

Advances in Statistical Modeling and Mapping of Groundwater Contaminants

Mindy L. Erickson, PhD, PE Hydrologist

U.S. Department of the Interior U.S. Geological Survey

Current USGS work

- Statistical probability mapping of several constituents across Glacial Aquifer System
 - Arsenic
 - Manganese
 - pH
 - Redox/DO
- Complex statistical modeling now possible
 - Machine learning (ML) approaches: Random Forest (RF) and Boosted Regression Tree (BRT)





Statistical predictive modeling

- Advantages of new techniques
 - Does not assume a relationship
 - Can fit nonlinear relations
 - Accommodates predictor variables of any type (e.g. categorical, numeric, binary)
 - Boosting' focuses on the unexplained deviance of previous tree



Motivation

2015 USGS circular about water quality in the US Glacial Aquifer System:

"Contaminants from geologic sources – in particular arsenic and manganese – in groundwater used for drinking are a potential concern for human health"

- MN arsenic and manganese
 - Arsenic \geq 10 µg/L in 10.7% of new potable wells
 - Manganese ≥ 100 µg/L in 50% of ambient wells;
 ≥ 300 µg/L in 22%
 - Unequally distributed across state
 - Redox- and pH-sensitive



Manganese in Minnesota's Groundwaters

Emphasizing the Health Risks of Manganese in Drinking Water Prepared for the Minnesota Ground Water Association

September 2015



Miles

120

0 15 30

60

90



Each dot represents a single well. Wells which were sampled and had less than 2 $\mu g/L$ arsenic are not shown on this map.

Current USGS statistical modeling and mapping work

- Types of variables considered for Machine Learning
- GW travel time: well depth, depth below the water table, recharge, hydrologic position
- Soil drainage: hydrologic groups, soil drainage class, soil texture
- Geology: Surficial geology, thickness of fine grained material, age of bedrock deposits
- Reduction potential: organic matter content
- Other: land use, soil chemistry
 USGS

Online geospatial data sources

SSURGO
USGS
States
Etc.





Current USGS statistical modeling and mapping work

Variable Selection for Machine Learning

- Include all variables (60 80 or more) in initial runs
- Variables ranked by importance
- Remove least important variables
- Final models include the 12 20 most important variables



Current USGS manganese work

- Reasonable performance for Manganese, DO, Arsenic prediction in Central Valley, CA
 - overall accuracy 90%
 - Top predictor variables: related to older groundwater and anoxic conditions
 - Lateral position
 - Depth to water table
 - Sum of poorly-drained soils
 - Average porosity
 - Average organic matter



Boosted Regression Tree Example





Mindy L. Erickson, USGS merickso@usgs.gov

