# Statistical analysis (II)
## -- complex experimental designs

Jessica Ewald, Postdoctoral Fellow

jessica.ewald@mcgill.ca
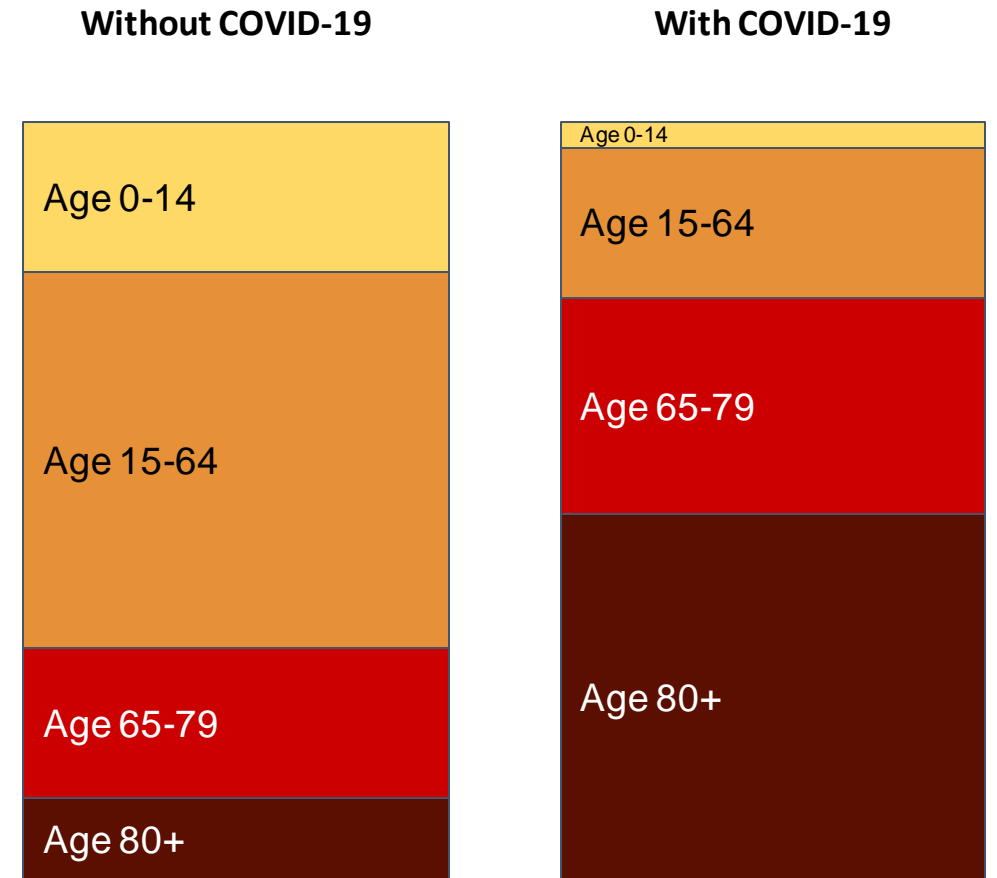
McGill University, Montreal, QC Canada

# Why complex meta-data?

- More studies are collecting samples from the "real world":
  - Epidemiology
  - Field studies

- These samples have many uncontrolled variables (covariates):
  - Biological: sex, age, disease status
  - Environmental: location, lifestyle, temperature

- Taking covariates into account can increase statistical power and reduce confounding relationships between primary variable of interest and 'omics data
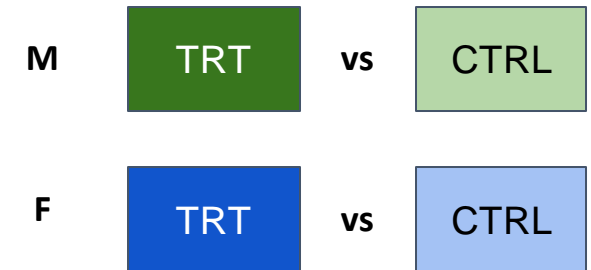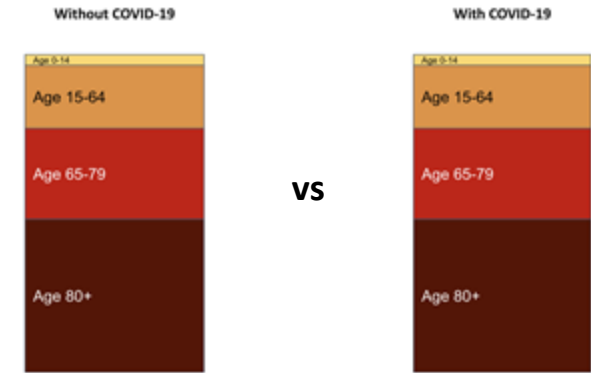
# Confounding factors: example

