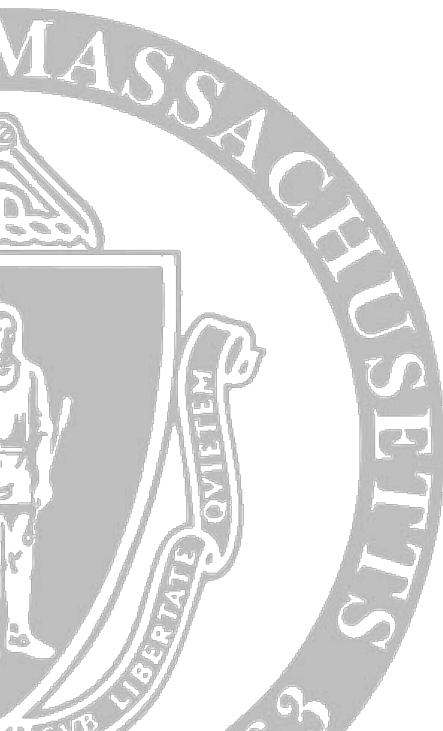


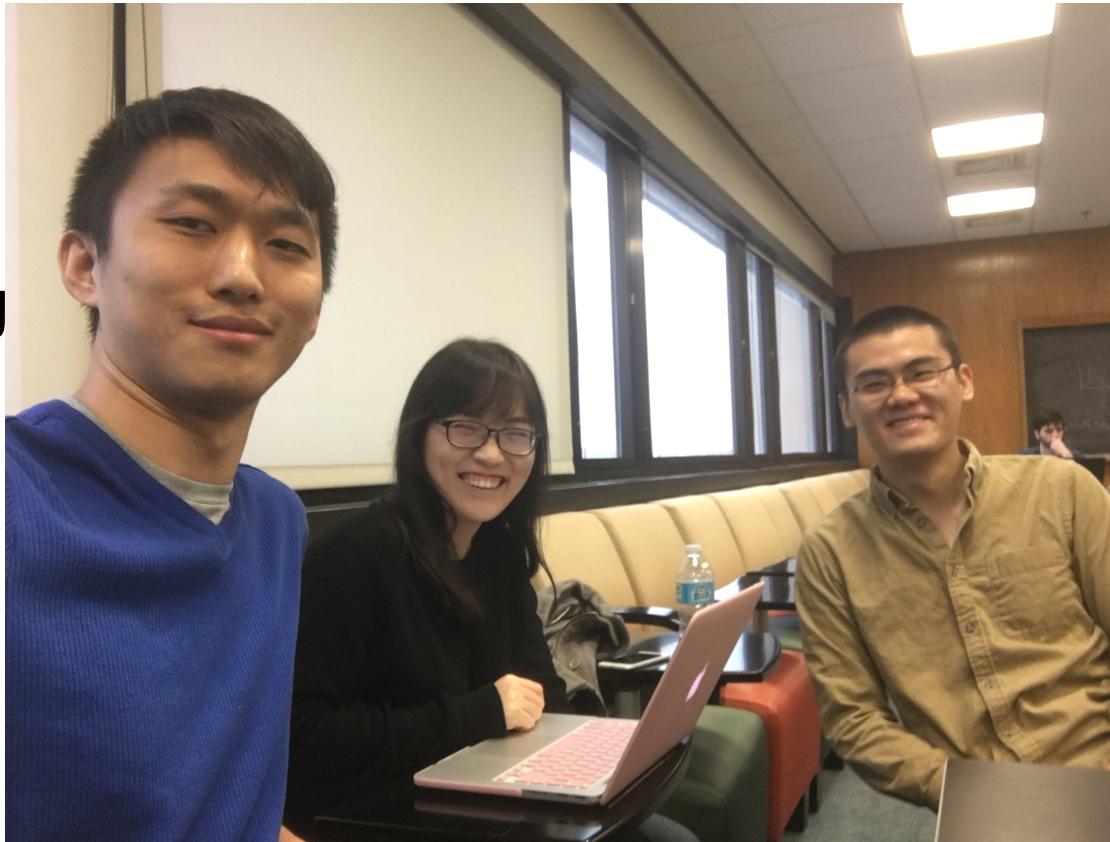
Lead Prediction at Flint

Team Pizza-Pita
Feb 26, 2017



Our Team: Pizza-Pita

Huan Chen
MS Stat



Shujian Liu
PhD MechEng

Tangxin Jin
PhD Math

Outline

- Background
- Exploration Datasets
- Regression
 - Azure ML
 - Point to Point Modeling
 - Sequence Modeling
- New Metric
- Results
- Conclusions
- Acknowledge

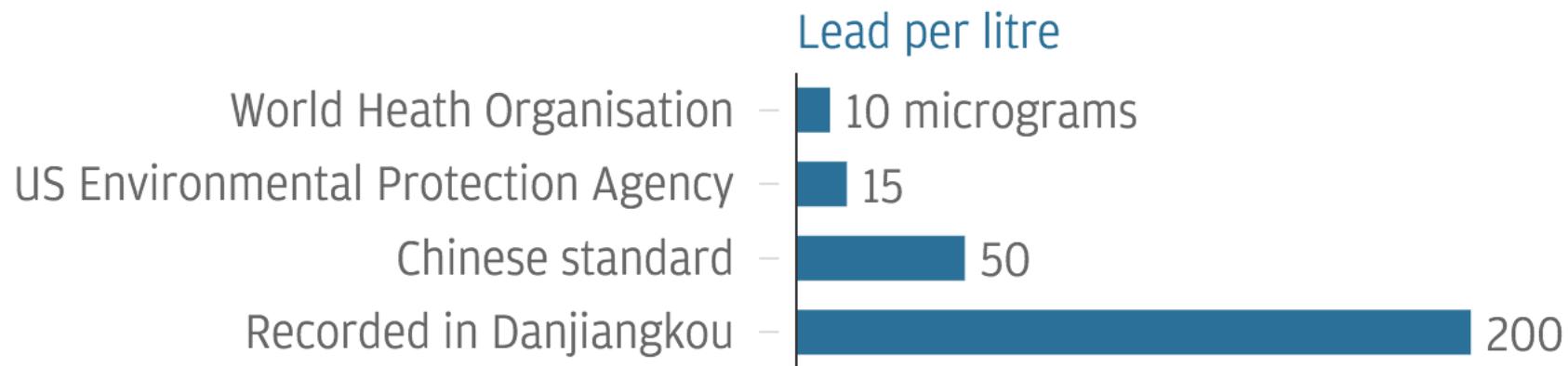
Background



Source: <http://northdallasgazette.com/2016/03/20/flint-water-crisis-move-children-babies/>

