

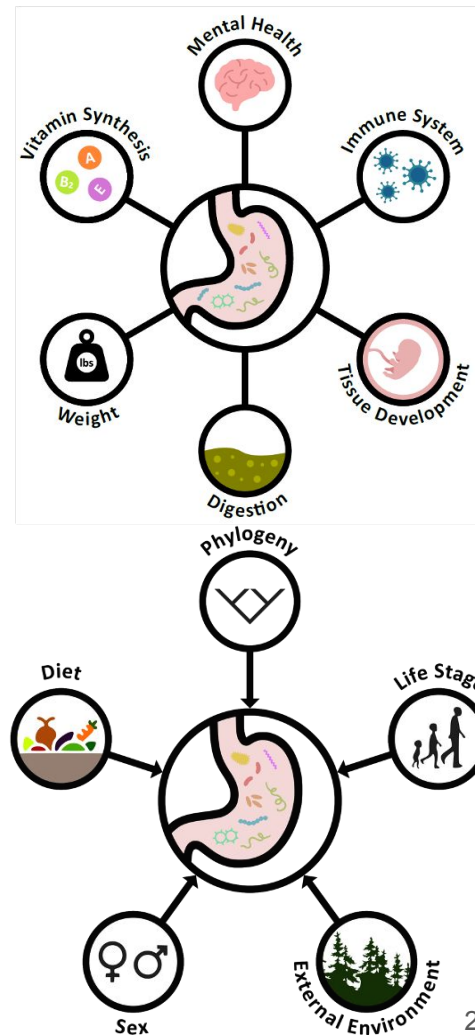
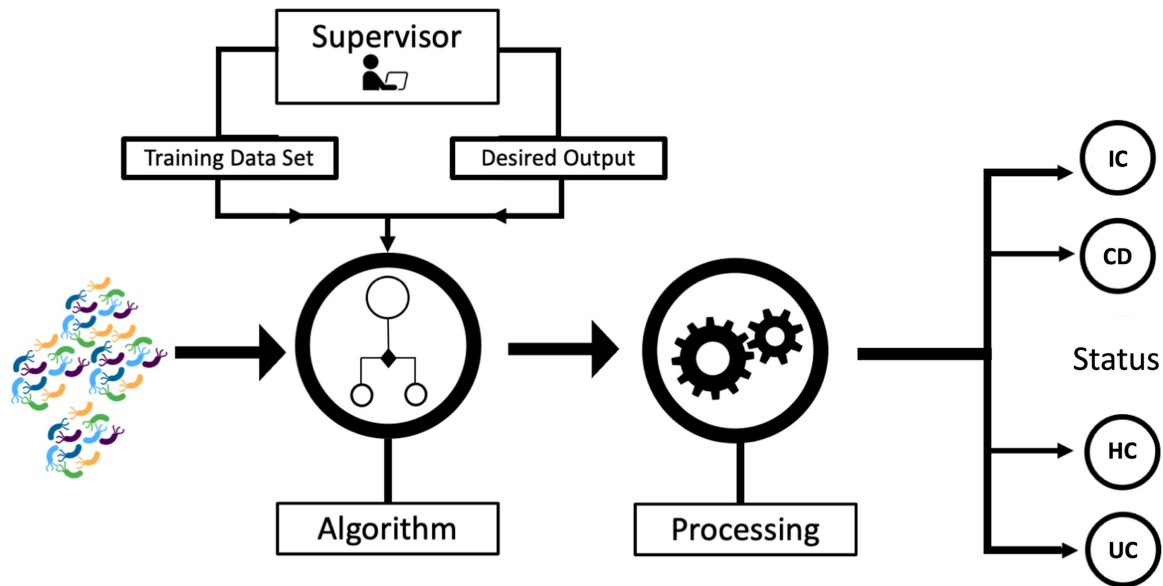
# Supervised Machine Learning for Microbiome Data

