

# The in<u>Telligence And Machine IEarning</u> (TAME) Toolkit for Introductory Data Science in Environmental Health

#### Julia E. Rager, PhD, MSEE

Assistant Professor | University of North Carolina at Chapel Hill (UNC) Department of Environmental Sciences and Engineering (ENVR) Institute for Environmental Health Solutions (IEHS) Center for Environmental Medicine and Lung Biology (CEMALB) Curriculum in Toxicology and Environmental Medicine (CiTEM)

#### Kyle Roell, PhD

Data Analyst | University of North Carolina at Chapel Hill (UNC) Department of Environmental Sciences and Engineering (ENVR) Institute for Environmental Health Solutions (IEHS)







- Introduction to TAME and its development Dr. Julia Rager
- Real-time demonstration of how to use TAME Dr. Kyle Roell
- Q&A is encouraged throughout and also at the end



































How can we make use of these growing resources to better understand environmental stressors???

### **DATA SCIENCE TRAINING**



# **Data Science Training through TAME Toolkit**



¥ f

- Multiple authors share knowledge and translate experiences gathered by multiple authors on data management and analysis methods to address real-world issues in environmental health through carefully developed training modules
- This training effort is currently being launched as the in<u>Telligence And</u> <u>Machine IEarning (TAME)</u> Toolkit
- Promotes didactic data generation, management, and analysis methods to "TAME" data in environmental health studies



≡ Q A ≟ i

Research

Preface

Background

Objectives

CHAPTER 1 INTRODUCTORY DATA SCIENCE

1.1 Introduction to Coding in R

1.2 Data Organization Basics 1.3 Finding and Visualizing Data Trends

1.4 High-Dimensional Data Visualizations

1.5 FAIR Data Management Practices CHAPTER 2 CHEMICAL-BIOLOGICAL ANALYSES AND PREDICTIVE MODELING

2.1 Dose-Response Modeling 2.2 Machine Learning and Predictive M.. 2.3 Mixtures Analysis 2.4 -Omics Analyses and Systems Biol.. 2.5 Toxicokinetic Modeling

2.6 Read-Across Toxicity Predictions CHAPTER 3 ENVIRONMENTAL HEALTH DATABASE MINING

3.1 Comparative Toxicogenomics Data...

3.2 Gene Expression Omnibus

ADDITIONAL RESOURCES

Published with bookdow

Resources

3.3 Database Integration: Air Quality St.

#### Module Development Overview

Training modules were developed to provide applications-driven examples of data organization and analysis methods that can be used to address environmental health questions. Target audiences for

The inTelligence And Machine lEarning (TAME)

**Biological Analyses, Predictive Modeling, and** 

**Database Mining for Environmental Health** 

Caviness, David M. Reif, Ilona Jaspers, Rebecca C. Fry, and Julia E. Rager

surrounding the training of researchers on these in silico methods

these modules can be found at the parent UNC-SRP Github Page.

Toolkit for Introductory Data Science, Chemical-

Kyle Roell, Lauren Koval, Rebecca Boyles, Grace Patlewicz, Caroline Ring, Cynthia Rider, Cavin Ward-

Research in exposure science, toxicology, and environmental health is becoming increasingly reliant upon data science and computational methods that can more efficiently extract information from

to chemicals in the environment and human disease outcomes. Still, there remains a critical gap

We aimed to address this critical gap by developing the inTelligence And Machine IEarning (TAME) Toolkit, promoting trainee-driven data generation, management, and analysis methods to "TAME" data in

environmental health studies. This toolkit encompasses training modules, organized as chapters within

this Github Bookdown site. All underlying code (in RMarkdown), input files, and imported graphics for

complex datasets. These methods can be leveraged to better identify relationships between exposures

https://uncsrp.github.io/Data-Analysis-Training-Modules/

https://github.com/UNCSRP



### **Wonderful Team of TAME Contributors**





Kyle Roell Data Analyst, UNC-SRP, UNC



Lauren Koval PhD Candidate, Ragerlab, UNC



**Rebecca Boyles** Data Management Director and Scientist, RTI

Headshot not preferred





**Caroline Ring** Computational Exposure Scientist, US EPA



Cynthia Rider Mixtures Toxicologist, NIEHS / NTP



Cavin Ward-Caviness Computational Biologist, US EPA



**David Reif** Professor and NCSU-SRP Data Management and Analysis Lead, NCSU



**IIONA JASPERS** Professor and Curriculum in Environmental Medicine and Toxicology Director, UNC



**Rebecca Fry** Professor and Superfund Research Program (UNC-SRP) Director, UNC



Julia Rager Assistant Professor and UNC-SRP Data Management and Analysis Co-Lead, UNC Rage