Safe Levels of Lead in Water

Maximum safe levels of lead in drinking water supply

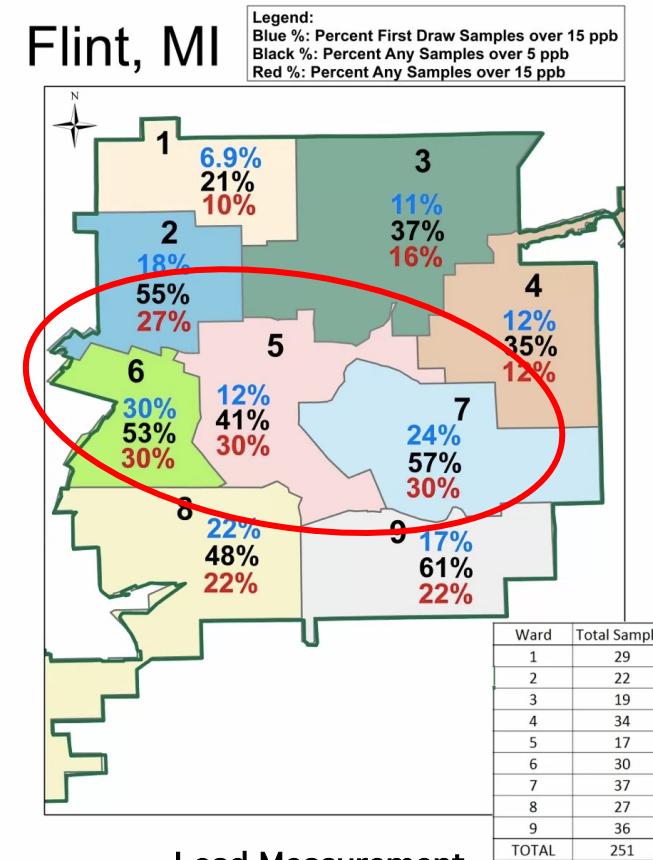
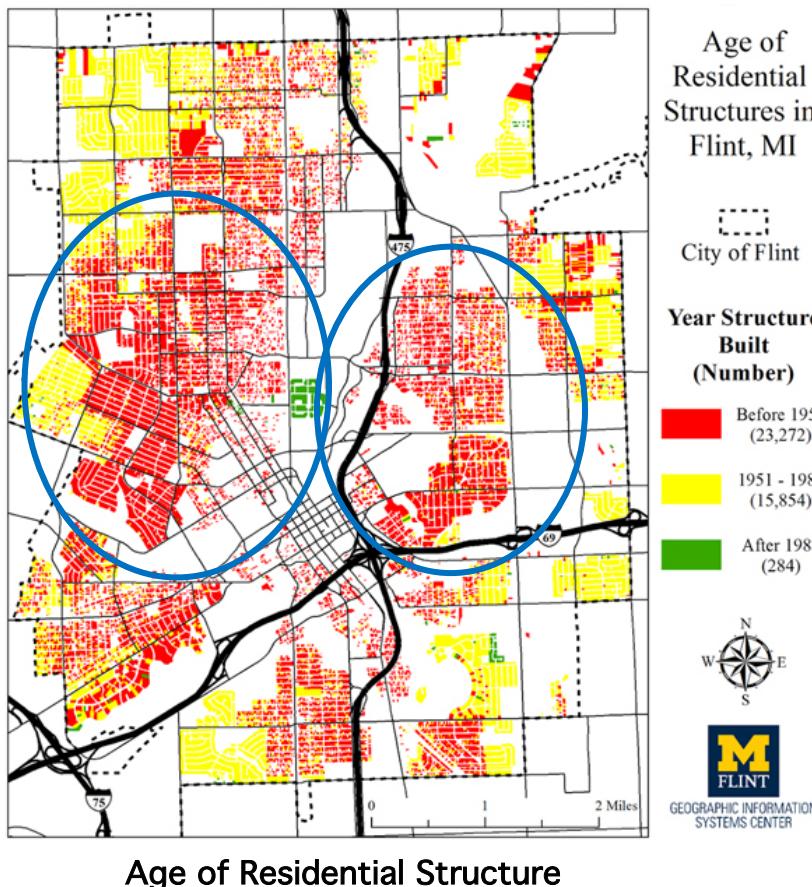


Source: WHO, EPA, MWR, Journal of Environmental Informatics

SCMP

House Age and Lead Measurement

Old
House

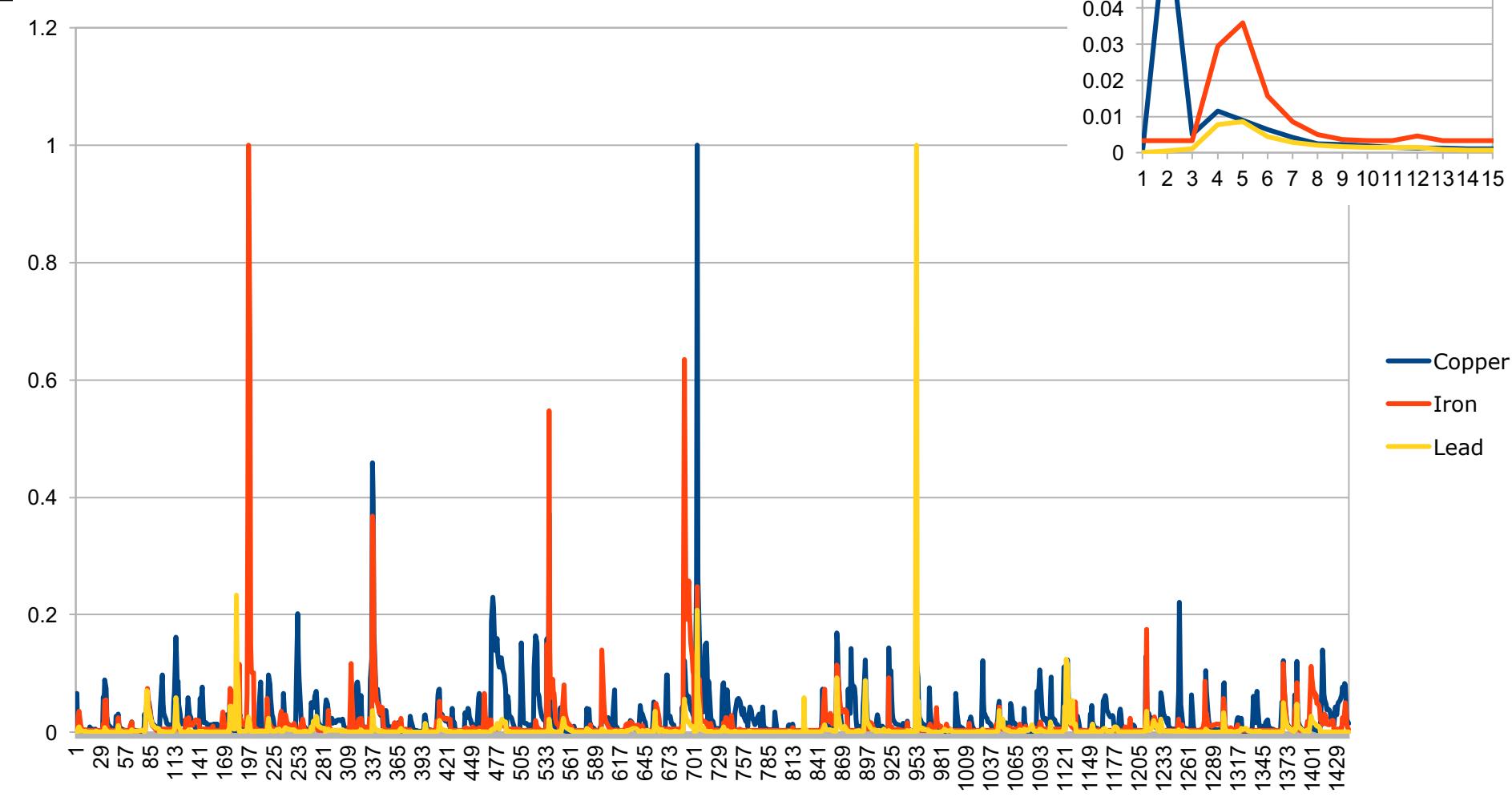


Source:

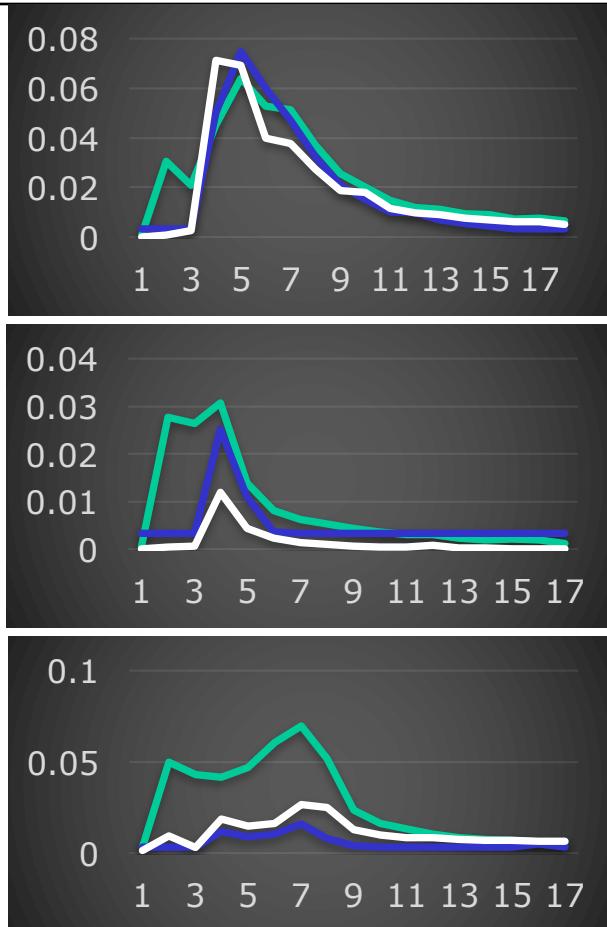
<http://www.fondriest.com/news/amid-flint-water-crisis-gis-effort-maps-citys-pipes.htm>

