

# An Adaptive Association Test for Microbiome Data

Chong Wu<sup>1</sup>, Jun Chen<sup>2</sup>, Junghi<sup>1</sup> Kim and Wei Pan<sup>1</sup>

<sup>1</sup> Division of Biostatistics, School of Public Health, University of Minnesota;  
<sup>2</sup> Division of Biomedical Statistics and Informatics, Mayo Clinic

Aug. 4<sup>th</sup>, 2016



# Table of Contents

## ① Background

## ② Microbiome Based Sum of Powered Score Tests

## ③ Numerical Examples

Simulation Results

Application to a Gut Microbiome Data

# Human Microbiome Data

- Microbe: Tiny living organism, such as bacterium, fungus, or virus
- Microbiome: The genomes of human microbes; the way they interact with the human host
- Why human microbiome is important?
  - More than 10 times the number of microbes lives in the human body than cells.
  - Play an important part in our overall health.

# Human microbiome association studies

- Seen as an "extended" human genome
- Detect an association of the human microbiome diversity with a phenotype of interest
- Improve our understanding of the non-genetic component of complex traits and diseases

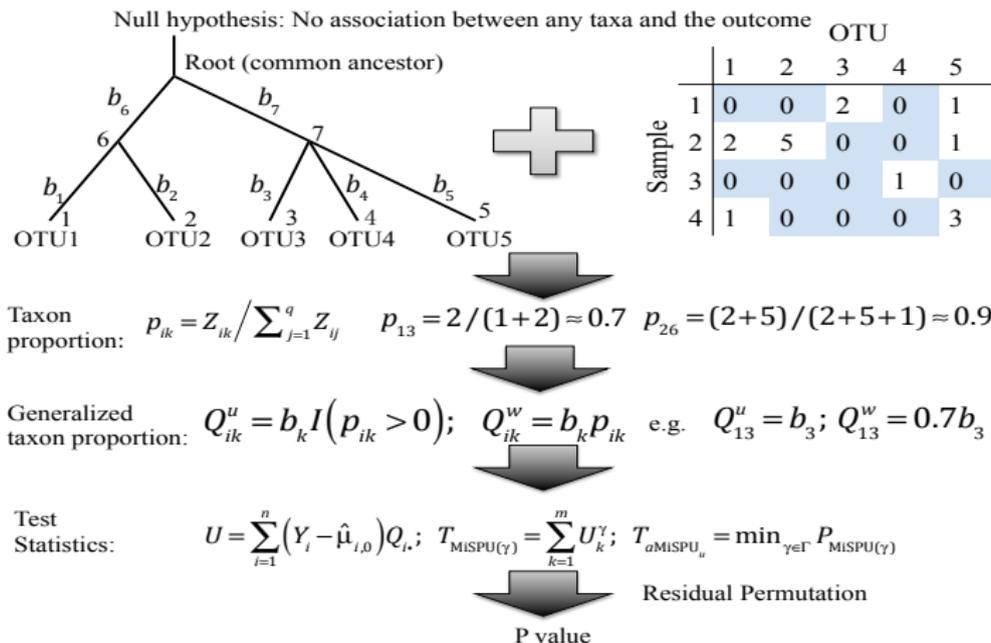
## Goal

Testing the association between the whole microbiome composition and the outcome of interest

# The Feature of The Human Microbiome Data

- Operational taxonomic units (OTUs): surrogates for biological taxa
  - High dimensional: the number of OTUs are usually much larger than the sample size
  - Overdispersion: most OTUs are rare
- Phylogenetic Tree
  - A branching diagram or "tree" showing the inferred evolutionary relationships among various biological species

# Outline



# MiSPU

- Generalized taxon proportion  $Q_{ik}$ :

$$Q_{ik}^W = b_k p_{ki}; \quad Q_{ik}^U = b_k I(p_{ki})$$

- Raw weighted UniFrac distance is exactly the same as the  $L_1$  distance of the  $Q_{ik}^W$
- For a binary outcome, we use a logistic regression model:

$$\text{Logit}[Pr(Y_i = 1)] = \beta_0 + \beta' X_i + \sum_{k=1}^m Q_{ik} \varphi_k$$