- Objective: identify metabolites associated with COVID-19

- Data: blood samples from subjects with and without COVID-19
  - Not a controlled experiment
  - Early in pandemic, more COVID-19 diagnoses in older subjects

- Problem: difficult to distinguish metabolites related to age vs. related to COVID-19

**Without COVID-19**

Age 0-14

Age 15-64

Age 65-79

Age 80+

**With COVID-19**

Age 0-14

Age 15-64

Age 65-79

Age 80+

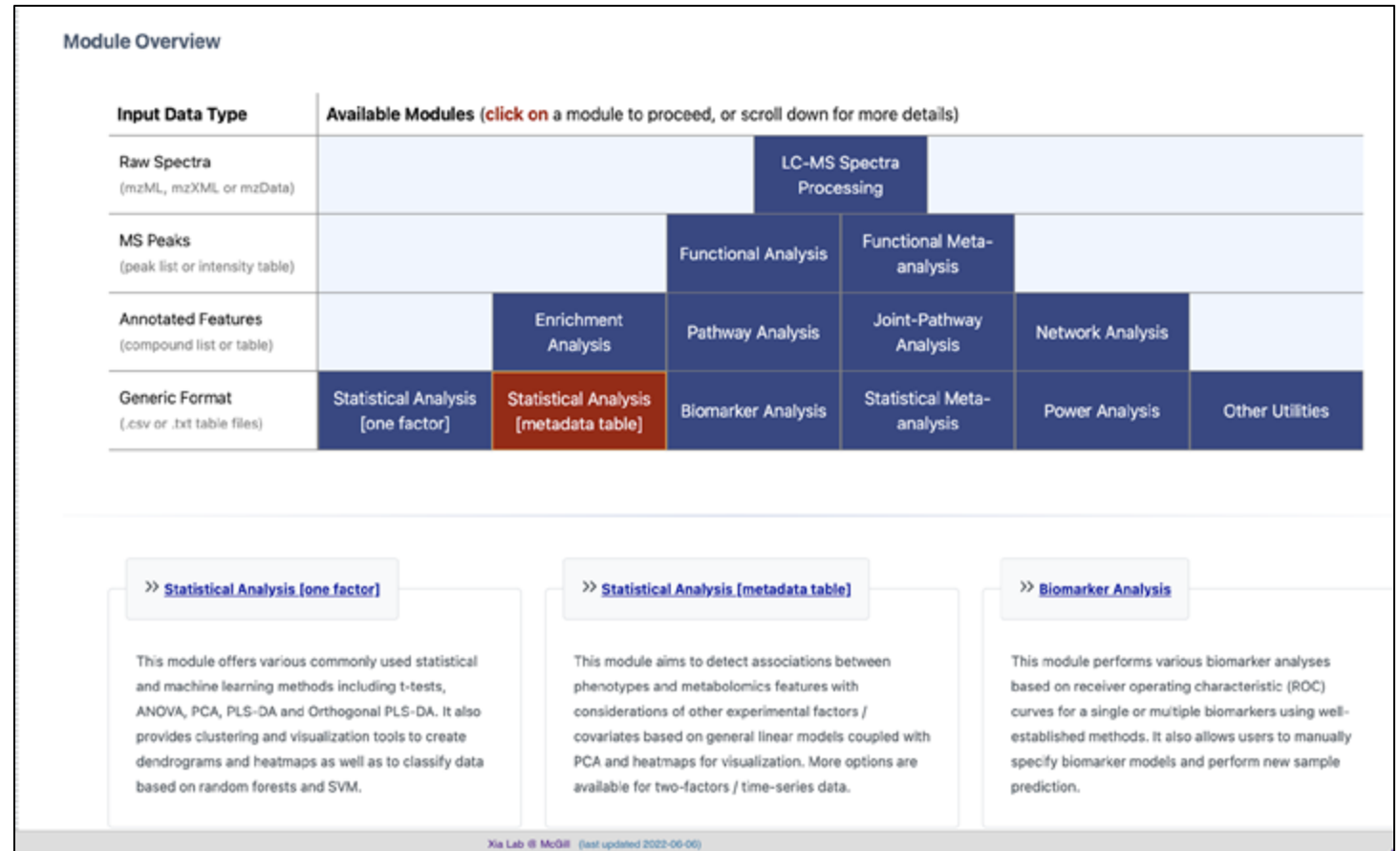# Strategies for dealing with meta-data

- Experimental design:
  - **Control** group that 'matches' your group of interest
- Data analysis (often limited by sample size):
  - **Stratification** - for factors with few classes, split up the data and analyze separately
  - **Take into account** - statistical Analysis [metadata table] module in MetaboAnalyst