Binnan Yu, Lillian Tatka, Kristina Herman, David Lee, Sierra Gillman

# Introduction


## *The problem being addressed*

Increasing accessibility to machine learning for microbiome research




# The Data

## 4 tabundance ables of Operational Taxonomic Units (OTUs) and their metadata files




**The Treatment-Naive Microbiome in New-Onset Crohn's Disease**


Dirk Gevers,<sup>1</sup> Shihua Kugathasan,<sup>2,3,4</sup> Lee A. Denson,<sup>5,6</sup> Yoshiki Vázquez-Baeza,<sup>7</sup> Will Van Treuren,<sup>1</sup> Boqin Ren,<sup>8</sup> Emma Schwager,<sup>9</sup> Dan Knights,<sup>1,10</sup> De Jin Kong,<sup>1</sup> Marco Yassani,<sup>1</sup> Kaceli C. Morgan,<sup>1</sup> Aleksandar D. Kostic,<sup>1</sup> Chongpei Luo,<sup>1</sup> Andrew Gonzalez,<sup>1</sup> Daniel McDonald,<sup>1</sup> Yael Hershman,<sup>1</sup> Thomas Walters,<sup>1</sup> Susan Baker,<sup>1</sup> Joel Rhee,<sup>11</sup> Michael Stephens,<sup>12</sup> Malen Heyman,<sup>13</sup> James Markowitz,<sup>14</sup> Robert Baldassano,<sup>15</sup> Anne Griffiths,<sup>16</sup> Francisco Sylvester,<sup>17</sup> David Mack,<sup>18</sup> Bianca Kim,<sup>19</sup> Wallace Crowell,<sup>20</sup> Jeffrey Ayres,<sup>21</sup> Curtis Huttenhower,<sup>1</sup> Rob Knight,<sup>1,2,3</sup> and Harnik J. Xavier<sup>1,2,3\*</sup>



**BRIEF COMMUNICATION**


<https://doi.org/10.1038/s41592-018-0141-9>






**Multi-omics of the gut microbial ecosystem in inflammatory bowel diseases**

Jason Lloyd-Paet,<sup>1</sup> Cesar Aron,<sup>2</sup> Ashwin N. Ananthakrishnan,<sup>3</sup> Melinda Schrimmer,<sup>4</sup> Julian Anhe-Pacheco,<sup>5</sup> Tiffany W. Peat,<sup>6</sup> Elizabeth Andrews,<sup>7</sup> Nadine I. Ajami,<sup>8</sup> Kevin S. Busham,<sup>9</sup> Colin J. Brislawn,<sup>10</sup> David Gower,<sup>11</sup> Holly Gourtirey,<sup>12</sup> Antonio Gonzalez,<sup>13</sup> Thomas C. Crowder,<sup>14</sup> A. Bradley Karp,<sup>15</sup> Katherine Lake,<sup>16</sup> Camille I. Leclercq,<sup>17</sup> James Markowitz,<sup>18</sup> Danyal R. Plicht,<sup>19</sup> Mahadev Prasad,<sup>20</sup> Ghodsiyah Rahnavard,<sup>21</sup> Jonny Saik,<sup>22</sup> Dmitry Shugart,<sup>23</sup> Yoshiki Vázquez-Baeza,<sup>24</sup> Richard A. White III,<sup>25</sup> William K. Hargrett-Anderson,<sup>26</sup> Jonathan Bramer,<sup>27</sup> Lee A. Denson,<sup>28</sup> James K. Janssen,<sup>29</sup> Rob Knight,<sup>30,31</sup> Shihua Kugathasan,<sup>32</sup> Dennis P. B. McGovern,<sup>33</sup> Joseph F. Petrosino,<sup>34</sup> Thaddeus S. Stappenbeck,<sup>35</sup> Harland S. Winter,<sup>36</sup> Clary B. Chisholm,<sup>37</sup> Eric A. Franzosa,<sup>38</sup> Hana Vlamacki,<sup>39</sup> Harnik J. Xavier,<sup>40,41</sup> and Curtis Huttenhower<sup>42,43\*</sup>



**QIIME2**



**ARTICLES**

PUBLISHED: 13 FEBRUARY 2017 | VOLUME: 2 | ARTICLE NUMBER: 17004

**QIITA: rapid, web-enabled microbiome meta-analysis**

Antonio Gonzalez<sup>1,12</sup>, Jose A. Navas-Molina<sup>1,2,10,12</sup>, Tomasz Kosciolk<sup>1</sup>, Daniel McDonald<sup>1</sup>, Yoshiki Vázquez-Baeza<sup>1</sup>, Gail Ackermann<sup>1</sup>, Jeff DeReus<sup>1</sup>, Stefan Janssen<sup>1</sup>, Austin D. Swafford<sup>3</sup>, Stephanie B. Orchanian<sup>3</sup>, Jon G. Sanders<sup>1</sup>, Joshua Shorenstein<sup>1,11</sup>, Hannes Holste<sup>1,2</sup>, Semar Petrus<sup>4</sup>, Adam Robbins-Pianka<sup>5</sup>, Colin J. Brislawn<sup>6</sup>, Mingxun Wang<sup>7</sup>, Jai Ram Rideout<sup>8</sup>, Evan Bolyen<sup>8</sup>, Matthew Dillon<sup>8</sup>, J. Gregory Caporaso<sup>8,9</sup>, Pieter C. Dorrestein<sup>1,3,7</sup> and Rob Knight<sup>1,2,3\*</sup>