<http://flintwaterstudy.org/2015/09/distribution-of-lead-results-across-flint-by-ward-and-zip-codes/>

Normalized Raw Data

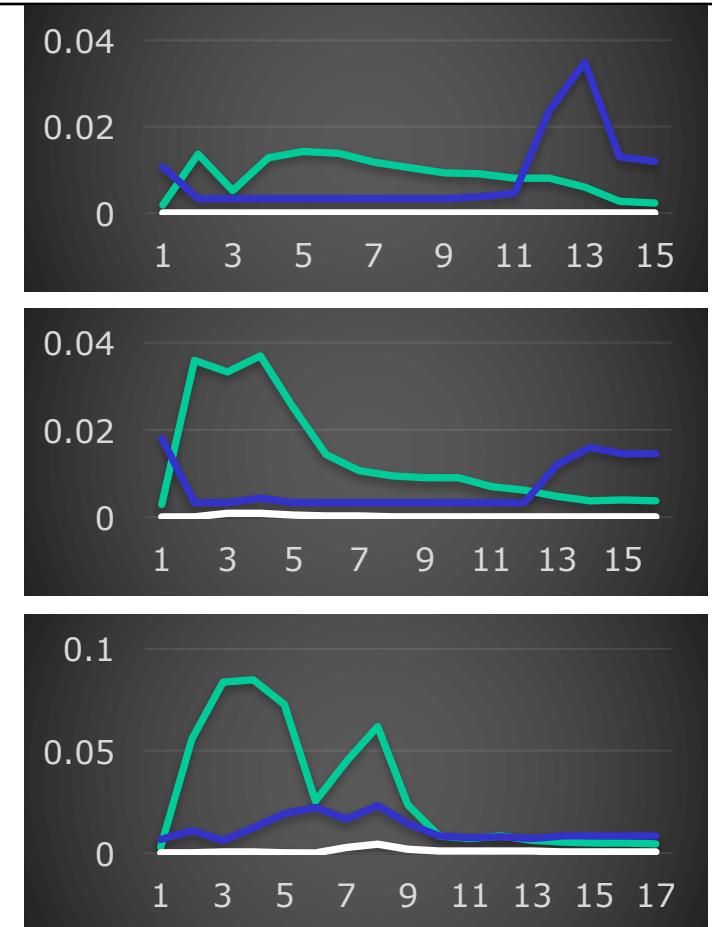


- Good

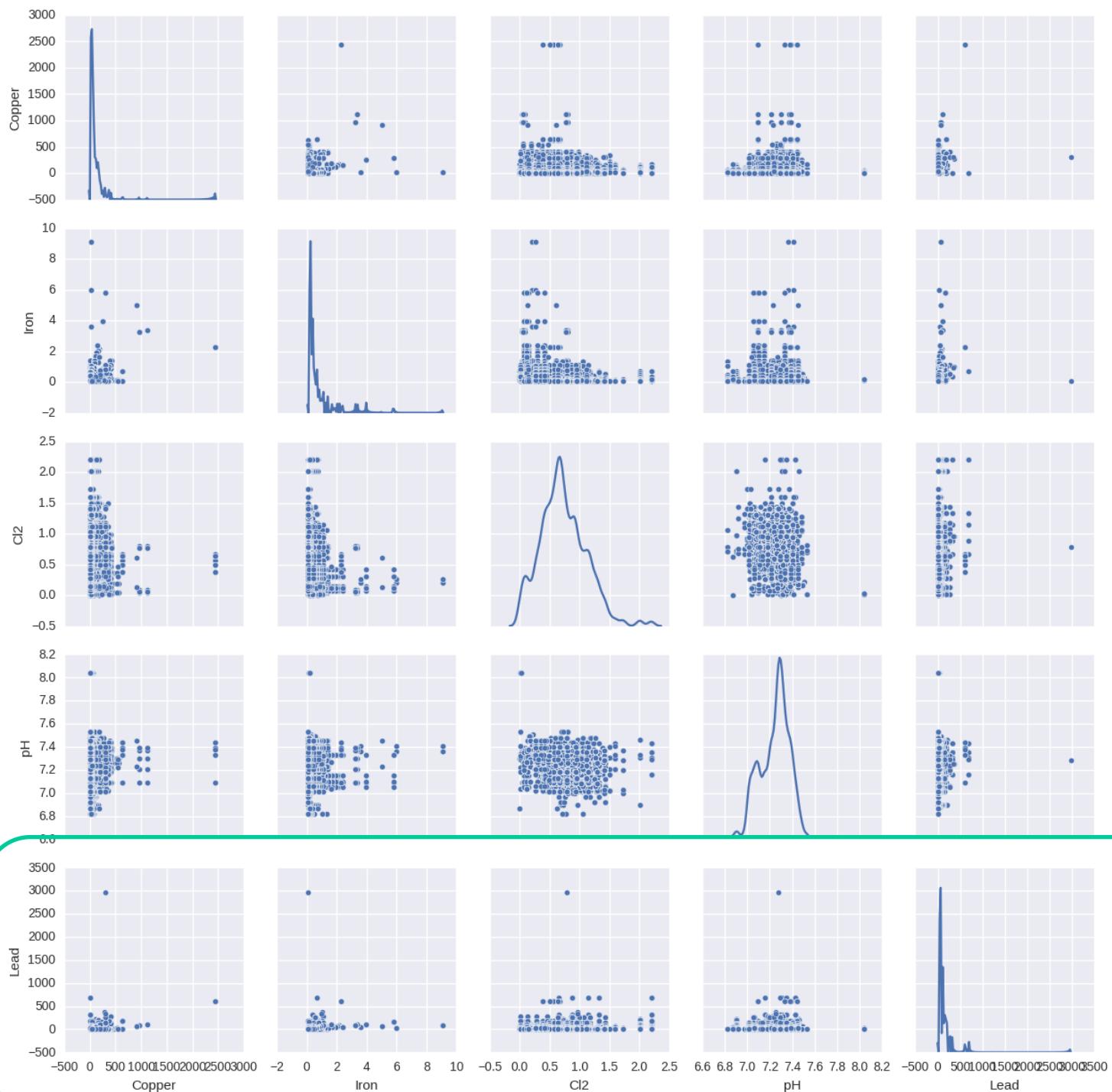


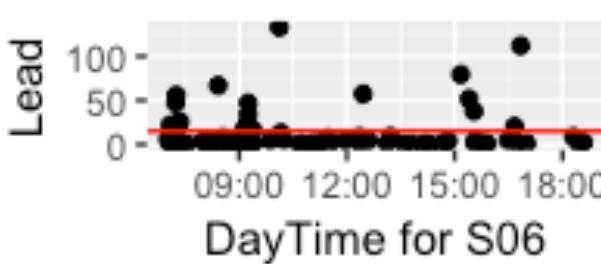
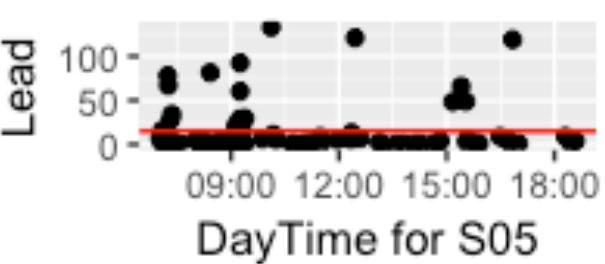
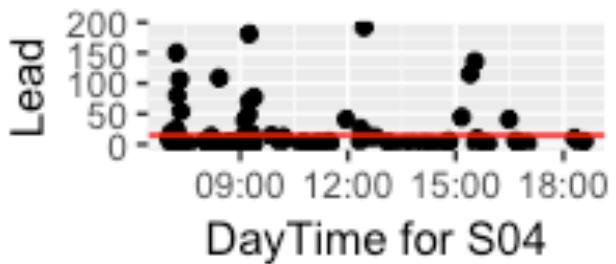
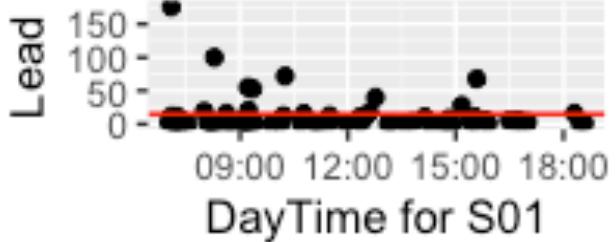
Green: Copper
Blue: Iron
White: Lead

- Hum...

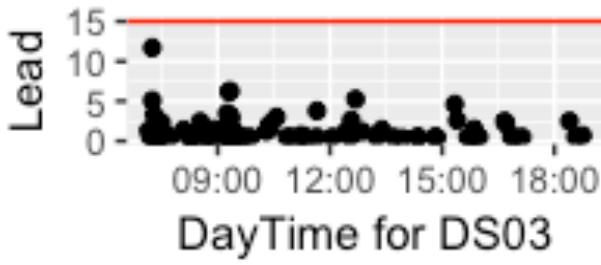
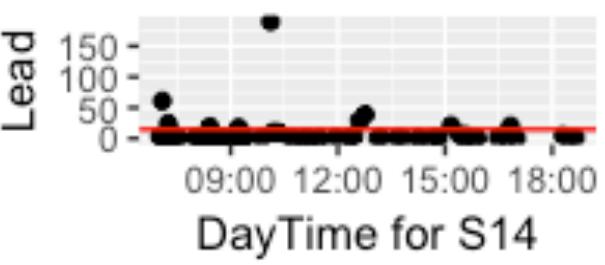
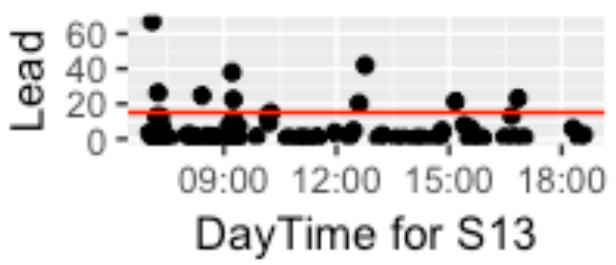
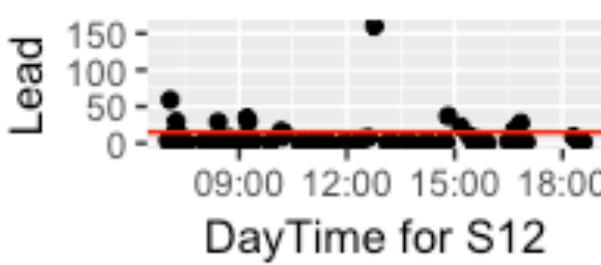
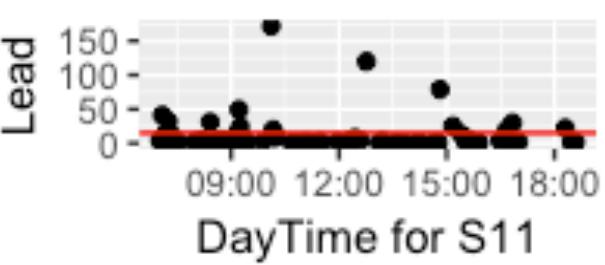
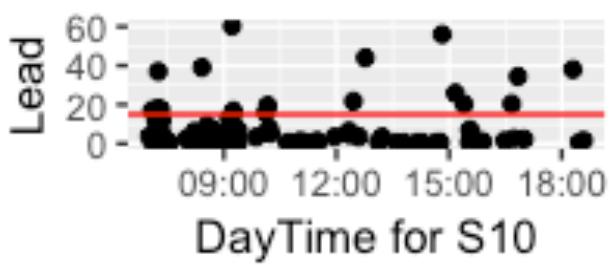
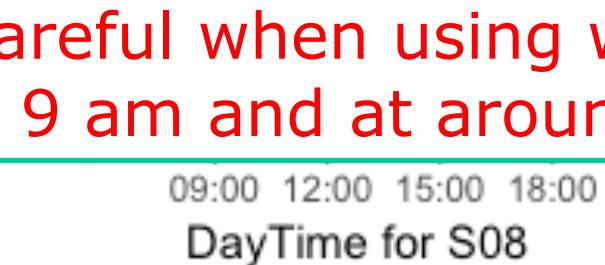
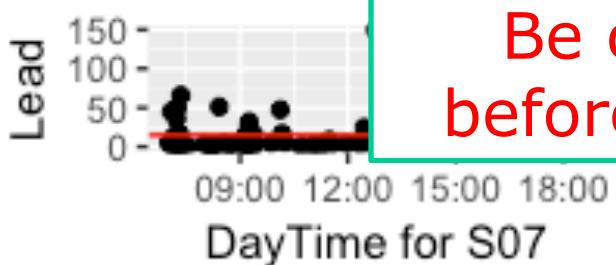


