

Schonlau M. **Boosted Regression (Boosting): An introductory tutorial and a Stata plugin.** *The Stata Journal*, 5(3), 330-354.

Boosted Regression (Boosting): An introductory tutorial and a Stata plugin

Matthias Schonlau
RAND

Abstract

Boosting, or boosted regression, is a recent data mining technique that has shown considerable success in predictive accuracy. This article gives an overview over boosting and introduces a new Stata command, *boost*, that implements the boosting algorithm described in Hastie et al. (2001, p. 322). The plugin is illustrated with a Gaussian and a logistic regression example. In the Gaussian regression example the R^2 value computed on a test data set is $R^2=21.3\%$ for linear regression and $R^2=93.8\%$ for boosting. In the logistic regression example stepwise logistic regression correctly classifies 54.1% of the observations in a test data set versus 76.0% for boosted logistic regression. Currently, *boost* accommodates Gaussian (normal), logistic, and Poisson boosted regression. *boost* is implemented as a Windows C++ plugin.

1 Introduction

Economists and analysts in the data mining community differ in their approach to regression analysis. Economists often build a model from theory and then use the data to estimate parameters of their model. Because their model is justified by theory economists are sometimes less inclined to test how well the model fits the data. Data miners tend to favor the “kitchen sink” approach in which all or most available regressors are used in the regression model. Because the choice of x-variables is not supported by theory, validation of the regression model is very important. The standard approach to validating models in data mining is to split the data into a training and a test data set.

The concept of training versus test data set is central to many data mining algorithms. The model is fit on the training data. The fitted model is then used to make predictions on the test data. Assessing the model on a test data rather on the training data ensures that the model is not overfit and is generalizable. If the regression model has tuning parameters (e.g. ridge regression, neural networks, boosting), good values for the tuning parameters are usually found running by the model several times with different values for the tuning parameters. The performance of each model is assessed on the test data set and the best model (according to some criterion) is chosen.

In this paper I review boosting or boosted regression and supply a Stata plugin for Windows. In the same way that generalized linear models include Gaussian, logistic and other regressions, boosting also includes boosted versions of Gaussian, logistic and other regressions. Boosting is a highly flexible regression method. It allows the researcher to specify the x-variables without specifying the functional relationship to the response. Traditionally, data miners have used boosting in the context of the “kitchen sink” approach to regression but it is also possible to use boosting in a more targeted manner, i.e. using only variables motivated by theory. Because it is more flexible, a boosted model will tend to fit better than a linear model and therefore inferences made based on the model may have more credibility.

There is mounting empirical evidence that boosting is one of the best modeling approaches ever developed. Bauer and Kohavi (1999) performed an extensive comparison of boosting to several other competitors on 14 datasets, and found boosting to

be “the best algorithm”. Friedman et al. (2000) compare several boosting variants to the CART (classification and regression tree) method and find that all the boosting variants outperform the CART algorithm on eight datasets.

The success of boosting in terms of predictive accuracy has been subject to much speculation. Part of the mystery that surrounds boosting is due to the fact that different scientific communities, computer scientists and statisticians, have contributed to its development. I will first give a glimpse into the fascinating history of boosting.

The remainder of the paper is structured as follows: Section 2 contains the syntax of the *boost* command. Section 3 contains options for that command. Section 4 explains boosting and gives a historical overview. Section 5 uses a toy example to show how the *boost* command can be used for normal (Gaussian) regression. Section 6 features a logistic regression example. Section 7 gives runtime benchmarks for the boosting command. Section 8 concludes with some discussion.

2 Installation and Syntax

To install the *boost* plugin copy the *boost.hlp* and *boost.ado* files in one of the *ado* directories, e.g. `c:\ado\personal\boost.ado` . A list of valid directories can be obtained by typing “*adopath*” within Stata. Copy the “*boost.dll*” file into a directory of your choosing (e.g. the same directory as the *ado* file). Unlike an *ado* file, a plugin has to be explicitly loaded:

```
capture program boost_plugin, plugin using("C:\ado\personal\boost.dll")
```

The command "capture" prevents this line resulting into an error in case the plugin was already loaded. The command syntax is as follows:

```
boost varlist [if exp] [in range] , DISTribution(string) maxiter(int)
      [ INfluence PREDict(varname) shrink(real=0.01)
      bag(real=0.5) INTERaction(int=5) seed(int=0) ]
```

3 Options

boost determines the number of iterations that maximizes the likelihood, or, equivalently, the pseudo R squared. The pseudo R2 is defined as $R2=1-L1/L0$ where L1 and L0 are the log likelihood of the full model and intercept-only model respectively.

Unlike the R2 given in *regress*, the pseudo R2 is an out-of-sample statistic. Out-of-sample R2's tend to be lower than in-sample-R2's.

Output and Return values

The standard output consists of the best number of iterations, *bestiter*, the R squared value computed on the test data set, *test_R2*, and the number of observations used for the training data, *trainn*. *trainn* is computed as the number of observations that meet the in/if conditions times *trainfraction*. These statistics can also be retrieved using *ereturn*. In addition, *ereturn* also stores the training R squared value, *train_R2*, as well as the log likelihood values from which *train_R2* and *test_R2* are computed.

distribution(string) Currently, possible distributions are "normal", "logistic", and "poisson".

influence displays the percentage of variation explained (for non-normal distributions: percentage of log likelihood explained) by each input variable. The influence matrix is saved in *e(influence)*.

predict(varname) predicts and saves the predictions in the variable *varname*. To allow for out-of-sample predictions *predict* ignores *if* and *in*. For model fitting only observations that satisfy *if* and *in* are used, predictions are made for all observations.

trainfraction(int) Specifies the percentage of data to be used as training data. The remainder, the test data is used to evaluate the best number of iterations. By default this value is 0.8.

interaction(int) specifies the maximum number of interactions allowed. *interaction=1* means that only main effects are fit, *interaction=2* means that main effect and two way interactions are fitted, and so forth. The number of interactions equals the number of terminal nodes in a tree plus 1. If *interaction=1*, then each tree has 2 terminal nodes. If

interaction=2, then each tree has 3 terminal nodes, and so forth. By default *interaction=5*.

maxiter(int) specifies the maximal number of trees to be fitted. The actual number used, *bestiter*, can be obtained from the output as *e(bestiter)*. When *bestiter* is too close to *maxiter* the maximum likelihood iteration may be larger than *maxiter*. In that case it is useful to rerun the model with a larger value for *maxiter*. When *trainfraction=1.0* all *maxiter* observations are used for prediction (*bestiter* is missing because it is computed on a test data set).

shrink(#) specifies the shrinkage factor. *shrink=1* corresponds to no shrinkage. As a general rule of thumb, reducing the value for *shrink* requires an increase in the value of *maxiter* to achieve a comparable cross validation R². By default *shrink= 0.01*.

bag(real) Specifies the fraction of training observations that is used to fit an individual tree. *bag=0.5* means that half the observations are used for building each tree. To use all observations specify *bag=1.0*. By default *bag=0.5*.

seed(int) *seed* specifies the random number seed to generate the same sequence of random numbers. Random numbers are only used for bagging. Bagging uses random numbers to select a random subset of the observations for each iteration. By default (*seed=0*). The boost seed option is unrelated to Stata's *set seed* command.

4 Boosting

Boosting was invented by computational learning theorists and later reinterpreted and generalized by statisticians and machine learning researchers. Computer scientists tend to think of boosting as an “ensemble” method (a weighted average of predictions of individual classifiers), whereas statisticians tend to think of boosting as a sequential regression method. To understand why statisticians and computer scientists think about the essentially same algorithms in different ways, both approaches are discussed. Section 4.1 discusses an early boosting algorithm from computer science. Section 4.2 describes regression trees, the most commonly used base learner in boosting. Section 4.3 describes Friedman’s gradient boosting algorithm, which is the algorithm I have implemented for

the Stata plugin. The remaining sections talk about variations of the algorithm that are relevant to my implementation (Section 4.4), how to evaluate boosting algorithms via a cross validated R^2 (Section 4.5), the influence of variables (Section 4.6) and advice on how to set the boosting parameters in practice (Section 4.7).