**Dynamics of the human gut microbiome in inflammatory bowel disease**

Jonas Halfvarson<sup>1</sup>, Colin J. Brislawn<sup>2</sup>, Regina Lamendella<sup>3</sup>, Yoshiki Vázquez-Baeza<sup>4</sup>, William A. Walters<sup>5</sup>, Lisa M. Bramer<sup>6</sup>, Mauro D'Amato<sup>7,8</sup>, Ferdinando Bonfiglio<sup>9</sup>, Daniel McDonald<sup>10</sup>, Antonio Gonzalez<sup>11</sup>, Erin E. McClure<sup>12</sup>, Mitchell F. Dunkleberger<sup>3</sup>, Rob Knight<sup>4,10,11</sup> and Janet K. Jansson<sup>2\*</sup>

# The Users



## The Ecologist: Loves R

Is interested in seeing if habitat quality can be predicted based on fecal microbiome community composition in American marten (*Martes americana*). They want to determine if they can classify individuals from primary or disturbed habitat as this could be a powerful tool for conservation and management!



## The Medical Clinician: Can google

Wants to be able to determine if suspected patients have inflammatory bowel disease so they can begin effective treatment before their health deteriorate, and will use BioME to classify people!



## The Microbiologist: bioinformatic buff

Is here for the preliminary results/confirmation. The lab intern might have mislabeled some (~100) of the samples. Does say that a *Turdus turdis* or just *Turtle's turds*? They don't want to have to throw out all those samples..





# The Users



## The Ecologist: Loves R

The ecologist collects fecal samples from wild martens and puts the data into BioME where he sees that support vector classification is able to sort out the samples with high accuracy.



## The Medical Clinician: Can google

The clinician gives BioME data from healthy patients and patients with inflammatory bowel disease. She then inputs the patient's data for BioME to classify as healthy or diseased.

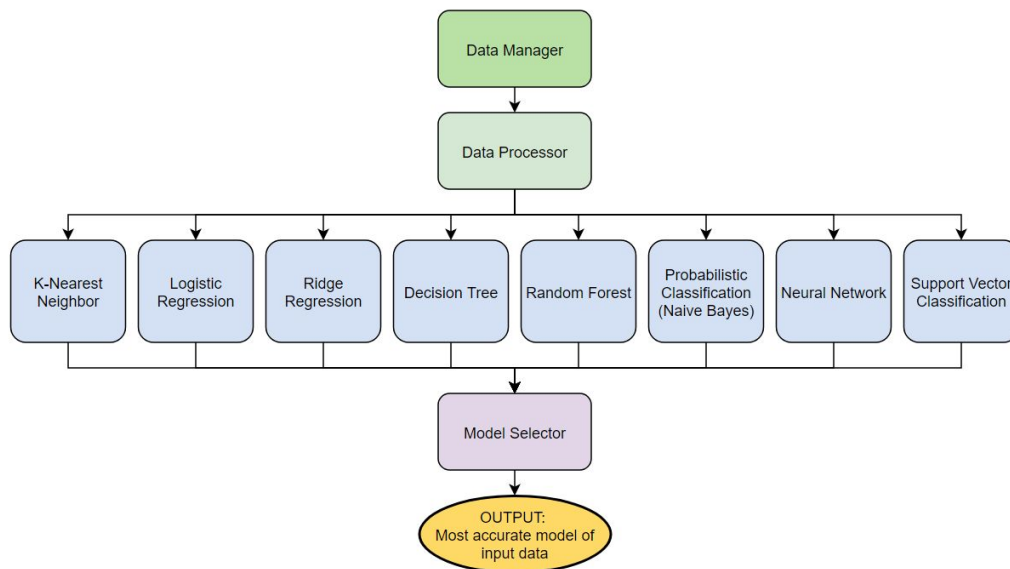


## The Microbiologist: bioinformatic buff

The microbiologist gives sample data to BioME to verify the identity of mislabeled samples. Later, he gives his verified data to BioME to compare models.