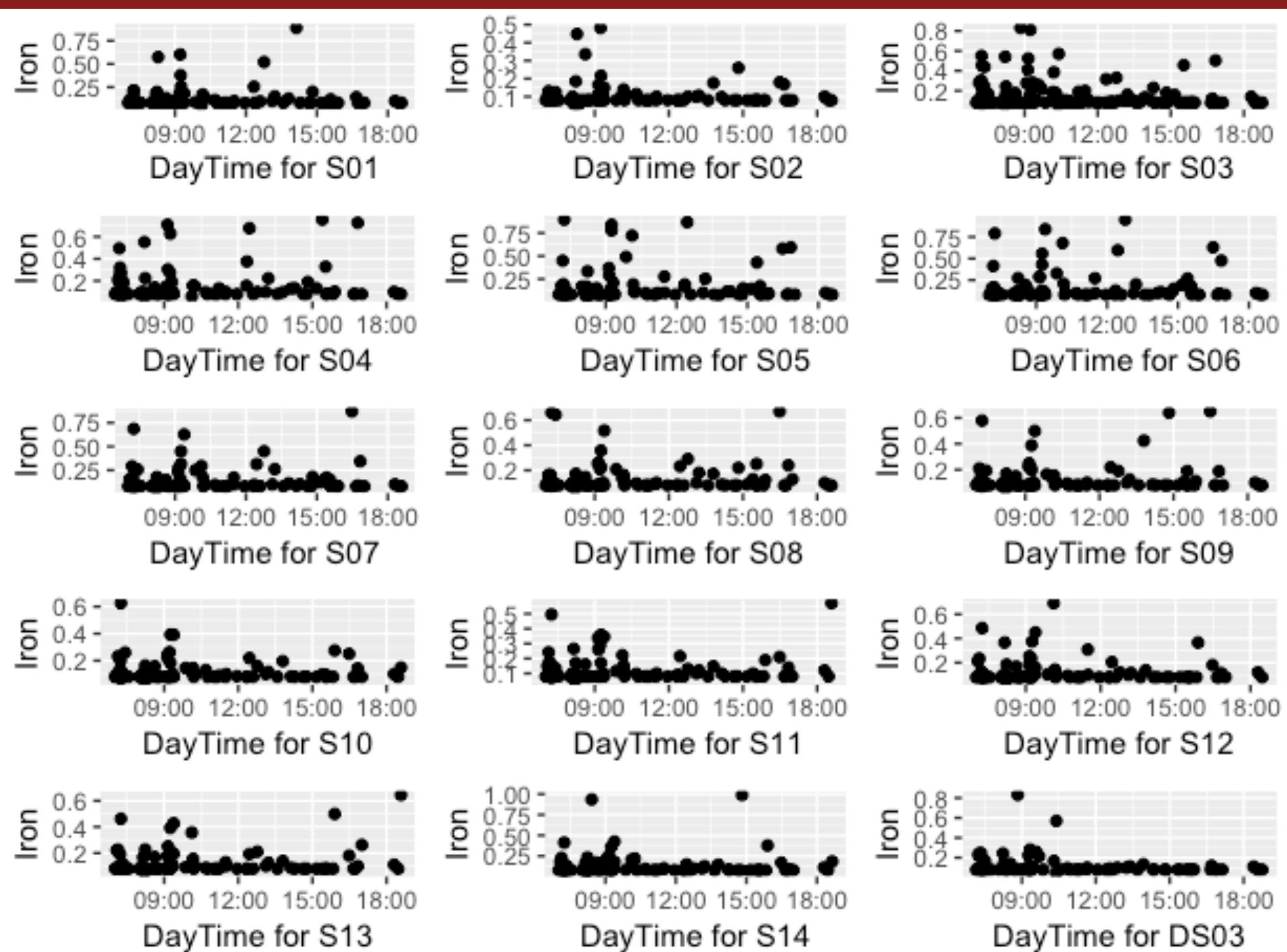
Lead





Be careful when using water before 9 am and at around 3pm

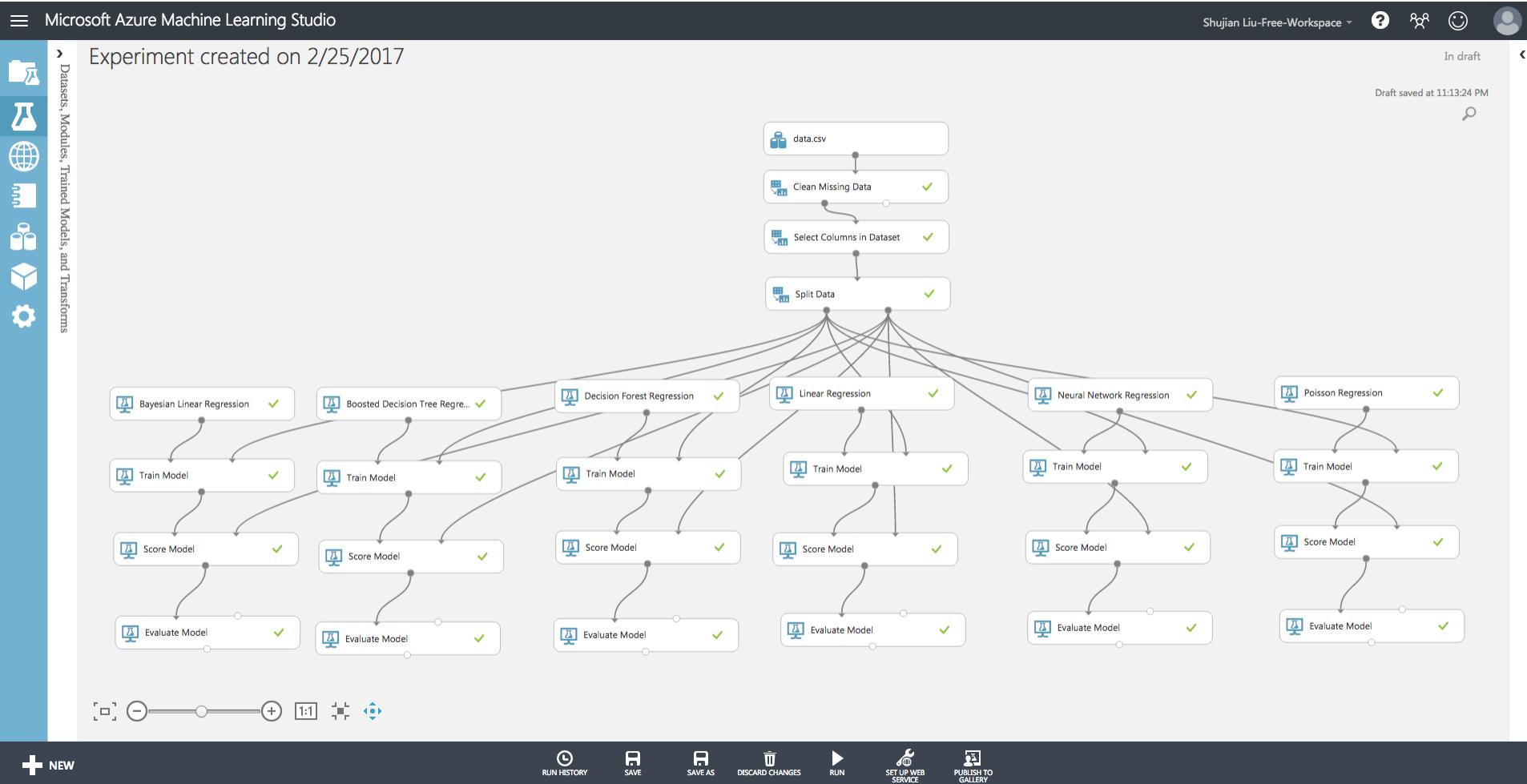




Data Manipulation

- Microsoft Azure Machine Learning Studio:
 - 4 variables
 - 75% of data for training
- Point-point Modeling (scikit-learn):
 - 2 variables: training-testing split 1000/443
 - 4 variables: training-testing split 900/368
- Sequence Modeling (Keras on TF):
 - Only used data from 15 samples (67 sequences)
 - 40 sequences data for training
 - 10% from training for cross-validation
 - 27 sequences data for testing

