# A Gradient Boosting Machine for Hierarchically Clustered Data

**Patrick J. Miller**, **Daniel B. McArtor**, and **Gitta H. Lubke**
University of Notre Dame

As increasingly larger data sets are collected (big data), exploratory data analysis to develop predictive models becomes more important but more challenging. Variable selection focuses on identifying important predictors and predictors with nonlinear effects. But simple correlation screens, stepwise regression, or ad-hoc model selection procedures can fail to detect these important nonlinear effects. When data are hierarchically clustered (e.g., observations on students within schools), the problem becomes even more challenging. When mixed effect models are used, model selection can produce misleading results when nonlinear effects are not included in the model. Additionally, mixed effect models can be difficult to estimate when there are a large number of predictors with group specific effects.

'Boosted decision trees' is an alternative approach to parametric modeling. It is an additive model of decision trees estimated by gradient descent (Friedman, 2001). This flexible approach can approximate complex and nonlinear functions of predictors without specifying their functional form beforehand; boasted decision trees are easily estimated (e.g., using the R package 'gbm', Ridgeway et al., 2015). However, it is unclear how to best use boosted decision trees with hierarchically clustered data. Including the clustering variable in the model makes the interpretation more difficult. Ignoring the clustering variable misses an opportunity for potentially improved prediction, potentially leading to biased variable selection.

We propose an extension to boosted decision trees that works by allowing the terminal node means of each tree in the model to vary by group at each iteration. We implement the approach using the popular R package 'lme4' (Bates et al., 2014). The method can be used with thousands of predictors, and allows the predictors to contain missing values. Initial results show that this procedure can improve prediction and variable selection performance by as much as 25% compared to including the grouping variable as a candidate for splitting in each tree.

Correspondence concerning this abstract should be addressed to Patrick J. Miller, 110 Haggar Hall, University of Notre Dame, Notre Dame, IN 46656. pmille13@nd.edu.

## Acknowledgments

## References

Friedman J. Greedy function approximation: a gradient boosting machine. Annals of Statistics. 2001; 29(5):1189–1232.

Ridgeway, G. with contributions from others. gbm: Generalized Boosted Regression Models. R package version 2.1.1. 2015. https://CRAN.R-project.org/package=gbm

Bates D, Maechler M, Bolker B, Walker S. lme4: Linear mixed-effects models using Eigen and S4. R package version. 2014; 1(7)