

Building Tools for Automatically Processing Lab Data

Rania Hakimi, REU Participant, Bunker Hill Community College

Wenjin Zhang, PhD Student, Northeastern University

Amy Mueller, Assistant Professor, Northeastern University



Northeastern University
College of Engineering



Abstract

Phosphorus (P) is a major nutrient that wastewater treatment plants are making all efforts to eliminate. The proposed cost-effective way to indirectly quantify P in real-time is through commercially available sensor array. This project aims to building algorithms to automatically extract useful information from available sensors data sheets, thus help build machine learning algorithms for P prediction.

Motivation

- ❖ Online sensors for P do not exist. The Mueller lab is bridging this gap by studying chemistry of a P removal reactor.
- ❖ This REU project creates tools that automates data extraction to support development of an online P sensor.

Goals

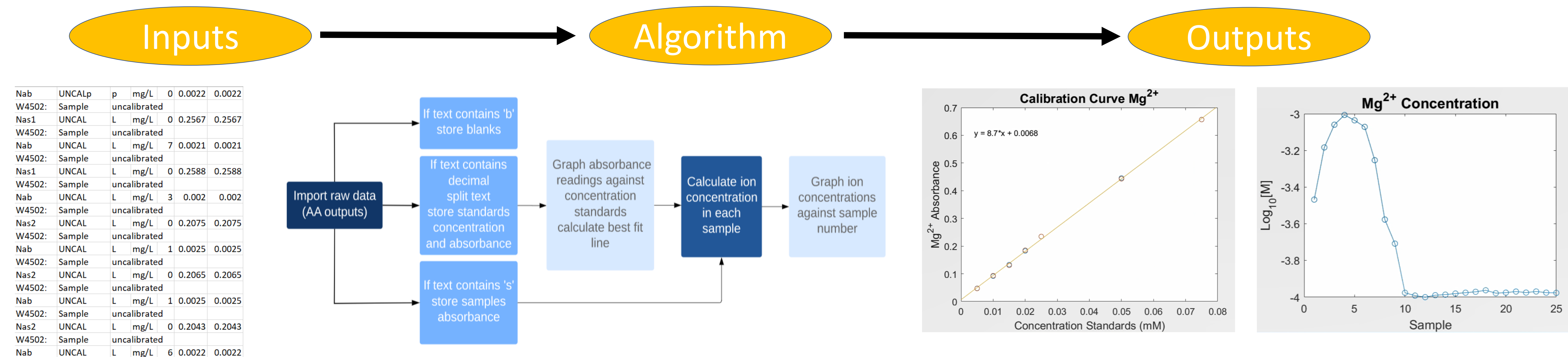
- ❖ Creation of a data processing/visualizing tool for Atomic Absorption Spectroscopy (AAS) log files.
- ❖ Creation of a tool to extract useful information from logged sensor data.
- ❖ Creation of a linear model for P prediction to assess data complexity.



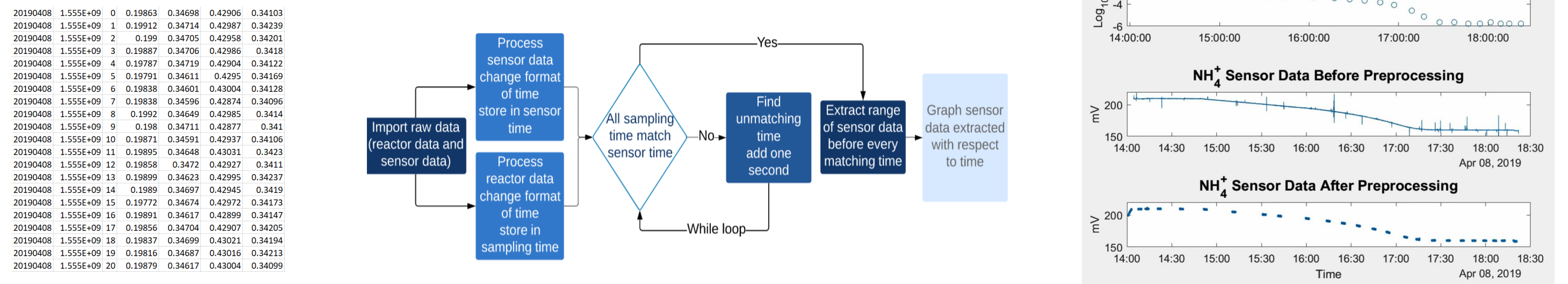
- ❖ Sensors include: electrodes for pH, DO, NH₄⁺, Na⁺, K⁺, Ca²⁺, ClO₄⁻, Cl⁻, Mg²⁺.