Azure ML



Results from Azure

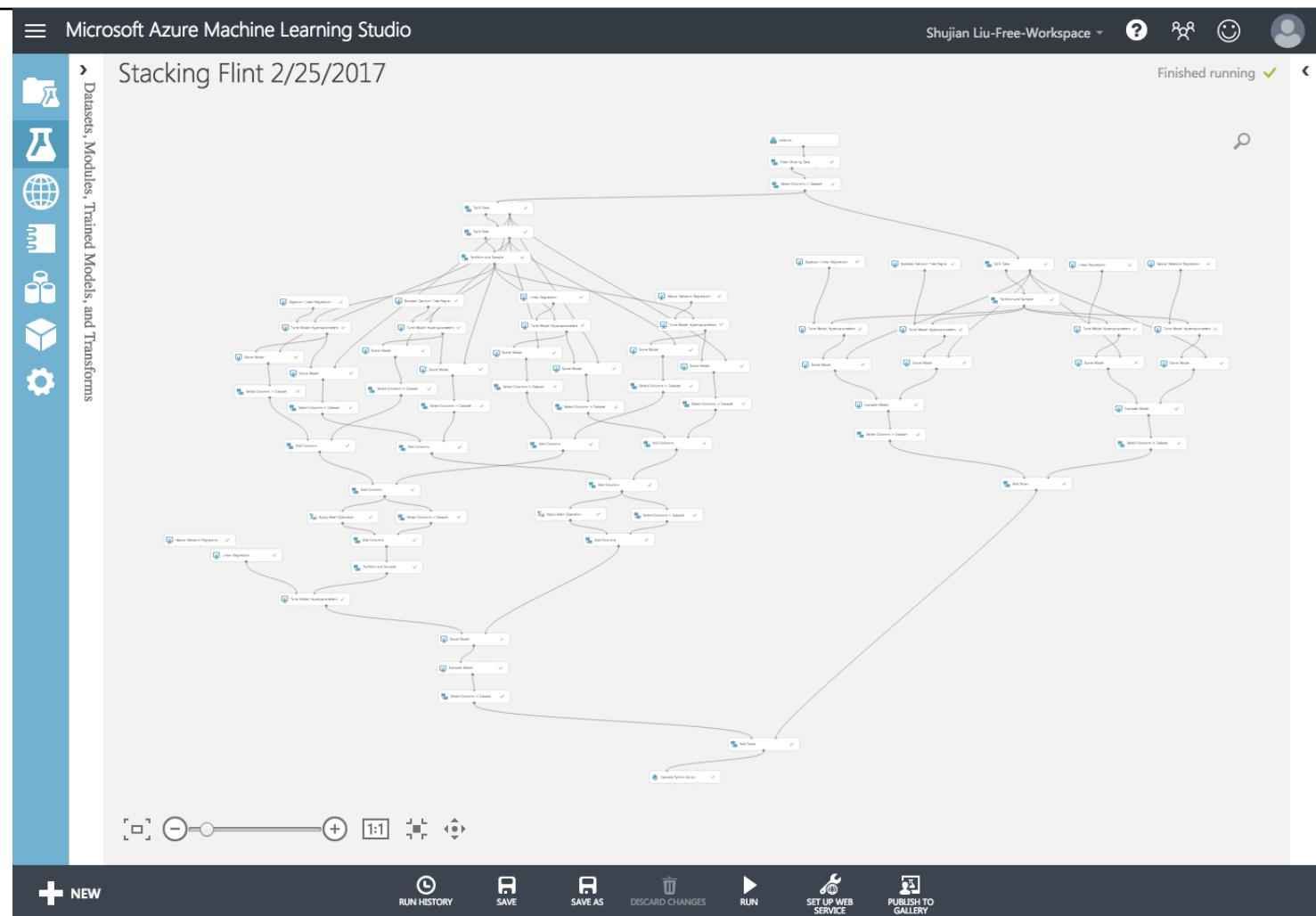


Microsoft Azure

Models	Mean Absolute Error	Root Mean Squared Error
Bayesian Linear Regression	21.263156	155.605892
Boosted Decision Tree	18.345794	157.411814
Decision Forest Regression	16.235807	157.512419
Linear Regression	21.581666	155.500647
Neural Network	22.395735	158.531654
Poisson Regression	23.83821	158.520784

Machine Learning in ONE HOUR

Ensemble Regressors using Stacking



Results from Stacking

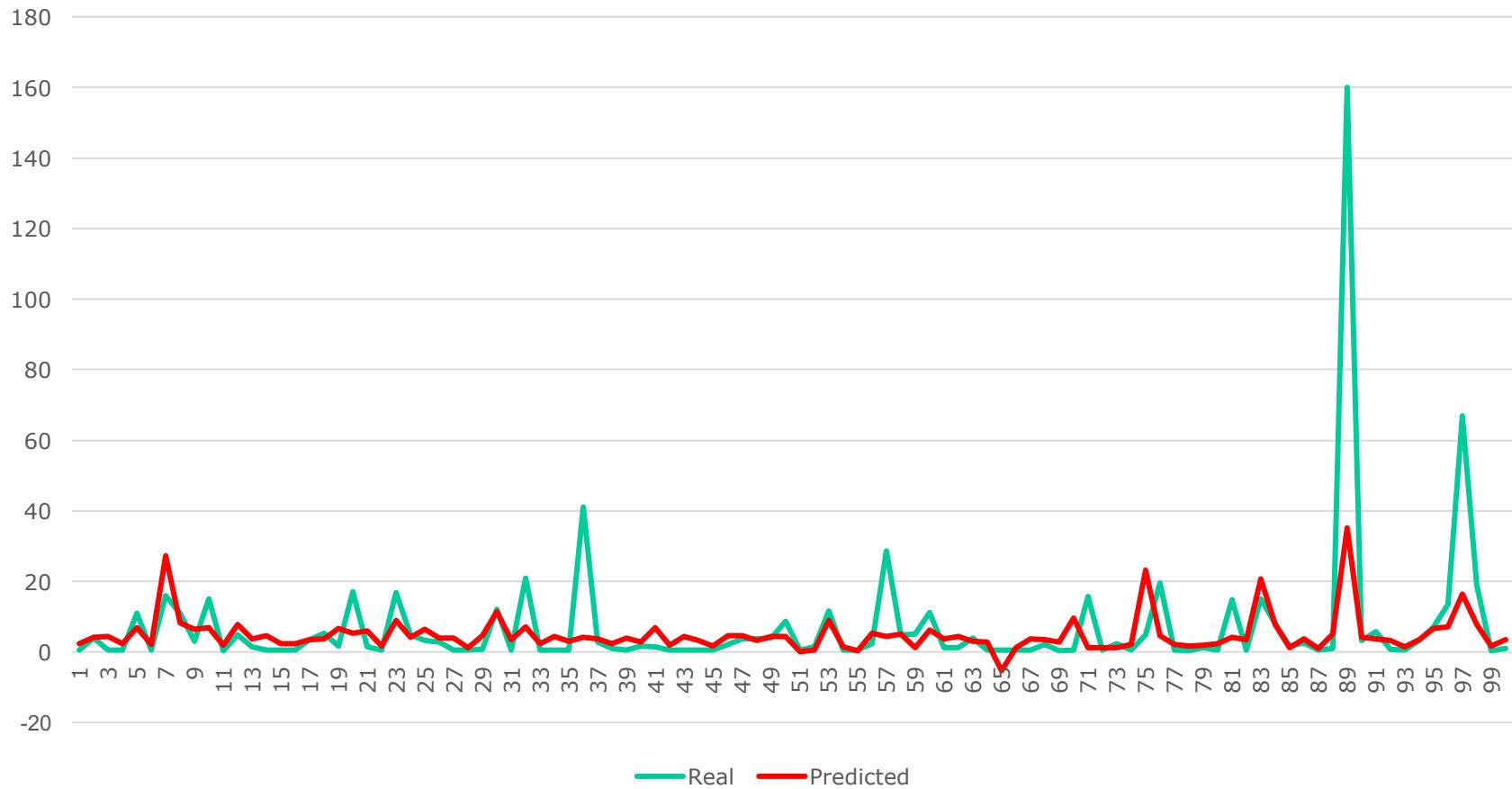
- Stacking of four basic regression algorithms:
 - Bayesian Linear Regression
 - Boosted Decision Tree Regression
 - Linear Regression
 - Neural Network Regression
- Metrics:

Mean Absolute Error	5.118008
Root Mean Squared Error	13.088811

- Compare to Decision Forest Regression with MAE of 16.2

Reference: <https://gallery.cortanaintelligence.com/Experiment/Building-Ensemble-of-Classifiers-using-Stacking-2>
Breiman, Leo. "Stacked regressions." *Machine learning* 24.1 (1996): 49-64.