Without COVID-19    With COVID-19

Age 0-14
Age 15-64
Age 65-79
Age 80+

vs

M | TRT | vs | CTRL
F | TRT | vs | CTRL

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + \ldots + b_n * x_n$$

# Complex meta-data with MetaboAnalyst

- We will cover how to:
  - Format meta-data
  - Analyze relationships between meta-data
  - Covariate adjustment with linear model
  - Supervised analysis with Random Forest

# Formatting meta-data

- Essential for downstream analysis
- Categorize as "Categorical" or "Continuous"
  - For categorical, must have at least 2 groups with at least 3 replicates each
- Missing values will be filtered out during analysis
- First meta-data column will be considered primary variable by default

| Sample | TCE_Exp_Category | TCE_Exp_Conc | Age | Sex | Smoking_Status | Alcohol_Use | BMI | Batch |
|--------|------------------|--------------|-----|-----|----------------|-------------|-----|-------|
| X1014 | Low | 0.025 | 28 | Male | Yes | Yes | 20.3 | 10 |
| X1049 | Low | 0.025 | 34 | Female | No | No | 33.7 | 1 |
| X1068 | Low | 0.025 | 30 | Male | Yes | No | 25.6 | 7 |
| X1070 | Low | 0.025 | 42 | Male | Yes | No | 21.6 | 1 |
| X1071 | Low | 0.025 | 41 | Female | No | No | 20.7 | 4 |
| X1073 | Low | 0.025 | 22 | Male | No | Yes | 20.6 | 4 |
| X1074 | Low | 0.025 | 26 | Female | No | No | 18.7 | 13 |
| X1075 | Low | 0.025 | 26 | Male | Yes | Yes | 23.1 | 12 |
| X1076 | Low | 0.025 | 16 | Male | Yes | Yes | 18.4 | 11 |
| X1078 | Low | 0.025 | 32 | Female | No | No | 19.2 | 7 |
| X1079 | Low | 0.025 | 22 | Male | Yes | No | 20.3 | 9 |
| X1080 | Low | 0.025 | 25 | Female | No | No | 21.3 | 7 |
| X1089 | Low | 0.025 | 30 | Male | No | No | 22.8 | 3 |
| X1090 | Low | 0.025 | 33 | Male | No | No | 25.5 | 10 |
| X1091 | Low | 0.025 | 27 | Male | No | Yes | 18.8 | 12 |
| X1092 | Low | 0.025 | 29 | Male | No | No | 23.6 | 2 |
| X1094 | Low | 0.025 | 35 | Male | Yes | No | 24 | 5 |
| X1095 | Low | 0.025 | 33 | Male | No | No | 21 | 12 |
| X1097 | Low | 0.025 | 27 | Male | No | No | 18.4 | 5 |
| X1098 | Low | 0.025 | 32 | Male | No | Yes | 23.2 | 6 |
| X1099 | Low | 0.025 | 20 | Male | No | No | 16.9 | 8 |
| X1100 | Low | 0.025 | 23 | Male | Yes | Yes | 19 | 12 |
| X1101 | Low | 0.025 | 25 | Male | Yes | Yes | 21.1 | 3 |
| X1106 | Low | 0.025 | 40 | Male | No | No | 30.5 | 2 |
| X1110 | Low | 0.025 | 18 | Male | No | Yes | 21 | 2 |
| X1112 | Low | 0.025 | 28 | Male | No | Yes | 21.3 | 5 |

| Name | Status | Type | Edit | Remove |
|------|--------|------|------|--------|
| TCE_Exp_Category | OK | Categorical ⌄ | Edit | 🗑 |
| TCE_Exp_Conc | OK | Categorical / Continuous | Edit | 🗑 |
| Age | OK | Continuous ⌄ | Edit | 🗑 |
| Sex | OK | Categorical ⌄ | Edit | 🗑 |
| Smoking_Status | OK | Categorical ⌄ | Edit | 🗑 |
| Alcohol_Use | OK | Categorical ⌄ | Edit | 🗑 |
| BMI | OK | Continuous ⌄ | Edit | 🗑 |
| Batch | OK | Categorical ⌄ | Edit | 🗑 |

# Method overview



**Overview to understand structure of data & metadata**

**Simple yet effective univariate statistical analysis**

**Advanced multivariate statistics & machine learning**

**Data and Metadata Overview**

Metadata Visualization

Users can explore the metadata patterns and correlations through intuitive graphics. It is very useful for users to identify highly dependent metadata and quickly assess the overall patterns of the metadata.

Interactive PCA Visualization

Users can visualize data using different colors or shapes based on selected metadata in an 2D and 3D (interactive) PCA plots. It is very useful to detect overall patterns of data with regard to different metadata.

Hierarchical Clustering and Heatmap Visualization

This method displays data and metadata in the form of colored cells. It provides direct visualization of feature abundances across different samples and metadata.