- $H_0 : \varphi = (\varphi_1, \dots, \varphi_m)' = 0$ ; that is, there is no association between any taxa and the outcome of interest

# MiSPU

- Score:

$$U = \sum_{i=1}^n (Y_i - \hat{\mu}_{i,0}) Q_i.$$

- MiSPU test statistic:

$$T_{\text{MiSPU}(\gamma)} = w' U = \sum_{k=1}^m U_k^\gamma$$

- Use permutation scheme (Pan et al., 2014) to calculate the p value

# The choice of $\gamma$

- As  $\gamma$  goes to infinity, we have

$$T_{\text{MiSPU}(\infty)} \propto \|U\|_{\infty} = \max_{k=1}^m |U_k|$$

- Intuition in the choice of  $\gamma$ :
  - the more sparse the signals, the larger  $\gamma$
  - if (most) associations in one direction, the use an odd  $\gamma$
- In practice, how to choose  $\gamma$ ?

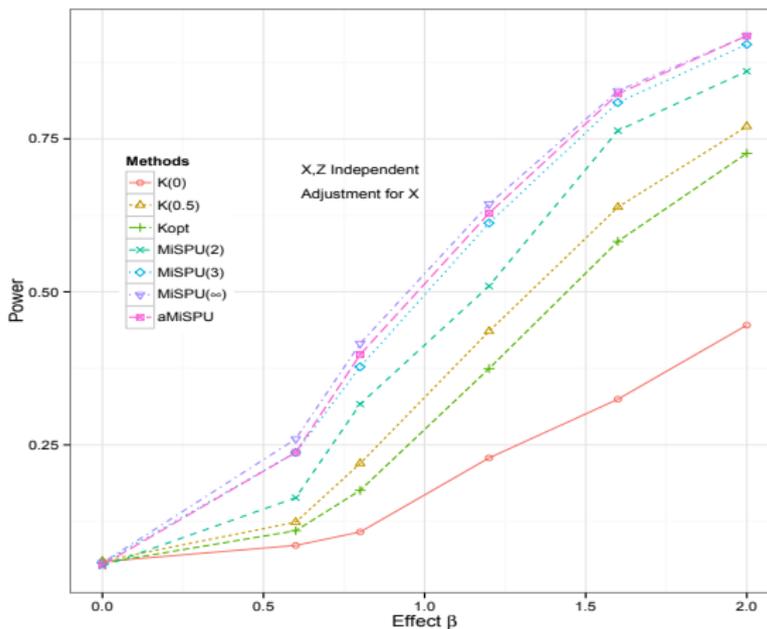
# aMiSPU

- Choose the one giving the most significant p value
- Use an adaptive test idea (Pan et al., 2014)
- aMiSPU test statistic:

$$T_{\text{aMiSPU}} = \min_{\gamma \in \Gamma} P_{\text{MiSPU}(\gamma)}.$$

- Use permutation or parametric bootstrap to estimate its p value

# Simulation Results



# Real Data: Gut Microbiome

- Diet strongly affects human health, partly by modulating gut microbiome composition
- In one cross-sectional study, 98 healthy volunteers were enrolled and habitual long-term diet information was collected using food frequency questionnaire
- Original study failed to detect the gender effect (p value **0.080**)
- Increasing evidence suggests that there is sex difference in the human gut microbiome (Bolnick et al., 2014)
- Our new method can detect it (p value **0.0058**)

# Real Data: Gut Microbiome

- A taxon in *Bacteroides* explains more than 90% relative contributions
- Top 4 taxa all come from the *Bacteroides*
- Gender status is likely associated with *Bacteroides*, but independent with other enterotypes