# Accurate prediction of spatial distribution of soil potentially toxic elements using machine learning and associated key influencing factors identification: A case study in mining and smelting area in southwestern China

Kai Li [a,b], Guanghui Guo [a,b,*], Degang Zhang [c], Mei Lei [a,b], Yingying Wang [d]

[a] Institute of Geographic Sciences and Natural Resources Research, CAS, Beijing 100101, China,
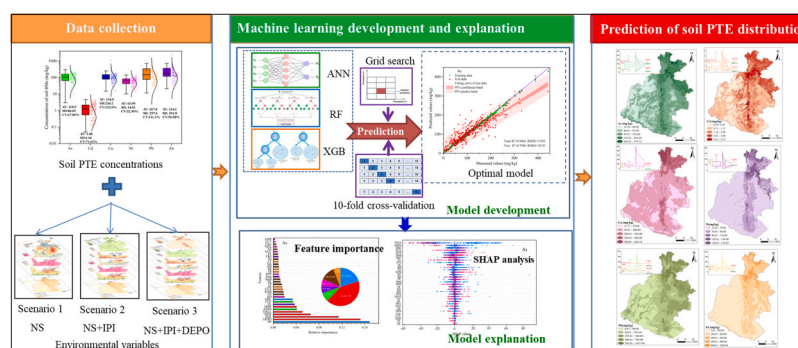[b] University of Chinese Academy of Sciences, Beijing 100049, China
[c] Honghe University, Mengzi 661100, China
[d] Sichuan Eco-environmental Monitoring Station, Chengdu 610091, China

HIGHLIGHTS

- Soil PTE pollution was predicted using ML methods based on environmental covariates.
- Introduction of DEPO can improve the prediction accuracy of soil PTEs.
- XGB had the best performance in predicting As, Cd, Cu, Pb and Zn pollution.
- RF was the best-performing model in predicting Ni pollution.
- Industrial and agricultural activities were the key factors.

GRAPHICAL ABSTRACT