### **TAME Modules Presented through Bookdown**



	≡ Q. A. Å. i 9 f ≺	vvitn u
	The inTelligence And Machine lEarning (TAME)	<> Code
TAME Toolkit	Toolkit for Introductory Data Science, Chemical-	
Preface	Biological Analyses, Predictive Modeling, and	Data
CHAPTER 1 INTRODUCTORY DATA SCIENCE	Database Mining for Environmental Health	រះ main
1.1 Introduction to Coding in R	Research	
1.2 Data Organization Basics		🤮 jrag
1.3 Finding and Visualizing Data Trends	Kyle Roell, Lauren Koval, Rebecca Boyles, Grace Patlewicz, Caroline Ring, Cynthia Rider, Cavin Ward-	
1.4 High-Dimensional Data Visualizations	Caviness, David M. Reif, Ilona Jaspers, Rebecca C. Fry, and Julia E. Rager	Cha
1.5 FAIR Data Management Practices		📄 Cha
CHAPTER 2 CHEMICAL- BIOLOGICAL ANALYSES	Preface	Cha
2.1 Dece-Response Modeling	Background	L .DS
2.2 Machine Learning and Predictive M	Research in exposure science, toxicology, and environmental health is becoming increasingly reliant	🗋 REA
2.3 Mixtures Analysis	upon data science and computational methods that can more efficiently extract information from	
2.4 -Omics Analyses and Systems Biol	complex datasets. These methods can be leveraged to better identify relationships between exposures	README
2.5 Toxicokinetic Modeling	surrounding the training of researchers on these in silico methods.	
2.6 Read-Across Toxicity Predictions	Objectives	Th
CHAPTER 3 ENVIRONMENTAL	We aimed to address this critical gap by developing the inTelligence And Machine (Earning (TAME)	Ch
3.1 Comparative Toxicogenomics Data	Toolkit, promoting trainee-driven data generation, management, and analysis methods to "TAME" data in	an
3.2 Gene Expression Omnibus	environmental health studies. This toolkit encompasses training modules, organized as chapters within this Github Bookdown site. All underlying code (in RMarkdown) input files and imported graphics for	
3.3 Database Integration: Air Quality St	these modules can be found at the parent UNC-SRP Github Page.	Re
ADDITIONAL RESOURCES	Module Development Overview	These
Resources	Training modules were developed to provide applications-driven examples of data organization and	Mach
Published with bookdown	analysis methods that can be used to address environmental health questions. Target audiences for	https://

https://uncsrp.github.io/Data-Analysis-Training-Modules/

#### 1 4 1:11 inderlying script and datasets posted on Github:

UNCSRP / Data-Analysis-Training-Modules		Q Type 🛛 to search	
<> Code 💿 Issues 🏦 Pull reque	sts 🕑 Actions 🗄 Projects	🖽 Wiki 🙂 Security	✓ Insights
🕮 Data-Analysis-Training-N	Nodules (Public)		⊙ Watch 1
<pre> % main - % 2 branches  \$ 0 </pre>	tags	Go to file Add file *	<> Code -
jrager Update README.md		644c444 on Mar 11, 2022	3 41 commits
Chapter 1	Update README.md		last year
Chapter 2	Update README.md		last year
Chapter 3	Update README.md		last year
DS_Store	Update Folders		2 years ago
C README.md	Update README.md		last year
README.md			Ø

.md

e TAME Toolkit for Introductory Data Science, nemical-Biological Analyses, Predictive Modeling, d Database Mining for Environmental Health search

e files represent the underlying script and input files associated with the inTelligence And nine IEarning (TAME) Toolkit, which is located at the following associated bookdown site: ://uncsrp.github.io/Data-Analysis-Training-Modules/index.html.

https://github.com/UNCSRP



# **Parent Manuscript**



> Front Toxicol. 2022 Jun 22;4:893924. doi: 10.3389/ftox.2022.893924. eCollection 2022.