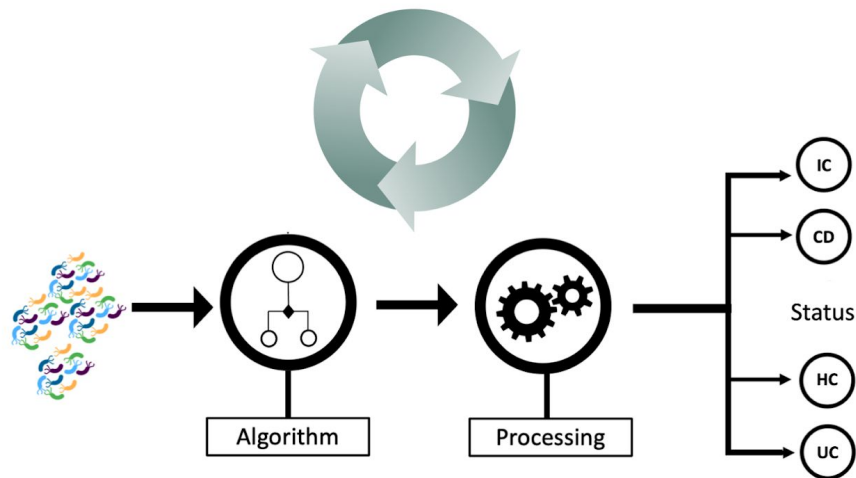
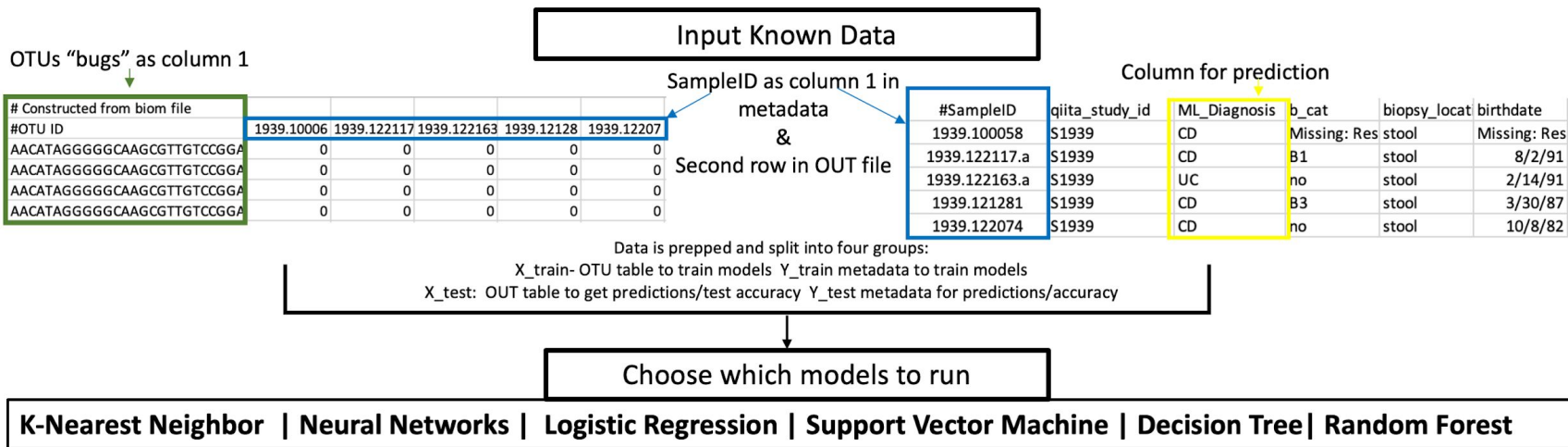


# Functional Specifications



BioME/

- | - README.md
- | - biome/
  - | - \_\_init\_\_.py
  - | - prep\_split\_data.py
  - | - select\_model.py
  - | - train\_mlp.py
  - | - knn.py
  - | - dtree.py
  - | - logistic.py
  - | - ridge.py
  - | - random.py
  - | - SVC.py
  - | - naive\_bayes.py
  - | - scripts/
    - | - biome\_run.py
  - | - tests/
    - | - ...
  - | - Data/
    - | - bug\_OTU\_rel.tsv
    - | - bug\_OTU\_raw.tsv
    - | - FecesMeta.txt
    - | - query\_point.tsv
- | - docs/
  - | - FunctionalSpec.md
  - | - ComponentSpec.md
  - | - Technology Review Presentation
  - | - images/
    - | - ...
- | - setup.py
- | - LICENSE
- | - BioME\_environment.yml