**Univariate Analysis**

Linear Models with Covariate Adjustment

This approach uses linear models (limma or lm) to perform significance testing with covariate adjustments. Users can choose different metadata to be included in the analysis.

Correlation and Partial Correlation Analysis

This approach allows users to explore the correlations or partial correlations (with covariate adjustments) between metabolomics features and different metadata of interest.

Two-way ANOVA (ANOVA2)

This approach provides classical two-way ANOVA based on the two factors selected by users. For time-series data, users should choose within-subjects ANOVA.

**Multivariate Analysis**

ANOVA Simultaneous Component Analysis (ASCA)

This approach is designed to identify major patterns with regard to the two given factors and their interaction. The implementation was based on the algorithm described by AK Smildle, et al. with additional improvements on feature selection and model validation.

Multivariate Empirical Bayes Analysis of Variance (MEBA) for Time Series

This approach is designed to compare temporal profiles across different biological conditions. It is based on the timecourse method described by YC Tai, et al.
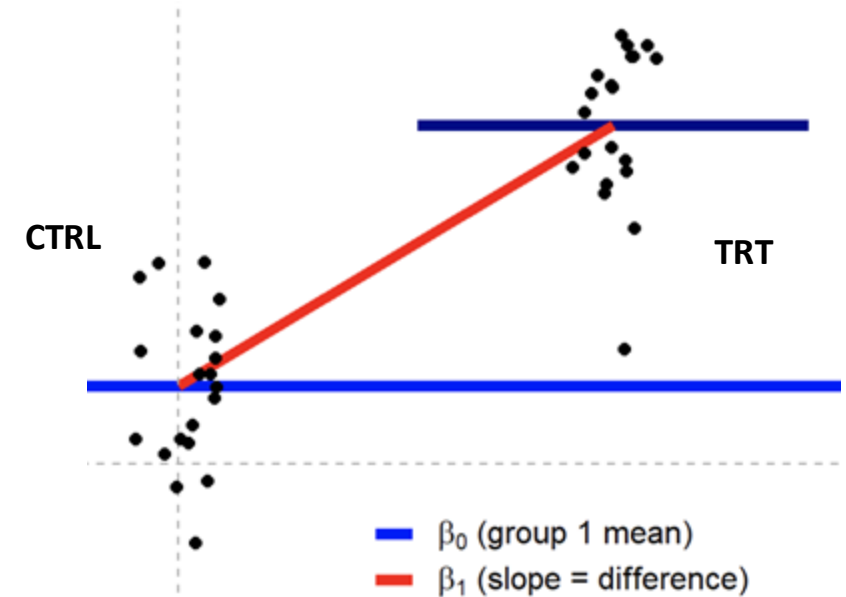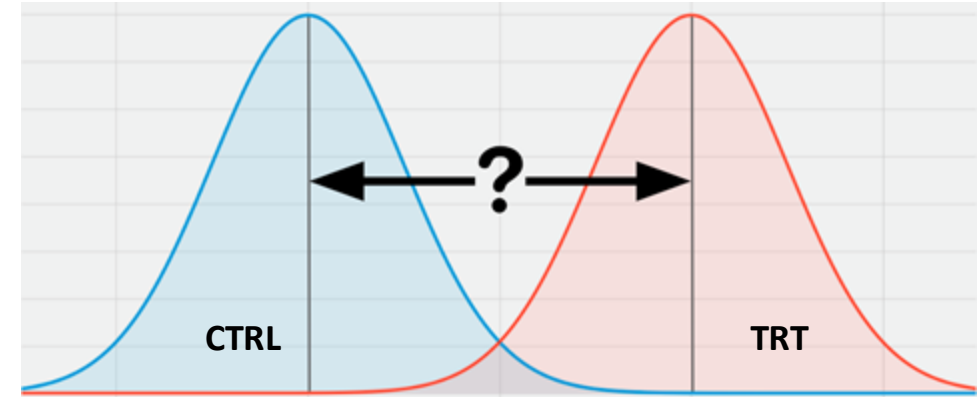
**Supervised Classification**

Random Forest

This machine learning approach is designed to perform classification and feature selection analysis. Users can also test contribution of meta-data to class prediction.