4.1 Boosting and its roots in computer science

Boosting was invented by two computer scientists at AT&T Labs (Freund and Schapire, 1997). Below I describe an early algorithm, the “AdaBoost” algorithm, because it illustrates why computer scientists think of boosting as an ensemble method; that is, a method that averages over multiple classifiers.

Adaboost (see Algorithm 1) works only in the case where the response variable takes only one of two values: -1 and 1. (Whether the values are 0/1 or -1/1 is not important- the algorithm could be modified easily). Let C_1 be a binary classifier (e.g. logistic regression) that predicts whether an observation belongs to the class “-1” or “1”. The classifier is fit to the data as usual and the misclassification rate is computed. This first classifier C_1 receives a classifier weight that is a monotone function of the error rate it attains. In addition to classifier weights there are also observation weights. For the first classifier, all observations were weighted equally. The second classifier, C_2 (e.g.. the second logistic regression), is fit to the same data, however with changed observation weights. Observation weights corresponding to observations misclassified by the previous classifier are increased. Again, observations are reweighted, a third classifier C_3 (e.g. a third logistic regression) is fit and so forth. Altogether *iter* classifiers are fit where *iter* is some predetermined constant. Finally, using the classifier weights the classifications of the individual classifiers are combined by taking a weighted majority vote. The algorithm is described in more detail in Algorithm 1.

Initialize weights to be equal $w_i = 1/n$

For $m = 1$ to *iter* classifiers C_m):

- (a) Fit classifier C_m to the weighted data
- (b) Compute the (weighted) misclassification rate r_m
- (c) Let the classifier weight $\alpha_m = \log((1 - r_m)/r_m)$
- (d) Recalculate weights $w_i = w_i \exp(\alpha_m \mathbf{I}(y_i \neq C_m))$

Majority vote classification: $\text{sign} \left[\sum_{m=1}^M \alpha_m C_m(x) \right]$

Algorithm 1: The AdaBoost algorithm for classification into two categories.

Early on, researchers attributed the success and innovative element of this algorithm to the fact that observations that are repeatedly misclassified are given successively larger weights. Another unusual element is that the final “boosted” classifier consists of a majority-weighted vote of all previous classifiers. With this algorithm, the individual classifiers do not need to be particularly complex. On the contrary, simple classifiers tend to work best.

4.2 Regression trees

The most commonly used simple classifier is a regression tree (e.g. CART, Breiman et al.,1984). A regression tree partitions the space of input variables into rectangles and then fits a constant (e.g. estimated by an average or a percentage) to each rectangle. The partitions can be described by a series of if-then statements or they can be visualized by a graph that looks like a tree. Figure 1 gives an example of such a tree using age and indicator variables for race/ethnicity. The tree has 6 splits and therefore 7 leaves (terminal nodes). Each split represents an if-then condition. Each observation descends the set of if-then conditions until a leaf is reached. If the if-then condition is true then the observation goes to the right branch otherwise to the left branch. For example, in Figure 2 any observation over 61 years of age would be classified in the right most leaf regardless

of his/her race/ethnicity. The leaf with the largest values of the response, 0.8, corresponds to observations between 49 and 57 years of age, who are neither Hispanic nor Black. As the splits on age show, the tree can split on the same variable several times. If the 6 splits in the tree in Figure 1 had split on six different variables the tree would represent a 6-level interaction because six variables would have to be considered jointly in order to obtain the predicted value.

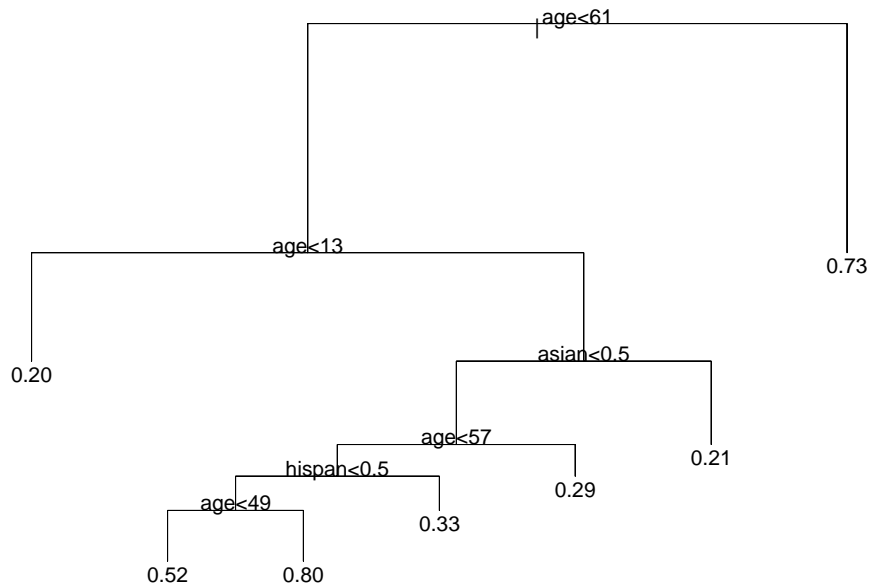


Figure 1: Example of a regression tree with 6 splits

A regression tree with only two terminal nodes (i.e., a tree with only one split) is called a tree stump. It is hard to imagine a simpler classifier than a tree stump – yet tree stumps work surprisingly well in boosting. Boosting with stumps fits an additive model and many datasets are well approximated with only additive effects. Fewer, complex classifiers can be more flexible and harbor a greater danger of overfitting (i.e., fitting well only to the training data).

4.3 Friedman's gradient boosting algorithm

Early researchers had some difficulty explaining the success of the AdaBoost algorithm. The computational learning theory roots of the AdaBoost puzzled statisticians who have traditionally worked from likelihood-based approaches to classification and, more generally, to regression. Then Friedman et al. (2000) were able to reinterpret this algorithm in a likelihood framework, enabling the authors to form a boosted logistic regression algorithm, a formulation more familiar to the statistical community. Once the connection to the likelihood existed, boosting could be extended to generalized linear models and further still to practically any loss criterion (Friedman 2001, Ridgeway, 1999). This meant that a boosting algorithm could be developed for all the error distributions in common practice: Gaussian, logistic, Poisson, Cox models, etc. With the publication of Hastie et al. (2001), a book on statistical learning, modern regression methods like boosting caught on in the statistics community.

The interpretation of boosting in terms of regression for a continuous, normally distributed response variable is as follows: The average y -value is used as a first guess for predicting all observations. This is analogous to fitting a linear regression model that consists of the intercept only. The residuals from the model are computed. A regression tree is fit to the residuals. For each terminal node, the average y -value of the residuals that the node contains is computed. The regression tree is used to predict the residuals. (In the first step, this means that a regression tree is fit to the difference between the observation and the average y -value. The tree then predicts those differences.) The boosting regression model - consisting of the sum of all previous regression trees - is updated to reflect the current regression tree. The residuals are updated to reflect the changes in the boosting regression model; a tree is fit to the new residuals, and so forth. This algorithm is summarized in more detail in Algorithm 2.

-
- 1) Initialization: Set initial guess to \bar{y}
 - 2) For all regressions trees $m=1$ to M :
 - 2a) Compute the residuals based on the current model

$$r_{mi} = y_i - f_{m-1}(x_i)$$

where i indexes observations.

Note that f_{m-1} refers to the sum of all previous regression trees.

2b) Fit a regression tree (with a fixed number of nodes) to the residuals

2c) For each terminal node of the tree, compute the average residual. The average value is the estimate for residuals that fall in the corresponding node.

