as communities of organisms rather than singular entities.
2. **Microbe-Disease:** Using CDNOD and do-calculus, we are able to quantify the effects of microbes on disease status, and they agree with microbiome-specific differential analysis methods such as BIRDMAN, and they also mostly agree with current literature.
3. **Prediction:** We find that baseline models like logistic regression perform better than those using representations from our VAE. We may require longitudinal data to leverage the mapped reconstruction of each microbe.

Our work can be improved to adjust for multiple testing and low statistical power, use different causal discovery algorithms for different data structures (e.g. longitudinal, meta-analyses, etc.) and account for compositionality and rareness in gut microbiome data.

We hope this project shows the potential of causal discovery and causal inference methods in human gut microbiome research, and can be generically applied to other diseases of interest. We would like to thank our mentors, Dr. Biwei Huang & Dr. Jelena Bradic, and Dr. Sam Degregori (Knight Lab) for guidance throughout this project.