# T-test vs. linear regression



- You can do t-test with linear regression:
- $y = B_0 + B_1 * x$
  - y: level of metabolite A
  - x: variable of interest



— $\beta_0$ (group 1 mean)
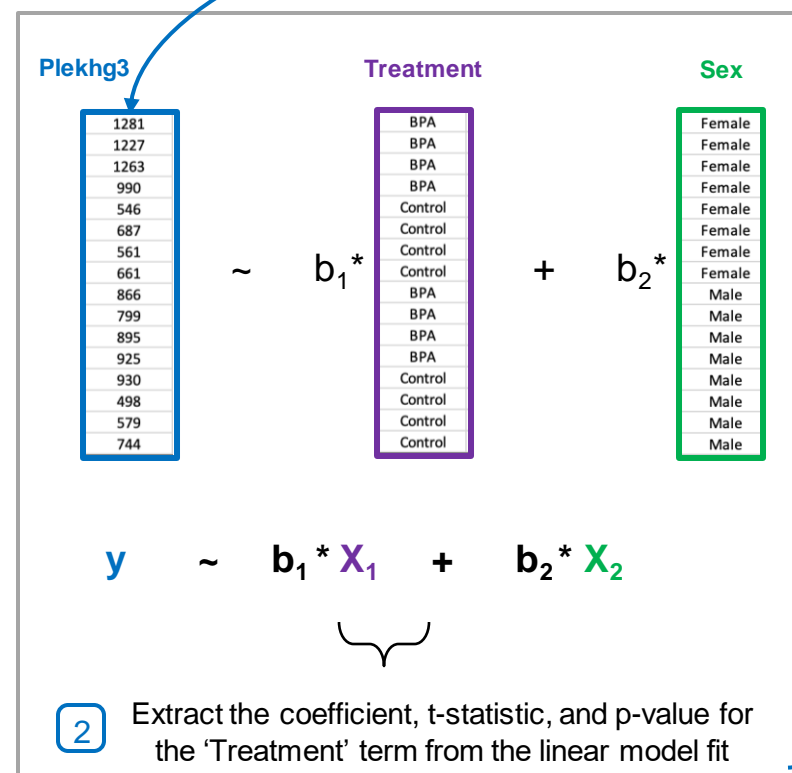— $\beta_1$ (slope = difference)

# More complex experimental design

- Treatment (Control/BPA)
- Sex (male/female)
- Features associated with *treatment* while considering *sex*
- Treatment = primary
- Sex = covariate
- Can include many covariates
- Primary/covariate can be continuous or discrete

# Interpretation of coefficient results

Multiple linear regression example:
- $y = B_0 + B_1 * x_{treatment} + B_2 * x_{sex}$

- By including $B_2 * x_{sex}$ in the model, we account for effects of sex

- Extract $B_1$ from the model:
  - $B_1$ value = magnitude & direction of relationship between metabolite 'y' and $x_{treatment}$
  - $B_1$ p-value = statistical significance of relationship



-log10(p-value): no covariate adjustment

Impact of controlling for sex

**Rpl39**

Female    Male

**S100a1**

Male

Female

| **Linear Model** | **Rpl39 Treatment P-value** | **S100a1 Treatment P-value** |
|---|---|---|
| $Y \sim$ **Treatment** | $3.7 \times 10^{-12}$ | $2.9 \times 10^{-6}$ |
| $Y \sim$ **Treatment** $+$ Sex | $9.9 \times 10^{-12}$ | $2.1 \times 10^{-11}$ |

# Functional Analysis

**Feature Details Table**

Click a feature name to edit its name and then click the next column to save the change. Click the view link to visualize a graphical summary of the distribution. The bar plots on the left show the original values (mean +/- SD). The box and whisker plots on the right summarize the normalized values. Note, positive infinite numbers are represented as 999999, and negative infinite numbers -999999.

To update a name suitable for graphical display, **click the name** to edit and then click the next column to save