2d) Add the regression tree of the residuals to the current best fit

$$f_m = f_{m-1} + \text{last regression tree of residuals}$$

Algorithm 2: Friedman's gradient boosting algorithm for a normally distributed response

Each term of the regression model thus consists of a tree. Each tree fits the residuals of the prediction of all previous trees combined. To generalize Algorithm 2 to a more general version with arbitrary distributions requires that "average y-value" be replaced with a function of y-values that is dictated by the specific distribution and that the residual (Step 2a) be a "deviance residual". (McCullagh and Nelder, 1989)

The basic boosting algorithm requires the specification of two parameters. One is the number of splits (or the number of nodes) that are used for fitting each regression tree in step 2b of Algorithm 2. The number of nodes equals the number of splits plus one. Specifying one split (tree stumps) corresponds to an additive model with only main effects. Specifying two splits corresponds to a model with main effects and two-way interactions. Generally, specifying J splits corresponds to a model with up to J -way interactions. When J x-variables need to be considered jointly for a component of a regression model this is a J -way interaction. Hastie et al. (2001) suggest that $J = 2$ in general is not sufficient and that $4 \leq J \leq 8$ generally works well. They further suggest that the model is typically not sensitive to the exact choice of J within that range. In the Stata implementation J is specified as an option: *interaction(J)*.

The second parameter is the number of iterations or the number of trees to be fit. If the number of iterations is too large the model will overfit; i.e., it will fit the training data well but not generalize to other observations from the same population. If the number of iterations is too small then the model is not fit as well either. A suitable value for *maxiter* can range from a few dozen to several thousand, depending on the value of a shrinkage parameter (explained below) and the data set. The easiest way to find a suitable number of iterations is to check how well the model fits on a test data set. In the Stata implementation the maximal number of iterations, *maxiter*, is specified and the number of iteration that maximizes the log likelihood on a test data set, *bestiter*, is automatically found. The size of the test data set is controlled by *trainfraction*. For example, if *trainfraction*=0.5 then the last 50% of the data are used as test data.

4.4 Shrinkage and bagging

There are two commonly used variations on Friedman's boosting algorithm: "shrinkage" and "bagging". Shrinkage (or regularization) means reducing or shrinking the impact of each additional tree in an effort to avoid overfitting. The intuition behind this idea is that it is better to improve a model by taking many small steps than a smaller number of large steps. If one step turns out to be a misstep, then the damage can be more easily undone in subsequent steps. Shrinkage has been previously employed, for example, in ridge regression where it refers to shrinking regression coefficients back to zero to reduce the impact of unstable regression coefficients on the model. Shrinkage is accomplished by introducing a parameter λ in step 2d of Algorithm 2:

$$f_m = f_{m-1} + \lambda * (\text{last regression tree of residuals})$$

where $0 < \lambda \leq 1$. The smaller λ , the greater the shrinkage. The value $\lambda = 1$ corresponds to no shrinkage. Typically λ is 0.1 or smaller with $\lambda = 0.01$ or $\lambda = 0.001$ being common. Smaller shrinkage values require larger number of iterations. In my experience, λ and the number of iterations typically satisfy $10 \leq \lambda * \text{bestiter} \leq 100$ for the best model. In other words, a decrease of λ by a factor of 10 implies an increase of the number of iterations by a similar factor. In the Stata implementation λ is specified through the option *shrink*(λ).

The second commonly used variation on the boosting algorithm is bagging. At each iteration only a random fraction, *bag*, of the residuals is selected. In that iteration,

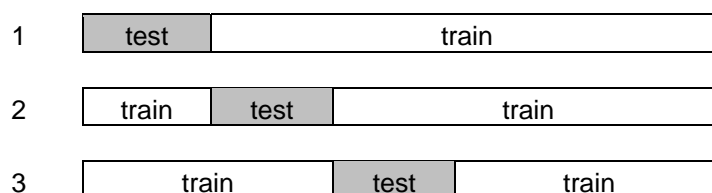
only the random subset of the residuals is used to build the tree. Non-selected residuals are not used in that iteration at all. The randomization is thought to reduce the variation of the final prediction without affecting bias. Different random subsets may have different idiosyncrasies that will average out. While not all observations are used in each iteration, all observations are eventually used across all iterations. Friedman (2002) recommends the bagging with 50% of the data. In the Stata implementation this can be specified as *bag(0.5)*.

4.5 Crossvalidation and the pseudo R^2

Highly flexible models are prone to overfitting. While it may still occur, overfitting is less of an issue in linear regression where the restriction of linearity guards against this problem to some extent. To assess predictive accuracy with highly flexible models it is important to separate the data the model was trained on from test data.

My boosting implementation splits the data into a training and a test data set. By default the first 80% of the data are used as training data and the remainder as test data. This percentage can be changed through the use of the option *trainfraction*. If *trainfraction=0.5*, for example, then the first 50% of the data are used as training data. It is important that the observations are in random order before invoking the *boost* command because otherwise the test data may be different from the training data in a systematic way.

Crossvalidation is a generalization of the idea of splitting the data into training and test data sets. For example, in five-fold crossvalidation the data set is split into 5 distinct subsets of 20% of the data. In turn each subset is used as test data and the remainder as training data. This is illustrated graphically in Figure 2. My implementation corresponds to the fifth row of Figure 2. An example of how to rotate the 5 groups of data to accomplish five-fold crossvalidation using the *boost* implementation is given in the help file for the *boost* command.



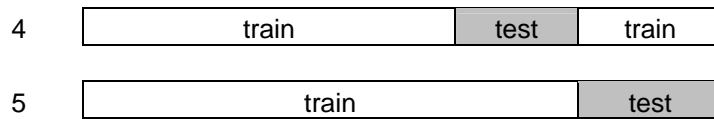


Figure 2: Illustration of five-fold cross validation: Data are split into training and test data in five different ways.

For each test data set a pseudo R^2 is computed on the test data set. The pseudo R^2 is defined as

$$\text{pseudo } R^2 = 1 - L1/L0$$

where $L1$ and $L0$ are the log likelihoods of the model under consideration and an intercept-only model, respectively (see Long and Freese, 2003, for a discussion on pseudo R^2). In the case of Gaussian (normal) regression, the pseudo R^2 turns into the familiar R^2 that can be interpreted as “fraction of variance explained by the model”. For Gaussian regression it is sometimes convenient to compute R^2 as

$$R^2 = \frac{\text{Var}(y) - \text{MSE}(y, \hat{y})}{\text{Var}(y)} \quad (1)$$

where $\text{Var}(\cdot)$ and $\text{MSE}(\cdot, \cdot)$ refer to variance and mean squared error respectively. To avoid cumbersome language this is referred to as “test R^2 ” in the remainder of this paper.

The training R^2 and test R^2 are computed on training and test data, respectively. The training R^2 is always between 0 and 1. Usually, the test R^2 is also between 0 and 1. However, if the log likelihood based on the intercept-only model is greater than the one for the model under consideration the test R^2 is negative. A negative test R^2 is a sign that the model is strongly overfit. The *boost* command computes both training and test R^2 values and makes them available in `e(test_R2)` and `e(train_R2)`.

4.6 Influence of variables and visualization

For a linear regression model the effect of x -variables on the response variable is summarized by their respective coefficients. The boosting model is complex but can be interpreted with the right tools. Instead of regression coefficients, boosting has the concept of “influence of variables” (Friedman, 2001). Each split on a variable in a regression tree increases the log likelihood of the regression tree model. In Gaussian regression the increase in log likelihood is proportional to the increase in sums of squares

explained by the model. The sum of log likelihood increases across all trees due to a given variable yields the influence of that variable. For example, suppose there are 1000 iterations or regression trees. Suppose that each regression tree has one split (i.e. a tree stump, meaning *interaction=1*, meaning main effects only). Further, suppose that 300 of the regression trees split on variable x_i . Then the sum of the increase in log likelihood due to these 300 trees represents the influence of x_i . The influences are standardized such that they add up to 100%.

Because a regression tree does not separate main effects and interactions, influences are only defined for variables - not for individual main effect or interaction terms. Influences only reveal the sum of squares explained by the variables; they say nothing about how the variable affects the response. The functional form of a variable is usually explored through visualization. Visualization of any one (or a any set of) variables is achieved by predicting over a range or grid of these variables. Suppose the effect of x_i on the response is of interest. There are several options. The easiest is to predict the response for a suitable range of x_i values while holding all other variables constant (for example at their mean or their mode). A second option is to predict the response for a suitable range of x_i values for each observation in the sample. Then the sample-averaged prediction for a give x_i value can be computed. A third option is to numerically integrate out other variables with a suitable grid of values to get what is usually referred to as a main effect.

