



## A hybrid machine learning model to predict and visualize nitrate concentration throughout the Central Valley aquifer, California, USA



Katherine M. Ransom<sup>a,\*</sup>, Bernard T. Nolan<sup>b</sup>, Jonathan A. Traum<sup>c</sup>, Claudia C. Faunt<sup>d</sup>, Andrew M. Bell<sup>e</sup>, Jo Ann M. Gronberg<sup>f</sup>, David C. Wheeler<sup>g</sup>, Celia Z. Rosecrans<sup>c</sup>, Bryant Jurgens<sup>c</sup>, Gregory E. Schwarz<sup>b</sup>, Kenneth Belitz<sup>h</sup>, Sandra M. Eberts<sup>i</sup>, George Kourakos<sup>a</sup>, Thomas Harter<sup>a</sup>

<sup>a</sup> University of California, Davis, Department of Land, Air, and Water Resources, United States

<sup>b</sup> U.S. Geological Survey National Water Quality Program, Reston, VA, United States

<sup>c</sup> U.S. Geological Survey California Water Science Center, Sacramento, CA, United States

<sup>d</sup> U.S. Geological Survey California Water Science Center, San Diego, CA, United States

<sup>e</sup> University of California, Davis, Center for Watershed Sciences, United States

<sup>f</sup> U.S. Geological Survey California Water Science Center, Menlo Park, CA, United States

<sup>g</sup> Virginia Commonwealth University, Department of Biostatistics, Richmond, VA, United States

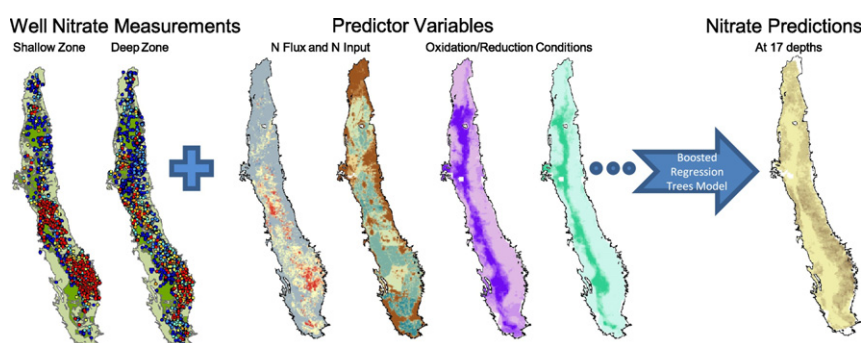
<sup>h</sup> U.S. Geological Survey New England Water Science Center, Northborough, MA, United States

<sup>i</sup> U.S. Geological Survey Ohio Water Science Center, Columbus, OH, United States

### HIGHLIGHTS

- Boosted regression tree model produced 3D map of nitrate concentration.
- Hybrid multi-modeling approach used numerical model outputs as predictors.
- Redox characteristics and field scale unsaturated zone N flux were most important.
- Nitrate concentrations <2 mg/L NO<sub>3</sub>-N generally conformed to basin subregion.
- Nitrate concentrations >10 mg/L NO<sub>3</sub>-N most common in eastern alluvial fans subregion

### GRAPHICAL ABSTRACT



### ARTICLE INFO

#### Article history:

Received 6 March 2017

Received in revised form 19 May 2017

Accepted 20 May 2017

Available online 9 June 2017

Editor: D. Barcelo

#### Keywords:

Groundwater

Nitrate

Boosted regression trees

### ABSTRACT

Intense demand for water in the Central Valley of California and related increases in groundwater nitrate concentration threaten the sustainability of the groundwater resource. To assess contamination risk in the region, we developed a hybrid, non-linear, machine learning model within a statistical learning framework to predict nitrate contamination of groundwater to depths of approximately 500 m below ground surface. A database of 145 predictor variables representing well characteristics, historical and current field and landscape-scale nitrogen mass balances, historical and current land use, oxidation/reduction conditions, groundwater flow, climate, soil characteristics, depth to groundwater, and groundwater age were assigned to over 6000 private supply and public supply wells measured previously for nitrate and located throughout the study area. The boosted regression tree (BRT) method was used to screen and rank variables to predict nitrate concentration at the depths of domestic and public well supplies. The novel approach included as predictor variables outputs from existing physically based models of the Central Valley. The top five most important predictor variables included two oxidation/

\* Corresponding author.

E-mail address: [kmlockhart@ucdavis.edu](mailto:kmlockhart@ucdavis.edu) (K.M. Ransom).