⬇ Download

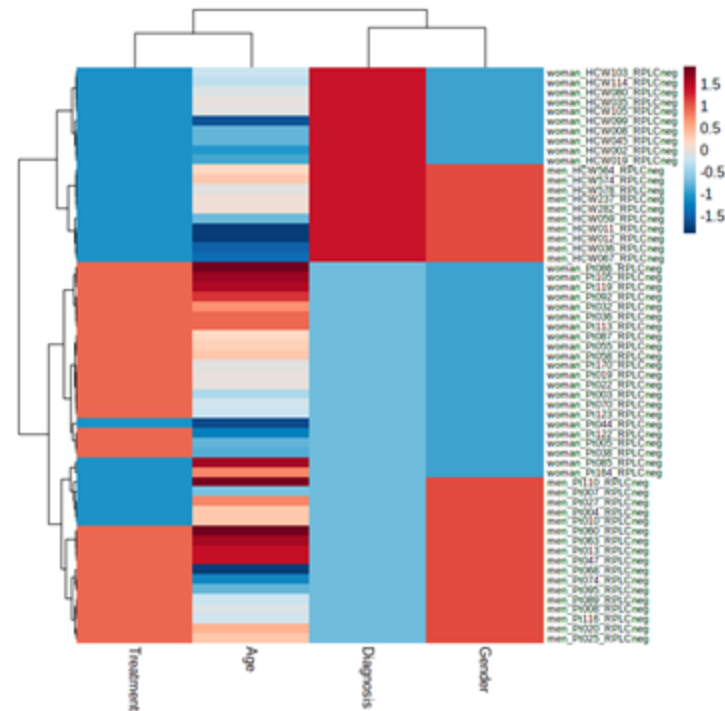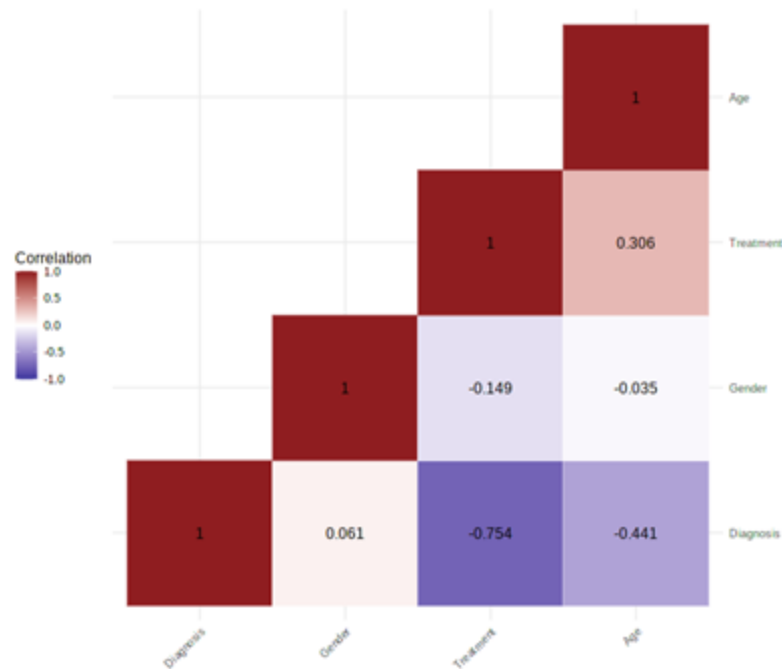| Name ↑↓ | | logFC ↑↓ | AveExpr ↑↓ | t ↑↓ | P.Value ↑↓ | adj.P.Val ↑↓ | B ↑↓ | |
|---|---|---|---|---|---|---|---|---|
| 319.227@344.95 | | 0.73798 | 0.30435 | 6.9298 | 4.2134E-9 | 8.6544E-6 | 10.605 | View |
| 347.2582@379.62 | | 0.59923 | -0.76139 | 6.6574 | 1.1944E-8 | 1.2266E-5 | 9.6281 | View |
| 173.9855@96.75 | | -1.793 | -2.1814 | -6.2058 | 6.6605E-8 | 4.5602E-5 | 8.0165 | View |
| 343.2271@340.17 | | 0.50703 | -0.45557 | 5.9527 | 1.7321E-7 | 8.8942E-5 | 7.1202 | View |
| 357.1785@488.19 | | 0.25935 | -0.70776 | 5.5101 | 9.0415E-7 | 3.2343E-4 | 5.5712 | View |
| 335.2944@537.41 | | 0.39263 | -0.40041 | 5.426 | 1.2334E-6 | 3.2343E-4 | 5.2803 | View |
| 309.2794@532.6 | | 0.34348 | 1.3126 | 5.4032 | 1.3415E-6 | 3.2343E-4 | 5.2017 | View |
| 310.2826@532.64 | | 0.34901 | 0.65576 | 5.3926 | 1.3946E-6 | 3.2343E-4 | 5.1653 | View |
| 303.2328@438.79 | | 0.29848 | 1.6434 | 5.3883 | 1.4172E-6 | 3.2343E-4 | 5.1503 | View |
| 385.2345@438.71 | | 0.26685 | -0.49809 | 5.3122 | 1.8737E-6 | 3.8486E-4 | 4.8889 | View |
| 304.2362@438.81 | | 0.30433 | 0.99595 | 5.2773 | 2.129E-6 | 3.9629E-4 | 4.7693 | View |

Use for functional analysis

# Important considerations

- Adding more variables into the model decreases statistical power
- Be strategic:
  - Do not include variables that have no impact on your data
  - Do not include correlated metadata
- Assess this with PCA, linear models, and metadata overview
  - We will go over this in the demo