In linear regression the coefficient of a variable need to be interpreted in the context of its range. It is possible to artificially inflate coefficients by, for example, changing a measurement from kilogram to gram. The influence percentages are invariant to any one-to-one rescaling of the variables. For example, measuring a variable in miles or kilometers does not affect the influence but it would affect the regression coefficient of a logistic regression. In the Stata implementation the influence of individual variables is displayed when option *influence* is specified. The individual values are stored in a Stata matrix and can be accessed through `e(influence)`. The help file gives an example.

4.7 Advice on setting boosting parameters

In what follows I give some suggestions of how the user might think about setting parameter values.

Maxiter: If $e(\text{bestiter})$ is almost as large as maxiter consider rerunning the model with a larger value for maxiter . The iteration corresponding to the maximum likelihood estimate may be greater than maxiter .

Interactions: Usually a value between 3 and 7 works well. If you know your model only has two-way interactions I still suggest trying larger values also for two reasons. One, typically, the loss in test R^2 for specifying too low a number of interactions is far worse than the one of specifying a larger number of interactions. The performance of the model only deteriorates noticeably when the number of interactions is much too large. Two, specifying 5-way interactions allows each tree to have 5 splits. These splits are not necessarily always used for an interaction (splitting on 5 different variables in the same tree). They could also be used for a nonlinearity in a single variable (5 splits on different values for a single variable) or for combinations of nonlinearities and interactions.

Shrinkage: Lower shrinkage values usually improve the test R^2 but they increase the running time dramatically. Shrinkage can be thought of as a step size. The smaller the step size, the more iterations and computing time are needed. In practice I choose a small shrinkage value such that the command execution does not take too much time. Because I am impatient, I usually start out with $\text{shrink}=0.1$. If time permits switch to 0.01 or lower in a final run. There are diminishing returns for very low shrinkage values.

Bagging: Bagging may improve the R^2 and it also makes the estimate less variable. Typically, the exact value is not crucial. I typically use a value in the range of 0.4-0.8.

5 Example: Boosted Gaussian regression

This section gives a simple example with four explanatory variables constructed to illustrate how to perform and evaluate boosted regressions. The results are also

compared to linear regression. Linear regression is a good reference model because many scientists initially fit a linear model. I simulated data from the following model

$$y = 30(x_1 - 0.5)^2 + 2x_2^{-0.5} + x_3 + \varepsilon$$

where $\varepsilon \sim \text{uniform}(0,1)$ and $0 \leq x_i \leq 1$ for $i \in 1, 2, 3$. To keep things relatively simple, the model has been chosen to be additive without interactions. It is quadratic in x_1 , nonlinear in x_2 , linear with a small slope in x_3 . The nonlinear contribution of x_2 is stronger than the linear contribution of x_3 even though their slopes are similar. A fourth variable, x_4 , is unrelated to the response but is used in the analysis in an attempt to confuse the boosting algorithm. Scatter plots of y vs x_1 through x_4 is shown in Figure 3.

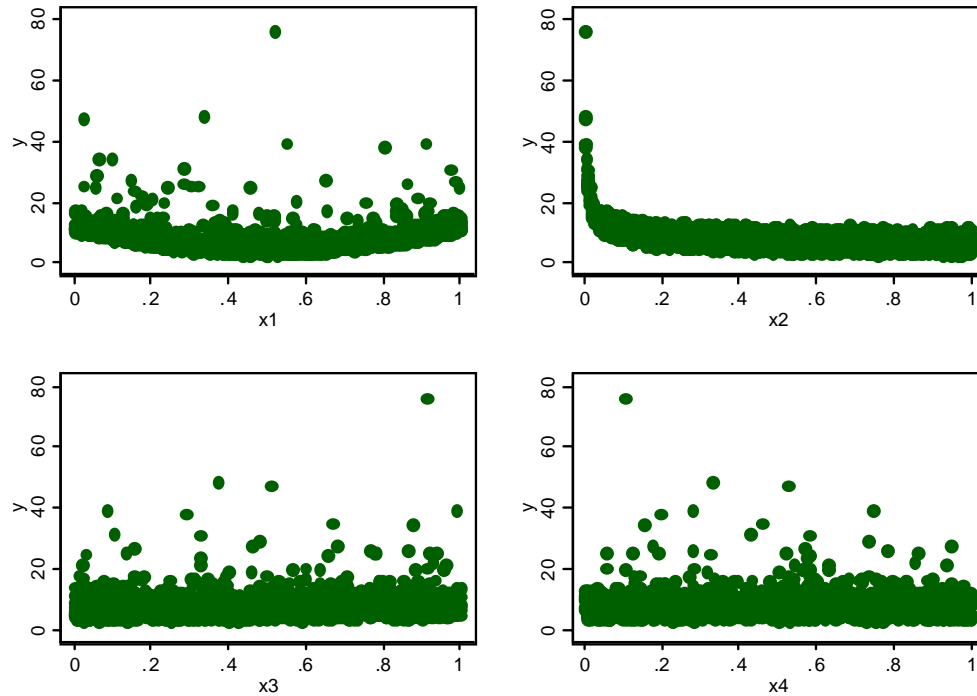


Figure 3: Scatter plots of y versus x_1 through x_4

I chose the following parameter values $shrink=0.01$, $bag=0.5$ and $maxiter=4000$. My rule of thumb in choosing the maximal number of iterations is that the shrinkage factor times the maximal number of iterations should be roughly between 10 and 100. In

my experience the cross-validated R^2 as a function of the number of iterations is unimodal, i.e. there is only one maximum. If *bestiter* is too close to *maxiter* then the number of iterations that maximizes the likelihood may be greater than *maxiter*. It is recommended to rerun the command with a larger value for *maxiter*. The command I am giving is:

```
boost y x1-x4, distribution(normal) train(0.5) maxiter(4000) seed(1)
      bag(0.5) interaction(`inter`) shrink(0.01)
```

where the commands only differ by *inter* ranging from 1 through 5. One of these *boost* commands runs in 8.8 seconds on my laptop. Fixing the seed is only relevant for bagging. Figure 4 shows a plot of the test R^2 versus the number of interactions. The test R^2 is roughly the same regardless of the number of interactions (note the scale of the plot). The fact that the test R^2 is high even for the main effect model (*interaction=1*) does not surprise because our model did not contain any interactions. The actual number of iterations that maximizes the likelihood, *bestiter*, varies. Here the number of iterations are (number of interactions in parenthesis): 3769 (1), 2171 (2), 2401 (3), 1659 (4), 1156 (5).

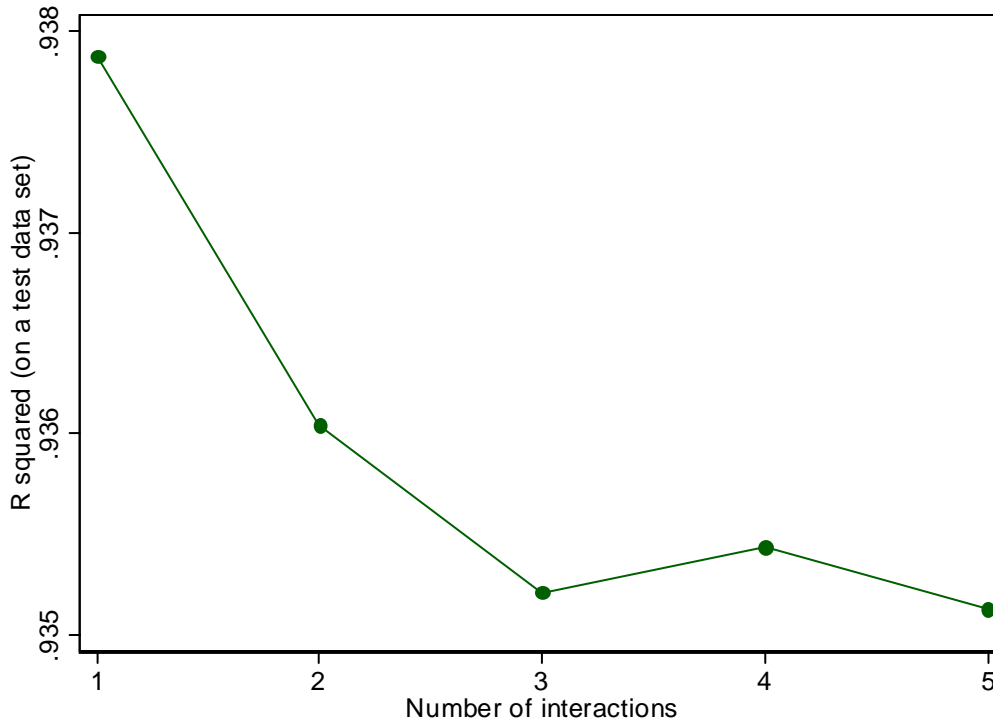


Figure 4: Scatter plot of the R^2 (computed on a test data set) versus the number of interactions. Note the scale on the vertical axis.