Results from Stacking



Point to Point Modeling (scikit-learn)

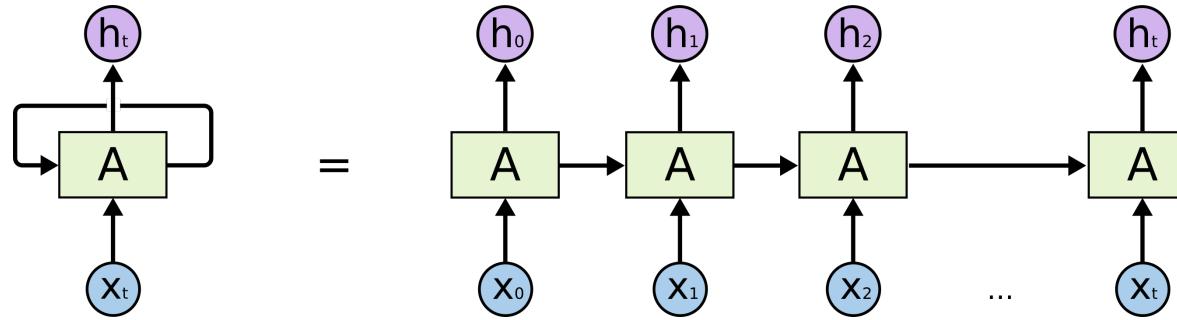
Training with 2 variables
(copper and iron)

Model	Linear regression	Ridge	Lasso	SVM (rbf)	Random Forest (5 trees)
MAE	14.21	14.21	16.90	15.86	20.70

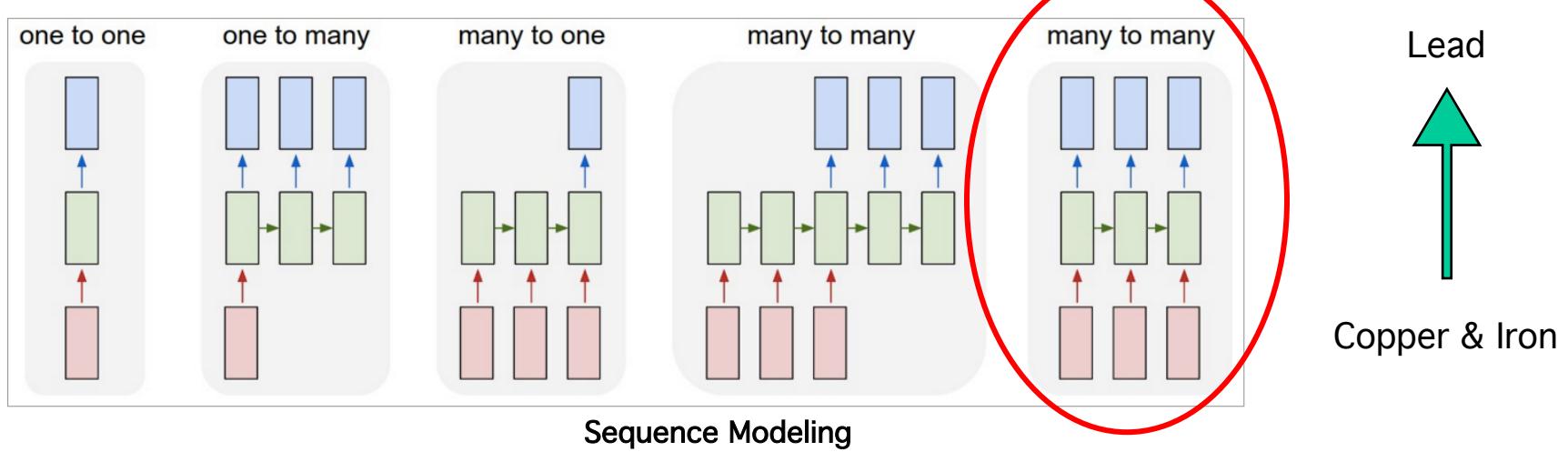
Training with 4 variables
(copper, iron, Cl2, pH)