# Relationships between metadata

- Must understand relationships between predictors
  - Know which to include in model (sample size)
  - Guide interpretation of the results
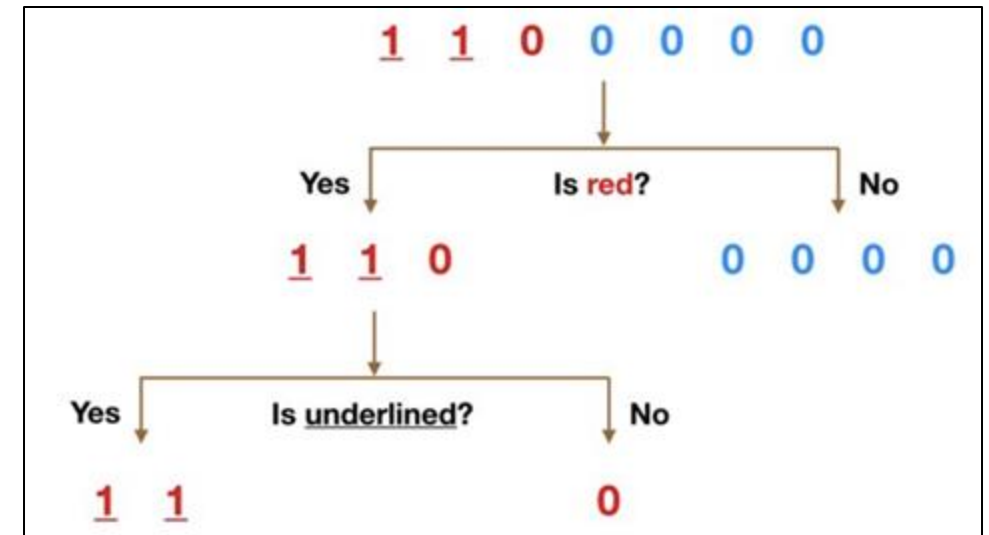
# Non-linear Strategies

- All methods so far rely on linear relationships
  - Very stable, robust against overfitting
  - Only option with small sample size
- Many non-linear relationships in complex data
- Random forest is one approach
  - Uses decision trees

**Example Decision Tree**

Predict whether digit is a "1" or a "0"
Two variables:
- Color (red, blue)
- Underlined (yes, no)

# Random Forest Classification

- In Random Forest, we build 1000s of trees
- For each tree, randomly select:
  - Random subset of metabolites + metadata
  - Random subset of samples
- Then, build a tree
- Use the rules from each tree to predict the class label for all samples (ie. COVID, healthy)
- Final prediction: majority wins

Example forest of trees
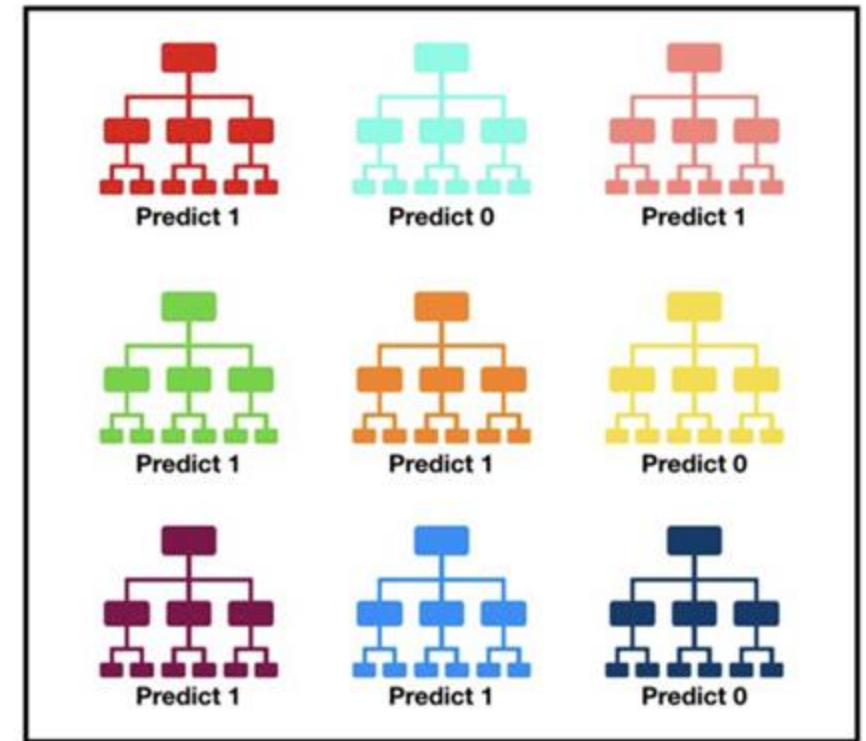


Predict 1   Predict 0   Predict 1
Predict 1   Predict 1   Predict 0
Predict 1   Predict 1   Predict 0