I often want to confirm that this model indeed works better than linear regression. Directly comparing the R^2 value for the boosted regression and the linear regression is not a fair comparison. The boosted regression R^2 refers is computed on a test data set, whereas the linear regression R^2 is computed on the training data set. Using equation (1) it is possible to compute an R^2 on a test data set for the linear regression. For the following set of Stata commands I assume that the first *trainn* observations in the data set constitute the training data and the remainder the test data. The predictions from the linear regression (or any other predictions) are denoted *regress_pred*.

```

global trainn=e(trainn) /* using e() from boost */
regress y x1 x2 x3 x4 in 1/$trainn
predict regress_pred
* compute Rsquared on test data
gen regress_eps=y-regress_pred
gen regress_eps2= regress_eps*regress_eps
replace regress_eps2=0 if _n<=$trainn
gen regress_ss=sum(regress_eps2)
local mse=regress_ss[_N] / (_N-$trainn)
sum y if _n>$trainn
local var=r(Var)
local regress_r2= (`var'-`mse')/^var'
di "mse=" `mse' " var=" `var' " regress r2=" `regress_r2'

```

The test R^2 will usually be lower than the R^2 in the output that Stata displays – but because of variability it may be larger on occasion. From the Stata output of “regress” the (training) R^2 value is $R^2= 24.1\%$. The test R^2 computed from the above set of stat commands is $R^2=21.3\%$. I compute the boosting predictions and the influences of the variables:

```

boost y x1 x2 x3 x4, distribution(normal) train(0.5) bag(0.5) maxiter(4000)
interaction(1) shrink(0.01) pred("boost_pred") influence seed(1)

```

(2)

Substituting the boosting predictions for the linear regression predictions in the above set of Stata commands the test boosting R^2 turns out to be $R^2= 93.8\%$. This is the same value as the test R^2 displayed in Figure 4. Figure 5 displays actual versus fitted y-values for both linear regression and boosting. The boosting model fits much better.

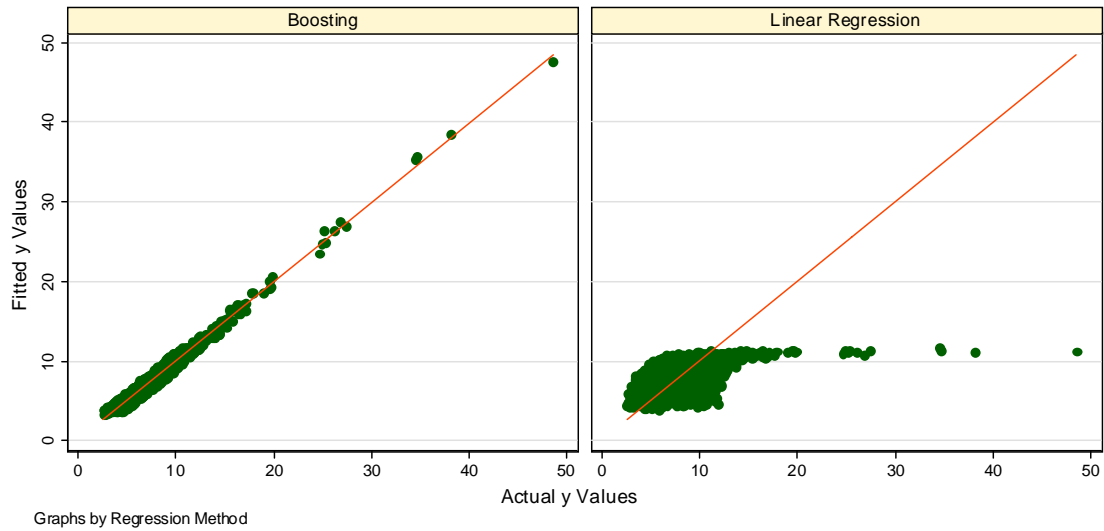


Figure 5: A Calibration plot for the linear regression example: Fitted versus actual values for both linear regression and boosting. The Boosting model is much better calibrated.

When working with linear regressions it is common to look at coefficients to assess how a variable affects the outcome. In boosted regression one looks at the influence of all variables. The influences are given in percentages. Specifying influence in the Stata command listed in equation (2) gives the following output:

x1	30.9%
x2	68.3%
x3	0.67%
x4	0.08%

Table 1: Influence of each variable (Percent)

Variables x_2 and x_1 are most influential. The other variables have almost no influence. Given that there is weak relation between y and x_3 and no relation between y and x_4 , it is nice to see that the influence of x_3 is larger than that of x_4 .

The influence shows one can learn how large the effect of individual variables is but not the functional form. To visualize the conditional effect of a variable x_1 : all variables except x_1 are set to a fixed value (here 0.5). For x_1 , values that cover its range are chosen. The new observations are fed to the model for predicting the response, and the predicted response is plotted against x_1 . This can be accomplished as follows:

```
drop if _n>1000
set obs 1400
replace x1=0.5 if _n>1000
replace x2=0.5 if _n>1000
replace x3=0.5 if _n>1000
replace x4=0.5 if _n>1000
replace x1= (_n-1000)/100 if _n>1000 & _n<=1100
replace x2= (_n-1100)/100 if _n>1100 & _n<=1200
replace x3= (_n-1200)/100 if _n>1200 & _n<=1300
replace x4= (_n-1300)/100 if _n>1300 & _n<=1400
boost y x1 x2 x3 x4 in 1/1000 , distribution(normal) /*
*/ maxiter(4000) bag(0.5) interaction(1) shrink(0.01) /*
*/ pred("pred")
line pred x1 if _n>1000 & _n<=1100
line pred x2 if _n>1100 & _n<=1200
line pred x3 if _n>1200 & _n<=1300
line pred x4 if _n>1300 & _n<=1400
```

The *boost* command uses the first 1000 observations to fit the model but uses all observations for prediction. Figure 6 displays all four conditional effects. All effects are step functions because the base learners, regression trees, can only produce step functions. The features of the smooth curves are well reproduced.

