

A Re-evaluation of Paleosol Elemental Proxies for Climate through Cross-validation and Machine Learning

WILLIAM E. LUKENS¹, GARY E. STINCHCOMB², LEE C. NORDT³, STEVEN G. DRIESE⁴, DAVID J. KAHLE⁵, JACK D. TUBBS⁶

¹Geosciences Department, Baylor University, Waco, TX 76798, bill_lukens@baylor.edu

²Watershed Studies Institute & Department of Geosciences, Murray State University, Murray, KY 42071, gstinchcomb@murraystate.edu

³Geosciences Department, Baylor University, Waco, TX 76798, lee_nordt@baylor.edu

⁴Geosciences Department, Baylor University, Waco, TX 76798, steven_driese@baylor.edu

⁵Department of Statistical Science, Baylor University, Waco, TX 76798, david_kahle@baylor.edu

⁶Department of Statistical Science, Baylor University, Waco, TX 76798, jack_tubbs@baylor.edu

The bulk elemental composition of soil subsurface (B) horizons is influenced by environmental, biological, geological, and climatic factors. Because fossil soils (paleosols) are common in the geologic record, quantitative models that link climate to paleosol geochemistry are highly desirable in the paleoclimate community. Over the last ca. 15 years, numerous transfer functions for climate have been developed using B horizon bulk elemental composition. Error for these models is typically reported as the root mean square error (RMSE) associated with regression analysis, and is understood to be the variance imparted by non-climatic influences on soil formation. However, prediction error must be estimated through cross-validation—a method frequently overlooked in simple regression analysis.

Here we re-evaluate the chemical index of alteration minus potassium (CIA-K) proxy for mean annual precipitation, perhaps the most widely applied pedotransfer function for paleosols, by performing cross-validation on two continental-scale soil data sets. Cross-validation of exponential regression models and data set inter-comparisons indicate that previously reported RMSE values underrepresent realistic prediction errors. We introduce promising results from recursive partitioning regression and random forest machine learning, which offer a new direction for automated variable selection and error estimation. These new models offer an optimized balance between wide applicability and relatively low prediction error, and are the necessary next steps in paleosol proxy development.