Tally: Six 1s and Three 0s
**Prediction: 1**

# Example Classification Results

**Primary metadata:** Diagnosis ⌄

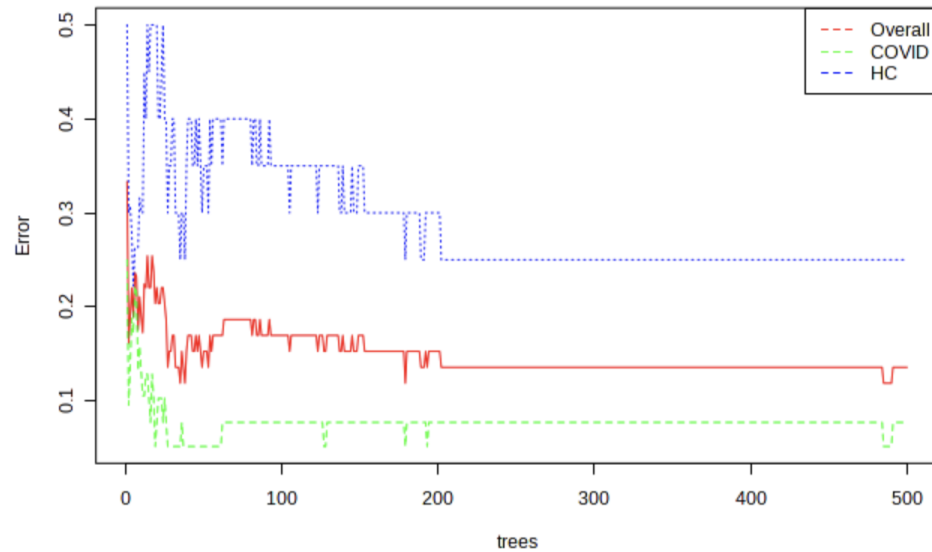**Choose metadata for predictors:** Gender ✕ Age ✕ ⌄   Update

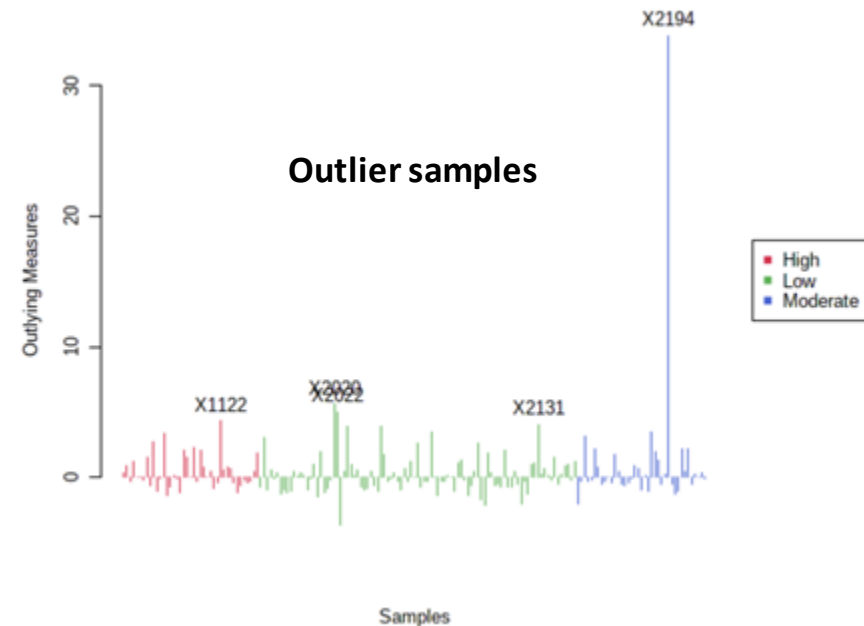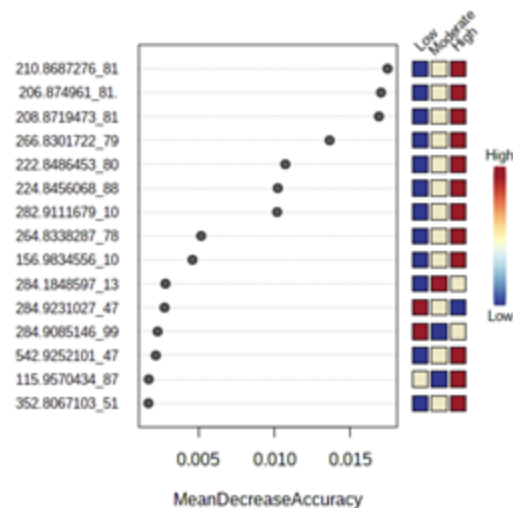**Randomness:** On ⌄



**Random Forest classification**

Legend:
- Overall
- COVID
- HC

X-axis: trees
Y-axis: Error

**The OOB error is 0.136**

|       | COVID | HC | class.error |
|-------|-------|-----|-------------|
| COVID | 36    | 3   | 0.0769      |
| HC    | 5     | 15  | 0.25        |

# Using Random Forest in MetaboAnalyst

➢ Using model in real life requires extensive validation & careful design

➢ Most MetaboAnalyst users: use as form of exploratory statistics

➢ Understand which variables have high predictive power

  ➢ Var. Importance tab

➢ Identify potential outlier observations

  ➢ "Outlier Detection" tab

**Most important variables**

**Outlier samples**

# Hands-On Demo