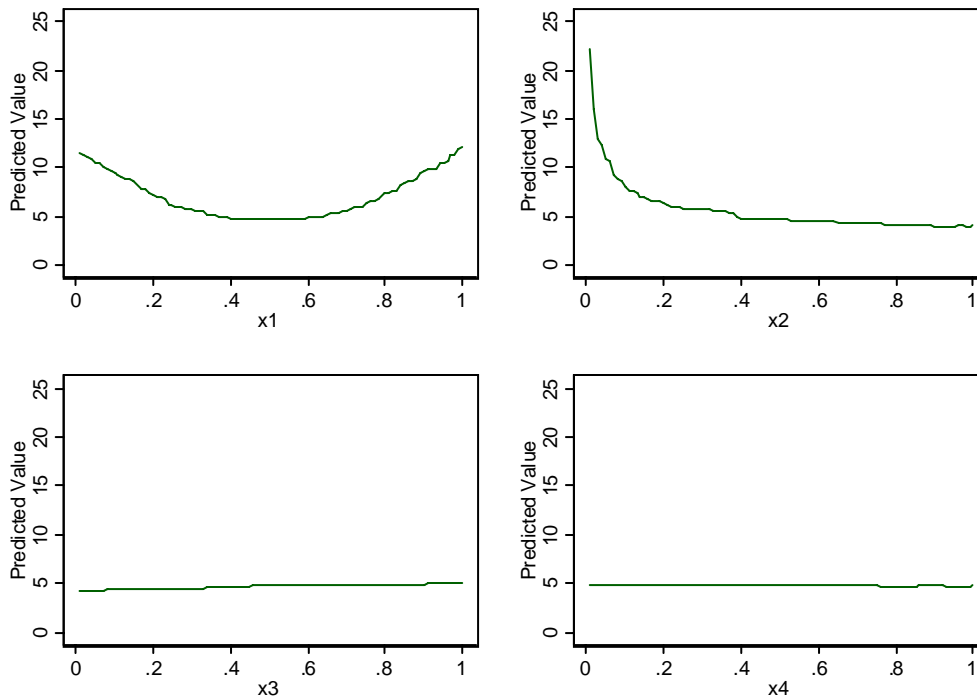


Figure 6: Conditional plots: Predictions for x_1 through x_4 while other variables are held constant at $x_i=0.5$.

6 Example: Boosted logistic regression

I compare the boosted logistic regression with a regular logistic regression model. I simulate data from the following model:

$$\log\left(\frac{p}{1-p}\right) = -35 + \sum_{j=1}^{10} \frac{1}{0.1 + \sum_{i=1}^3 (x_i - j/10)^2} - 100 I(x_4 > 0.95) + \varepsilon \quad (3)$$

where the indicator function $I(\text{arg})$ equals one if its argument is true and zero otherwise, $\varepsilon \sim \text{uniform}(0, 22.985)$ and $0 \leq x_i \leq 1$ for $i=1, 2, 3$. The value 22.985 corresponds roughly to one standard deviation of $\log(p/(1-p))$, i.e. the signal to noise ratio on the logit scale is 1. This model has a nonlinear 3-level interaction between x_1 , x_2 , and x_3 and a nonlinearity in form of a step function for x_4 . I simulate 46 additional variables x_5 through x_{50} , uniformly distributed across their support $0 \leq x_i \leq 1$ for $i=4, 5, \dots, 50$. All x variables are uncorrelated

to one another. The response y is a function of only the first four of the 50 variables. Of course, boosting still performs well when the variables are correlated. Boosting can also be used when there are more covariates than observations.

The data consist of the response y , 50 x -variables and 4000 observations. Half of these observations are used as training data and half as test data. I fit a regular, linear logistic regression model to the data. The odds ratios of the first 4 variables are shown in Table 2. The odds ratios were obtained by running the command “logistic y x_1 - x_{50} ”. Except for x_4 , none of the variables is significant. As expected by chance, a few of the remaining 46 variables were also significant at the 5% level.

Variable	Odds Ratio	p
x1	1.07	0.69
x2	1.26	0.15
x3	1.07	0.68
x4	0.54	<.001

Table 2: Odds ratios of the first 4 variables from the logistic regression

For the boosting model Figure 7 shows the R^2 value on a test data set as a function of the number of interactions. Clearly, the main effect model is not sufficiently complex. The slight dip in the curve for interaction=6 is just a reminder that these values are estimates and they are variable. The maximal R^2 is reached for interaction=8, but any number of interactions greater or equal to 4 would probably perform similarly. The large number of interactions is somewhat surprising because the model in equation (3) contains only a 3-level interaction. I speculate that the nonlinearities are more easily accommodated with additional nodes.

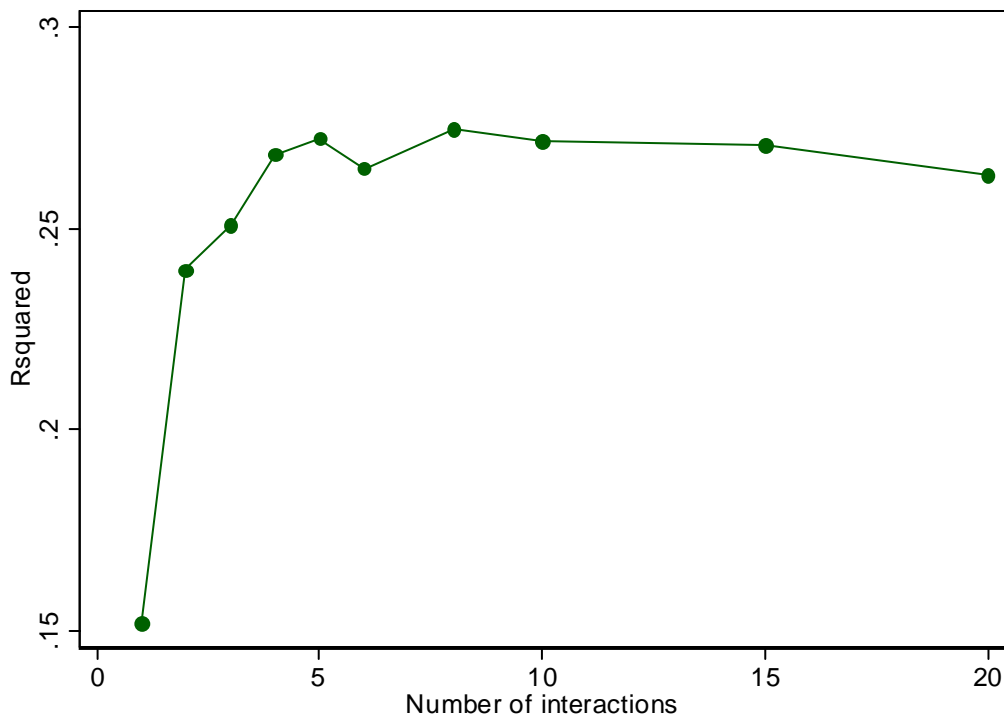


Figure 7: Scatter plot of the pseudo R^2 computed on a test data set versus the number of interactions.

I compare classification rates on the test data. Roughly 49% of the test data are classified as zero ($y=0$), and 51% as one ($y=1$). When using a coin flip to classify observations one would have been right about half the time. Using logistic regression 52.0% of observations in the test data set are classified correctly. This is just barely better than the rate one could have obtained by a coin flip. Because there were a lot of unrelated x-variables I also tried a backward regression using $p>0.15$ as criterion to remove a variable. Using the backward regression 54.1% of the observations were classified correctly. The boosting model with 8 interactions, shrinkage =0.5 and bag=0.5, classifies 76.0% of the test data observations correctly.

The Stata output displays the pseudo R^2 values for logistic regression (pseudo $R^2=0.02$) and backward logistic regression (pseudo $R^2=0.01$). Because the training data were used to compute the pseudo R^2 values, the backward logistic regression necessarily

has a lower value. Both values are much lower than the value obtained by boosted logistic regression (test $R^2=0.27$).

Because there are only two response values (0 and 1), I use a different plot for calibration than the scatter plot shown in Figure 5. If the predicted values are accurate one would expect that the predicted values are roughly the same as the fraction of response values classified as “1” that give rise to a given predicted value. The fraction of response values classified as “1” can be estimated by averaging or smoothing over response values with similar predictions. In Stata I use a lowess smoother to compare the predictions from the boosted logistic regression and the linear logistic regression:

```
twoway (lowess y logit_pred, bwidth(0.2)) (lowess y boost_pred,
bwidth(0.2)) (lfit straight y), xtitle("Actual Values")
legend(label(1 "Logistic Regression") label(2 "Boosting") label(3
"Fitted Values=Actual Values") ) xsize(4) ysize(4)
```

Calibration plots for the test data are shown in Figure 8. The near horizontal line for logistic regression in the test calibration plot implies that logistic regression classifies 50% of the observations correctly regardless of the actual predicted value. The logistic regression model does not generalize well.

