

Paranormal Distribution

Andy Chung

Kevin Gilbert

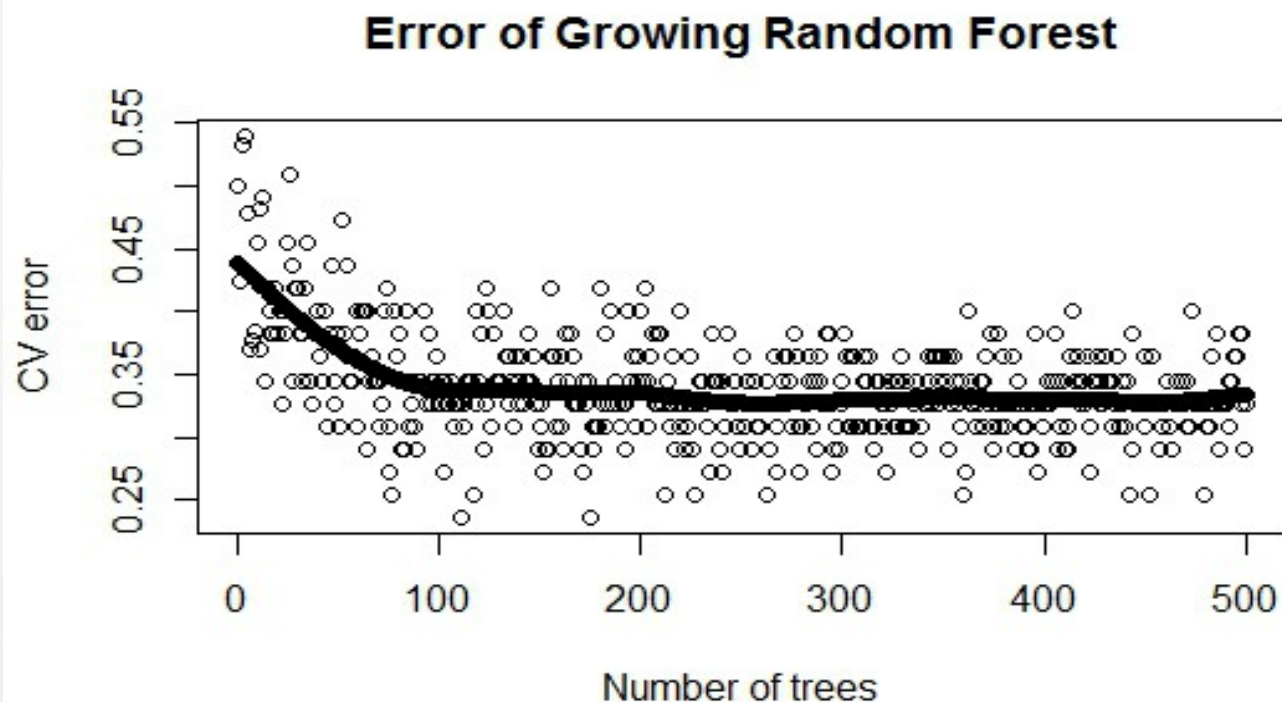
Jason Sun

First Step: Random Forest

One decision tree: One classification

500 decision trees: One classification

We bag these trees into a “random forest”

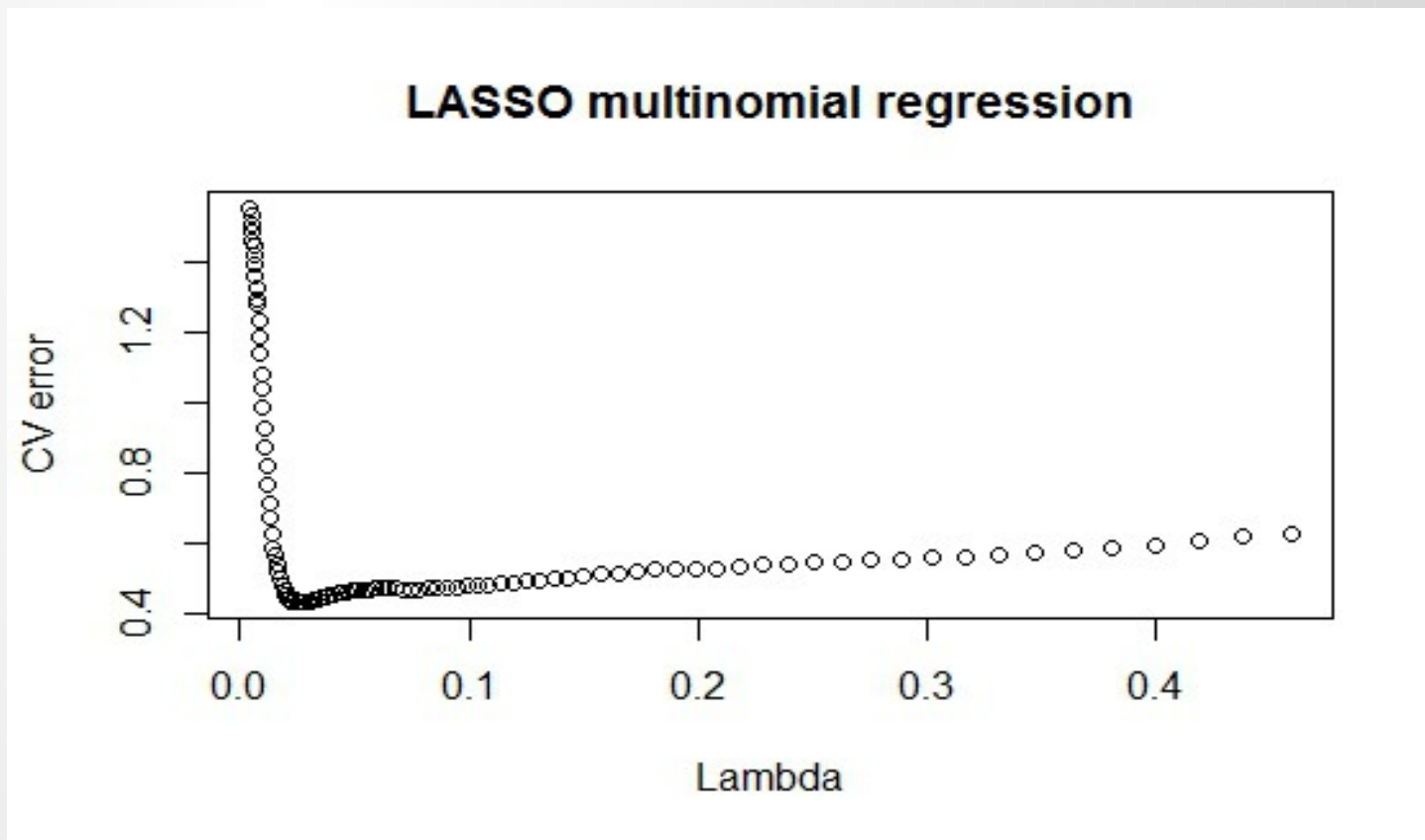


Second Step: LASSO (glmnet)

What if Random Forest isn't sure?

Recalculate unsure nodes with LASSO, ridge

Replacing decisions reduced our CV error



Results on training data

Random Forest prone to overfit

Individual forests vary between 25% and 40% error

Two-stage algorithm has only 30% error

Tested on 4-folds cross-validation

We're pretty sure map of neighborhoods is... Pittsburgh!

Scaled map of pair.dist

Pittsburgh, Pennsylvania

Neighborhoods