Development of the InTelligence And Machine LEarning (TAME) Toolkit for Introductory Data Science, Chemical-Biological Analyses, Predictive Modeling, and Database Mining for Environmental Health Research

```
Kyle Roell <sup>1</sup>, Lauren E Koval <sup>1</sup> <sup>2</sup>, Rebecca Boyles <sup>3</sup>, Grace Patlewicz <sup>4</sup>, Caroline Ring <sup>4</sup>,
Cynthia V Rider <sup>5</sup>, Cavin Ward-Caviness <sup>6</sup>, David M Reif <sup>7</sup>, Ilona Jaspers <sup>1</sup> <sup>2</sup> <sup>8</sup> <sup>9</sup> <sup>10</sup>,
Rebecca C Fry <sup>1</sup> <sup>2</sup> <sup>8</sup>, Julia E Rager <sup>1</sup> <sup>2</sup> <sup>8</sup> <sup>9</sup>
Affiliations + expand
PMID: 35812168 PMCID: PMC9257219 DOI: 10.3389/ftox.2022.893924
Free PMC article
```



### **Chapter—Based Content Organization**







# **Content within Each Specific Module**



- Introduction to the Topic
- Introduction to the Training Module
- List of Training Module's *Environmental Health Questions*
- Script Preparations
- Scripted Analysis with Applications-based Environmental Health Questions Sprinkled throughout to Keep Participants Engaged
- Concluding Remarks and Additional Resources



### Screenshots in the Mixtures Sufficient Similarity Module





2.6 Read-Across Toxicity Predictions

### **Screenshots in a Machine Learning Module**







### Screenshots in the -Omics Analyses Module

-

Preface

DATA SCIENCE

2.3 Mixtures Analysis

Transcriptomics

MA Plots

Volcano Plots







# **TAME Toolkit Dissemination**







- Much of this training content was tested through a new course at UNC, ENVR 730: Computational Toxicology and Exposure Science
- Included 17 graduate students from:
  - Department of Environmental Sciences and Engineering (ENVR)
  - Curriculum of Toxicology and Environmental Medicine (CiTEM)
- Received the Teaching Innovation Award from the UNC Gillings School of Global Public Health, which is student-nominated!
- We are teaching this course again this Fall



## **Other Dissemination Efforts**



# High dimensional data visualization training workshop (virtual, Spring 2022)



Workshop organized through the UNC Computational Biosciences Club (CBC) and

Curriculum in Toxicology & Environmental Medicine (CiTEM)

CITEN

# PRogramming for Environmental <u>HE</u>alth <u>And Toxicology</u> (PREHEAT) Retreat





## **International Engagement**



### Viewers have accessed the TAME toolkit world-wide

Google analytics (Sept 2022-current):

Global views across 91 countries:	Number of users per country with highest usage:	
	COUNTRY United States Germany Canada India	USERS 600 62 60 59
5,300 total site views 1,441 full site users	China France	39

### International collaborators

For example, the **European Union's ONTOX** project leaders are incorporating TAME materials into their training courses on data analysis and machine learning in toxicology



https://ontox-project.eu/consortium/



# **TAME Effectiveness and Feedback**





#### **Consistent Feedback across All Dissemination Efforts**

- Participants want more! More training time, more modules
- Training at the intro-level, mid-level, and advanced-level
- More basics, including best practices for wet lab data organization and sharing
- We are currently in the process of expanding these training materials...







Expanding and reorganizing content into 6 chapters:

- 1. Introduction to Data Science
- 2. Introduction to Coding in R
- 3. Basics of Data Analysis and Visualizations
- 4. Machine Learning and Artificial Intelligence
- 5. Applications in Toxicology and Exposure Science
- 6. Environmental Health Database Mining

Hoping to launch ~December 2023



# **Collaborations and Funding**



#### Rager lab

Celeste Carberry (PhD cand) Elise Hickman (post-doc) Lauren Koval (PhD cand) Elena McDermott (MSPH) Alexis Payton (data analyst) Sarah Miller (PhD cand) Kai Malone (undergrad) Raquel Winker (undergrad)

### UNC IEHS / SRP

Rebecca Fry Audrey Bousquet (res scientist) Lauren Eaves (res associate) Hadley Hartwell (lab manager) Kyle Roell (data analyst)

#### <u>NIH / NIEHS</u>

Nicole Niehoff Matthew Wheeler Alexis White

#### <u>UNC</u>

Stephanie Engel Ilona Jaspers Alex Keil Kun Lu Tracy Manuck Meghan Rebuli Gregory Smith

**<u>RTI</u>** Rebecca Boyles Shaun McCullough

### <u>NIH / DTT</u>

Scott Auerbach Stephen Ferguson Kyle Messier David Reif Cynthia Rider

#### **US EPA**

Kathie Dionisio M Ian Gilmour Kristin Isaacs Yong Ho Kim Grace Patlewicz Katie Paul-Friedman Caroline Ring John Wambaugh Cavin Ward-Caviness

Trainees in italics



# Kyle Roell's Walk Through TAME



• Dr. Roell, UNC-SRP lead data analyst