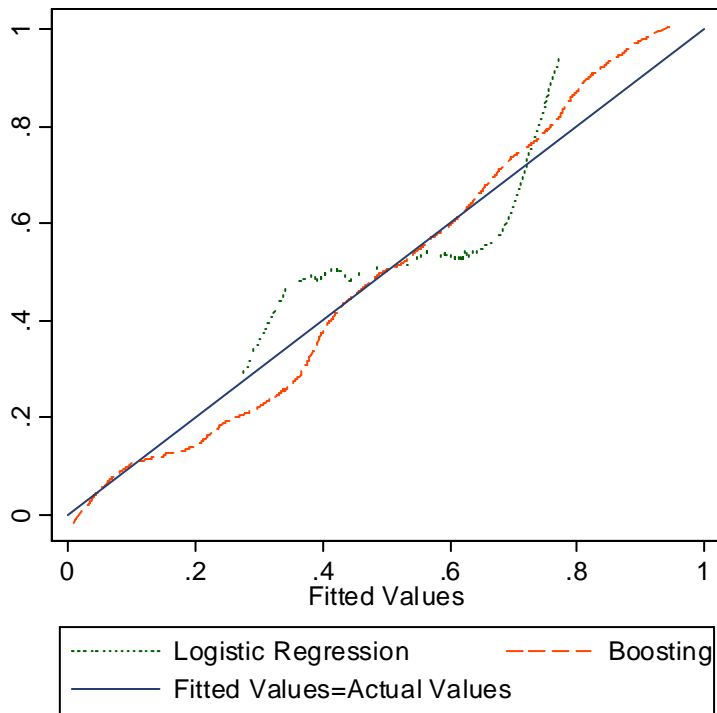


Figure 8: Calibration plot for the logistic regression example: Fitted versus actual values for logistic regression and boosting on the test data set.

A bar chart of the influence of each of the 50 variables is shown in Figure 9. The first bar corresponds to x_1 , the second to x_2 , and so forth. Boosting clearly discriminates between the important variables, the first 4 variables, and the remainder. The remaining 46 variables only explain a small percentage of the variation each. The concept of significance does not yet exist in boosting and there is no formal test that declares these variables “unimportant”. Interestingly, the influences of the noise variables, x_5 through x_{50} , depend on the number of interactions specified in the model. If the number of interactions is too large the influence of the noise variables increase.

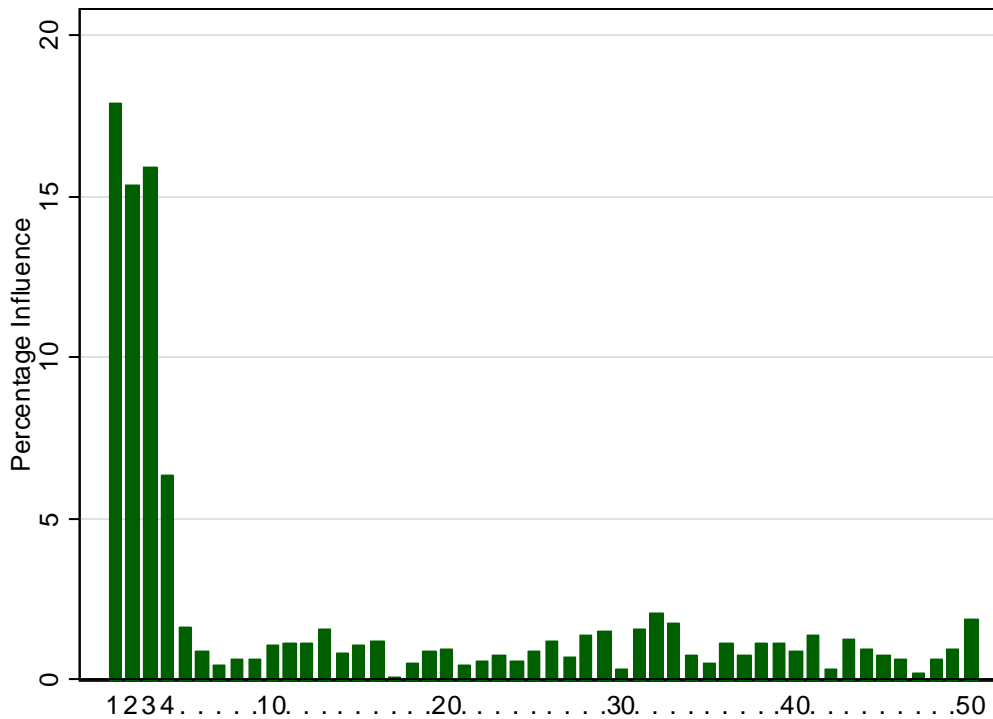


Figure 9: Percentage of influence of variables x_1 through x_{50} in the logistic regression example

To visualize the effect of x_4 on the response, I group the predictions (on the probability scale) from the training data set into 20 groups according to their corresponding x_4 value. The first group contains predictions where $0 < x_4 < .05$, the second where $0.05 < x_4 < 0.1$, and so forth. Figure 10 gives a boxplot for each of the 20 groups. Consistent with the model in equation (3) the last group with values $0.95 < x_4 < 1.0$ has much lower predictions than the other 19 groups. This approach to visualizing data is different from that in the previous example. In the previous example all other covariates were fixed at one value. All predictions in the training data set are used. This second method displays much more variation. Visualizing the nonlinear interaction between x_1 , x_2 and x_3 is of course much harder.

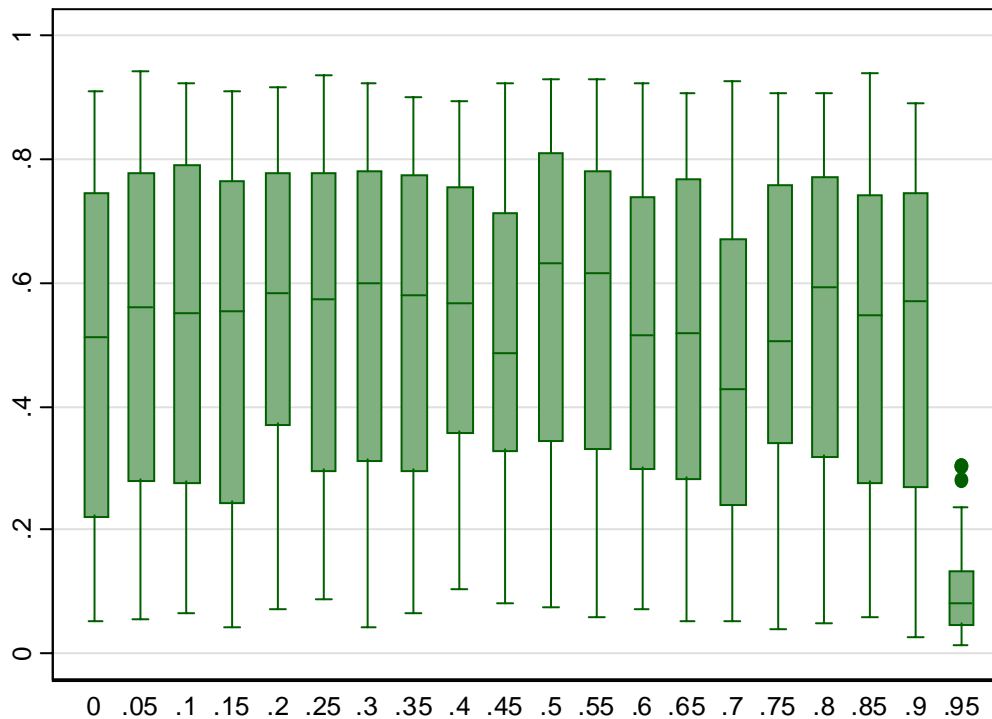


Figure 10: 20 Box plots of predictions for 20 non-overlapping subsets of x_4 . All predictions stem from the training data.