Methodology and Results

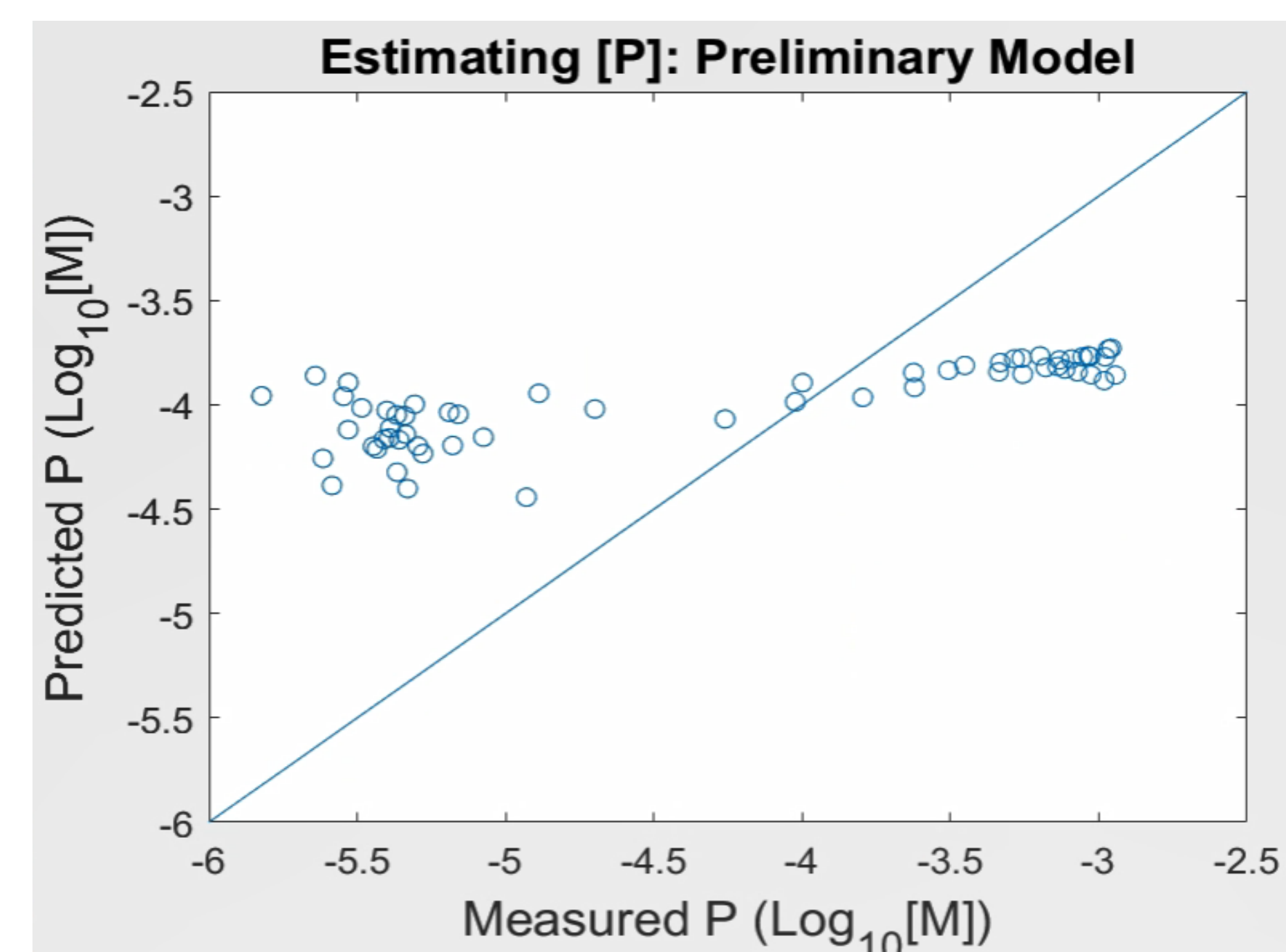
1. Data Processing/Visualizing Tool (AA Data)



2. Data Extraction Tool (Sensor Data)



3. Preliminary first-order analysis using MLR



- ❖ MLR is only valid when the input variables are uncorrelated – sensor data may be correlated.
- ❖ Relationship of P to other ions could be non-linear.
- ❖ Relationship of P to other ions may be a function of time in the reactor cycle.

Discussion and Conclusions

- ❖ AA sheets have different number of entries. Text processor has been built to extract useful AA data automatically.
- ❖ In order to match timestamps in sensor data to sampling data from 10 cycles, loops have been used to deliver this task.
- ❖ MLR is a baseline to study the relationship between P concentration and sensor data but it is not robust to capture the non-linear dynamic, in the future advanced models need to be used.

Acknowledgements: Department of Civil and Environmental Engineering at Northeastern University, Boston, MA
Department of Civil and Environmental Engineering at University of Washington, Seattle WA

This work was supported by the National Science Foundation Grant #1559894: NSF REU Data Driven Discovery (D3).