[illegible]

Please enter the relative path to the categorical data: biome/Data/FecesMeta.txt

What models would you like to test?  
See README.md for abbreviations. Type all if all models should be tested: mlp1,mlp3,dtree,knn,gnb

The best performing model is: MLP (single hidden layer)

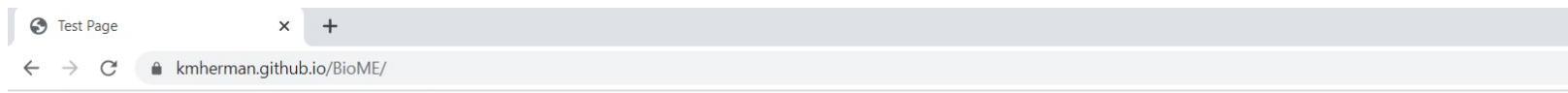
Please enter the path to the data that you would like to make a prediction for: biome/Data/query point.tsv

Would you like to make another prediction? no



# Future Work

Create a webapp that allows user to select a model they want to run

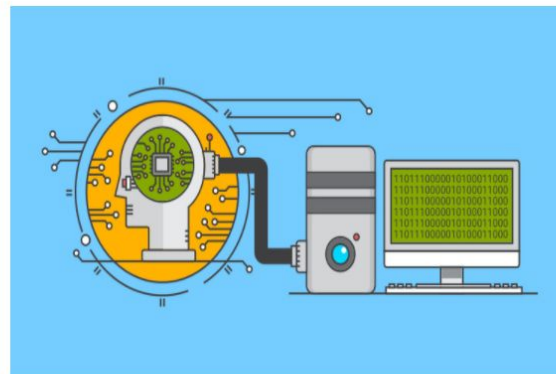
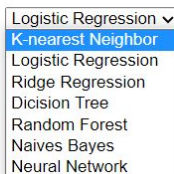


## Why Our Solution is Trustworthy?

### 1. Algorithm

Our product is developed based on 7 different algorithms, each of it have different characteristics so that our product can handle different dataset.

*Click the dropdown list to see algorithms option we provide*



# Supervised Machine Learning for Microbiome Data

Kristina Herman<sup>1</sup> Lillian Tatka<sup>2</sup> David Lee<sup>2</sup> Sierra Gillman<sup>3</sup> Binnan Yu<sup>4,5</sup>

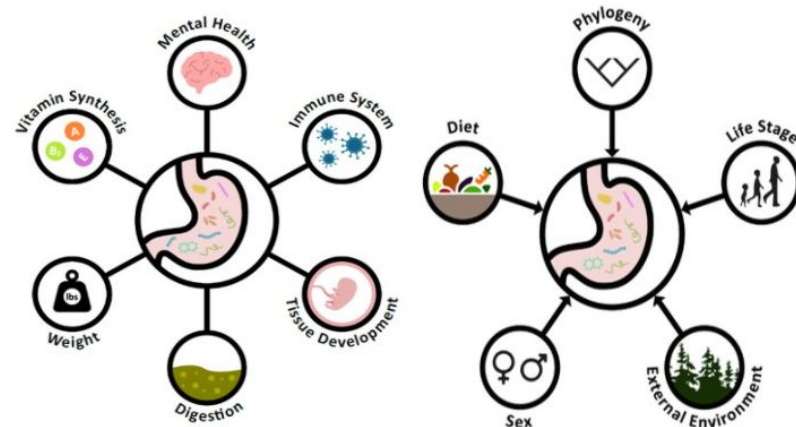
<sup>1</sup>Department of Chemistry; <sup>2</sup>Department of Bioengineering; <sup>3</sup>School of Environmental and Forest Sciences;

<sup>4</sup>Department of Rehabilitation Medicine; <sup>5</sup>Department of Applied Mathematics, University of Washington, Seattle, WA, 98105

## Statement of The Problem

Once thought only to be pathogenic, the microorganisms living on and within an animal host are now recognized as playing critical roles in host health. Factors that shape gut microbial communities are multifaceted and include the host's diet and life-stage.

While microbial shifts have been implicated in numerous human ailments (e.g., obesity, anxiety, inflammatory bowel disease), research has thus far been limited to differentiating microbial communities between groups and less for predictive uses. As a result of an ever expanding microbiome data availability, microbiome research lends itself to advancement with machine learning and data science. However, the implementation of machine learning in microbiome research might feel daunting and time consuming to those outside of the realm of data science.



Thank You!