7 Runtime benchmarks

To generate runtime benchmarks I ran boosting models on randomly generated data. A single observation is generated by generating random x -values, $x_{ij} \sim \text{Uniform}(0,1)$, $j=1, \dots, p$ where p is the number of x -variables, $i=1, \dots, n$ denotes the observations and the output y_i is computed as follows:

$$y_i = \sum_{j=1}^p x_{ij}$$

I generated a number of data sets with varying numbers of observations $N \in \{100, 1000, 10000\}$, varying number of x -variables $p \in \{10, 30, 100\}$, varying numbers of iterations $\text{maxiter} \in \{1000, 5000, 10000\}$ and varying numbers of interactions $\text{interaction} \in \{2, 4, 6\}$. I fit a boosting model to each of the data sets specifying a normal

distribution, trainfraction=0.5, bagging=0.5, shrink=0.01 and predict(varname). Table 3 displays runtimes for all 81 combinations of these four factors.

Because the runtimes range over several orders of magnitude, Table 3 displays the run time in both seconds (top) and hours (bottom) for each combination. The benchmark calculations were computed on a Dell D600 laptop with a 1.6 GHz processor and 0.5 GB of RAM.

Obs	Variables	Iterations								
		1000			5000			10000		
		Interactions			Interactions			Interactions		
2	4	6	2	4	6	2	4	6		
100		0.3	0.5	0.6	1.3	2.2	3	2.6	4.4	6.2
	10	0	0	0	0	0	0	0	0	0
	30	0.6	1.1	1.6	3	5.5	7.8	6.2	11.1	15.6
	100	0	0	0	0	0	0	0	0	0
1000		2	3.5	5	9.5	17.9	25	19.1	35.3	50
	100	0	0	0	0	0	0	0	0	0
	10	3.9	5.6	7.9	15.8	22.1	31.1	27.6	43	62.9
	30	0	0	0	0	0	0	0	0	0
10000		7	11.4	17.2	32.3	57.6	82.6	63.8	113.9	158.7
	30	0	0	0	0	0	0	0	0	0
	100	39.7	75.9	111.1	197	380.9	557.6	394.5	758.4	1118
	100	0	0	0	0.1	0.1	0.2	0.1	0.2	0.3
10000		29.8	64.5	91.1	156.5	328	422.1	377	599.5	862.5
	10	0	0	0	0	0.1	0.1	0.1	0.2	0.2
	30	97	178.6	249.6	479.2	906.3	1311.7	964.3	1683.4	2780.7
	100	0	0	0.1	0.1	0.3	0.4	0.3	0.5	0.8
10000		487.8	935.1	1382.6	2443.7	4668.2	6840.8	4878.4	9360.5	13615.2
	100	0.1	0.3	0.4	0.7	1.3	1.9	1.4	2.6	3.8

Table 3: Benchmark runtimes for boosting: various combinations of the number of observations (50% used for training, 50% for testing), number of variables, boosting iterations, and number of interactions chosen. The time is given both in seconds (top number) and in hours (bottom number).

The runtime range from 0.3 seconds (100 observations, 10 variables, main effects only, 1000 iterations) to 3.8 hours (10,000 observations, 100 variables, six-way interactions, 10,000 iterations). The time increases roughly linearly with the number of iterations, the number of interactions, and the number of variables. The time increases more than linearly with the number of observations. Because the observations are sorted the runtime is $O(n \log(n))$ where n is the number of observations (i.e. the runtime is bounded by a constant times $n \log(n)$; for linear increases the runtime would be bounded by a constant times n).

Shrinking does not affect runtime, except in the sense that a smaller shrinkage value will tend to require a larger number of iterations. Bagging improves runtime. Typically, the runtime with bagging with 50% of the observations is roughly 30% faster than the runtime without bagging. Running logistic regression instead of normal regression increases the runtime only a little (usually less than 10%).

Table 3 displays runtimes for up to 10,000 observations, but 10,000 is not an upper limit. I have used this implementation of boosting with 100,000 observations.

8 Discussion

Boosting is a powerful regression tool. Unlike linear regression, boosting will work when there are more variables than observations. I successfully performed a boosted logistic regression with 200 observations and 500 x-variables. Linear logistic regression will cease to run normally with more than about 50 x-variables and assuming 200 observations because individual observations are uniquely identified.

The question “When should I use boosting?” is not easy to answer, but there some general indicators that are outlined in Table 4.. A strength of the boosting algorithm is that interactions and nonlinearities need not be explicitly specified. Unless the functional relationship is highly nonlinear, there is probably little point in using boosted regression in a small data set with, for example, 50 observations and a handful of variables. In small data sets linearity is usually an adequate approximation. Large numbers of continuous variables make nonlinearities more likely. Indicator variables have only two levels. Nonlinearities cannot arise from indicator variables. Ordered categorical x- variables are awkward to deal with in regular regression. Because boosting uses trees as a base learner it is highly suited for the use of ordered variables. The separation of training and test data guards against overfitting that may arise in the context of correlated data.

Indicator	Indicator favors the use of boosting	Indicator against the use of boosting	Why?
-----------	--------------------------------------	---------------------------------------	------

small data set		x	linear approximation usually adequate
large data set	x		nonlinearities and interactions likely
more variables than observations (or close)	x		linear (Gaussian and logistic) regression methods fail
suspected nonlinearities	x		nonlinearities need not be explicitly modeled
suspected interactions	x		interactions need not be explicitly specified
ordered categorical x-variables	x		awkward in parametric regression
correlated data	x		potential for overfitting
x-variables consist of indicator variables only		x	nonlinearities cannot arise from indicator variables, interactions still might

Table 4: Some indicators in favor and against the use of boosting.

I would like to point out a few issues that the reader may run into as he or she starts experimenting with the boosting plugin. The “tree fully fit” error means that more tree splits are required than are possible. This arises, for example, if the data contain only 10 observations but *interaction*=11 (or if *bag*=0.5 and *interaction*=6) is specified. It will also tend to arise when the number of iterations (*maxiter*) and the number of interactions (*inter*) are accidentally switched. Reducing the number of interactions always solves this problem. A second issue are missing values, which are not supported in the current version of the boosting plugin. I suggest imputing variables ahead of time, for example using a stratified hotdeck imputation (my implementation of hotdeck, “hotdeckvar”, is available from www.schonlau.net or type “net search hotdeckvar” within Stata). Another option that is popular in the social sciences is to create variables that flag missing data and add them to the list of covariates. Discarding observations with missing values is a less desirable alternative. I am working on a future release that allows saving the boosting model in Stata.

Acknowledgement

I am most grateful to Gregory Ridgeway, developer of the GBM boosting package in R, for many discussions on boosting, advice on C++ programming and for comments on earlier versions of this paper. I am equally grateful to Nelson Lim at the

RAND Corporation for his support and interest in this methodology and for comments on earlier versions of this paper.

References

- Breiman, L., Friedman, J., Olshen, R., and Stone C. 1984. *Classification and Regression Trees*. Belmont, CA: Wadsworth.
- Elkan, C. Boosting and Naïve Bayes Learning. Technical Report No CS97-557. September 1997, UCSD.
- Friedman, J., Hastie, T. and R. Tibshirani. 2000. Additive logistic regression: a statistical view of boosting. *The Annals of Statistics* 28(2):337-407.
- Friedman, J.H. (2001). "Greedy Function Approximation: A Gradient Boosting Machine," *Annals of Statistics* 29(5):1189-1232.
- Hastie, T., Tibshirani, R. and J. Friedman. 2001. *The Elements of Statistical Learning*. New York: Springer-Verlag.
- Long, J. S. and J. Freese. 2003. *Regression Models for Categorical Dependent Variables Using Stata*, Revised Edition. College Station, TX: Stata Press.
- McCaffrey, D., Ridgeway, G., Morral, A. (2004, to appear). "Propensity Score Estimation with Boosted Regression for Evaluating Adolescent Substance Abuse Treatment," *Psychological Methods*.
- McCullagh, P. and Nelder, J. A. (1989). *Generalized Linear Models*. New York: Chapman & Hall.
- Ridgeway, G. (1999). "The state of boosting," *Computing Science and Statistics* 31:172-181. Also available at <http://www.i-pensieri.com/gregr/papers.shtml> .