Model	Linear regression	SVM (rbf)	Random Forest (5 trees)
MAE	19.21	16.78	18.94

Recurrent Neural Network / Long Short Term Memory



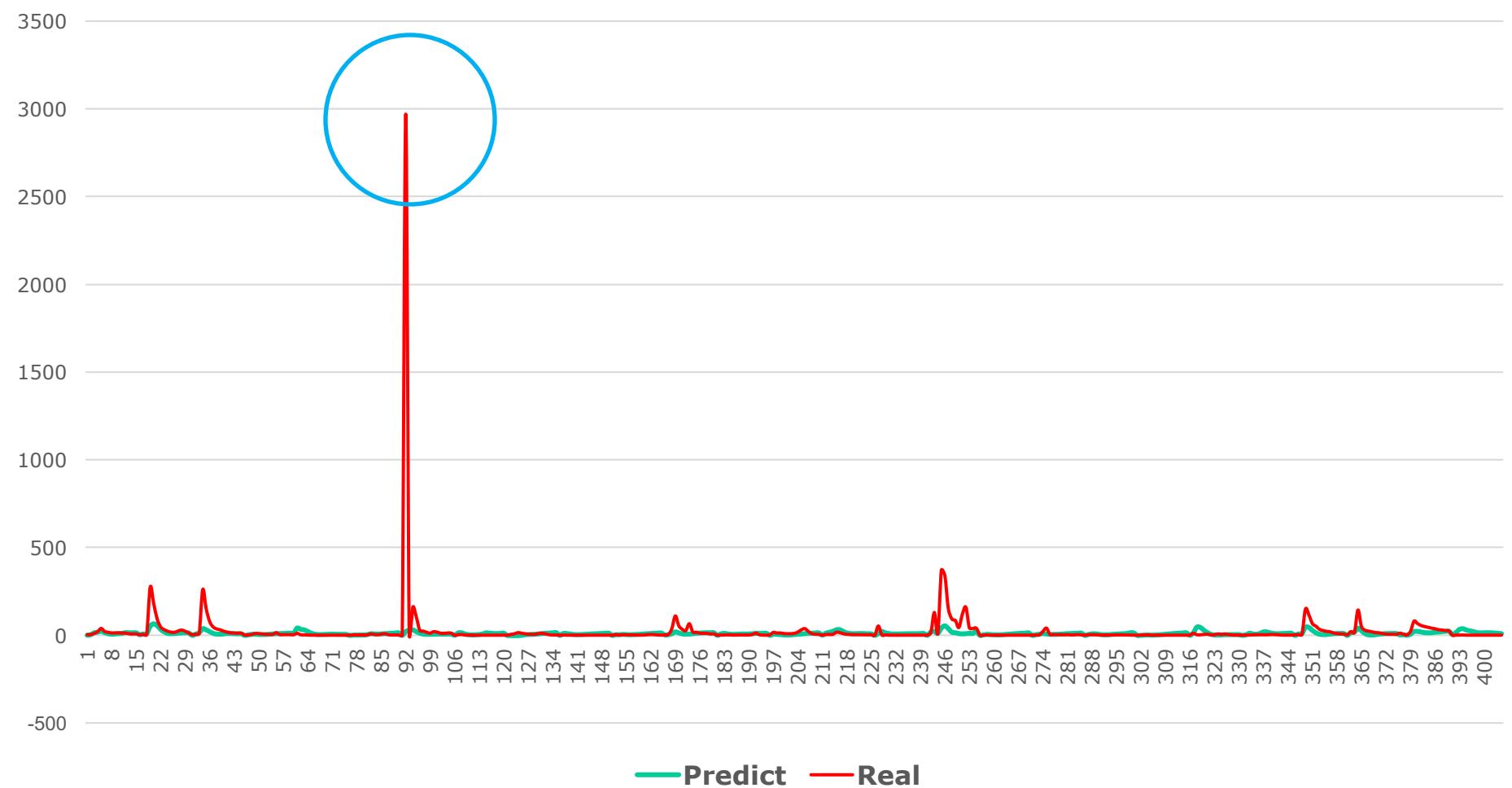
An unrolled recurrent neural network



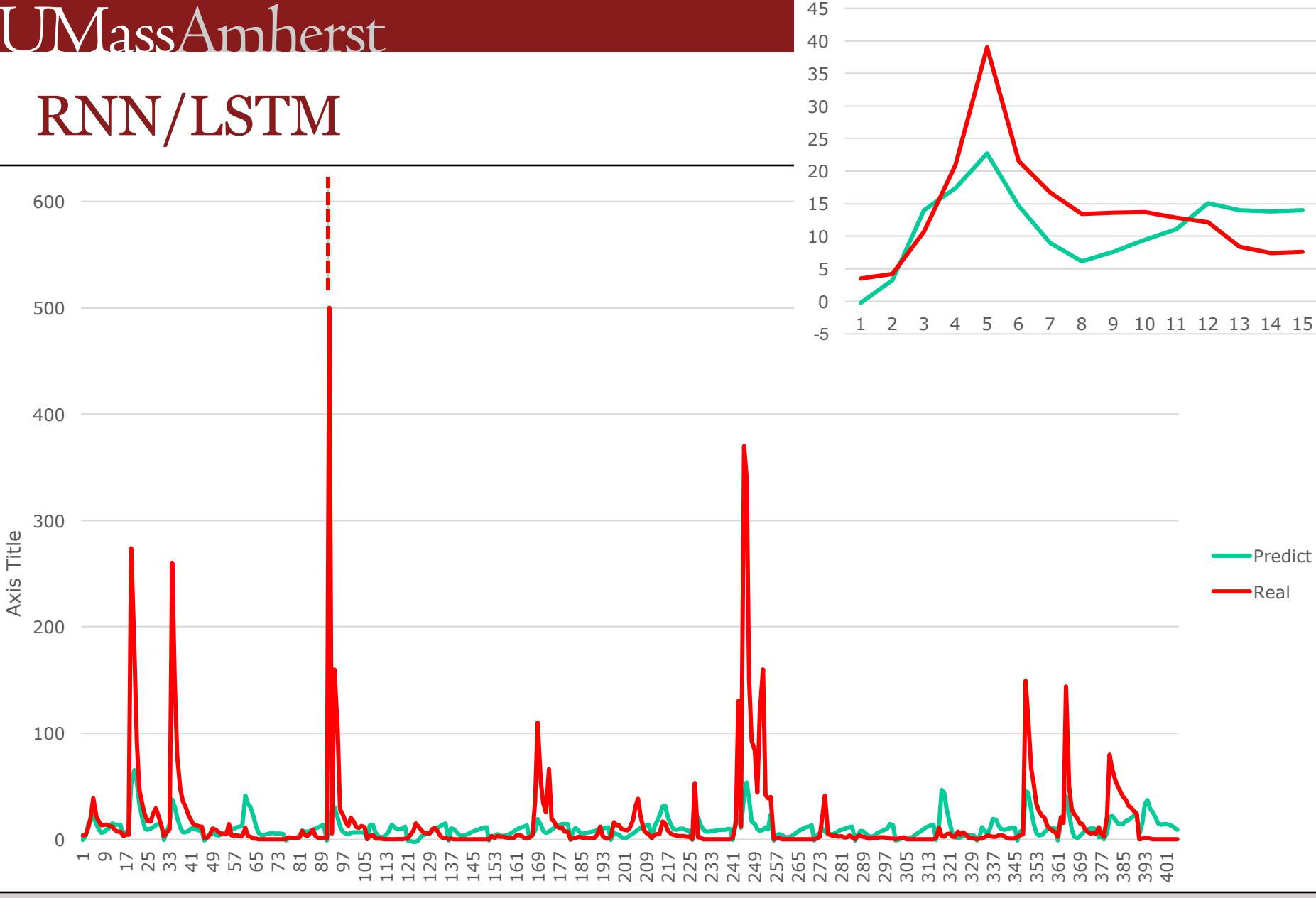
Source:

<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>
<http://karpathy.github.io/2015/05/21/rnn-effectiveness/>

Lead Prediction by RNN/LSTM



RNN/LSTM



New Metric

Maximum safe levels of lead in drinking water supply



Source: WHO, EPA, MWR, Journal of Environmental Informatics

SCMP

New Metric:

Instead of MAE/RMSE loss, we tried **rate of successfully predicted safety level (15 ug/L).**
~ in classification

New Metric – 15 ug/L

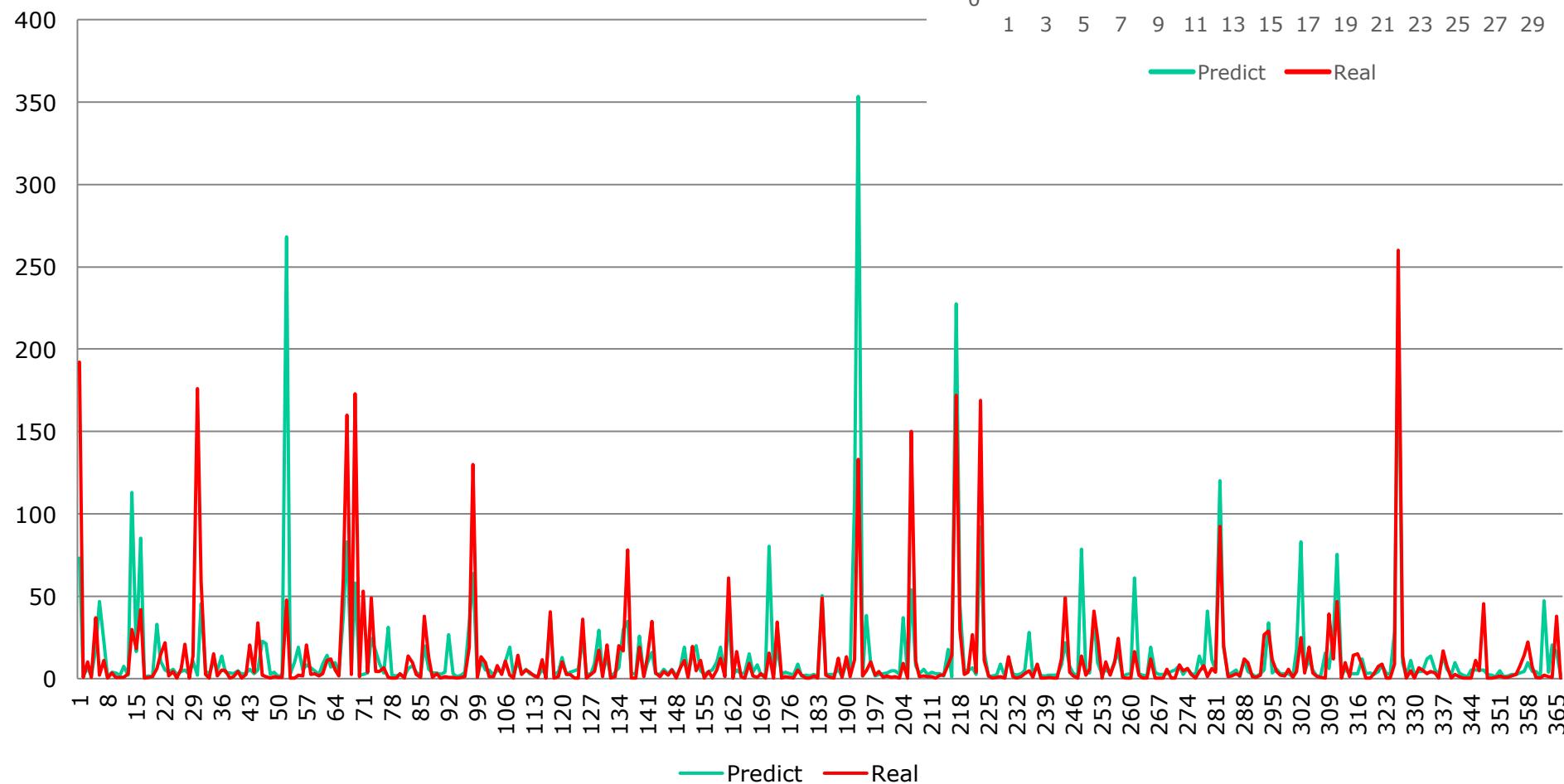
Training with 2 Variables (Copper and Iron)

Model	Linear regression	SVM (rbf)	Random Forest (5 trees)	RNN/LSTM (10 epochs)
Succ. Rate	0.72	0.80	0.84	0.82

Training with 4 Variables

Model	Linear regression	SVM (rbf)	Random Forest (10 trees)	Stacking Azure ML
Succ. Rate	0.69	0.78	0.87	0.89 !

Random Forest



Conclusion

- No strong linear correlation between copper/iron with lead is found.
- Be careful when using water before 9 am and at around 3pm.
- Training with 4 variables (copper, iron, Cl₂ and pH) tends to behave better than with only copper and iron.
- Stacking model has the best performance. Random Forest algorithm comes second.
- RNN/LSTM has great potential but large dataset is needed.

Acknowledge

- Grid Club
 - Great food!
- Mentors
- Sponsors
- Microsoft Azure