

BOOSTING 3: IMPLEMENTATIONS

-STATISTICAL MACHINE LEARNING-

Lecturer: Darren Homrighausen, PhD

OUTLINE

Now we will discuss two current, popular algorithms and their **R** implementations

- **GBM**
- **XGBoost**

GBM

GRADIENT BOOSTING MACHINES (GBM)

RECALL: AdaBoost effectively uses forward stagewise minimization of the exponential loss function

GBM takes this idea and

- generalizes to other loss functions
- adds subsampling
- includes methods for choosing B
- reports variable importance measures

GBM: LOSS FUNCTIONS

- **gaussian**: squared error
- **laplace**: absolute value
- **bernoulli**: logistic
- **adaboost**: exponential
- **multinomial**: more than one class (unordered)
- **poisson**: Count data
- **coxph**: For right censored, survival data

GBM: SUBSAMPLING

Early implementations of AdaBoost randomly sampled the weights (w)

This wasn't essential and has been altered to use deterministic weights

Friedman (2002) introduced stochastic gradient boosting that uses a new subsample at each boosting iteration to find and project the gradient

This has two possible benefits

- Reduces computations/storage
(But increases read/write time)
- Can improve performance

GBM: SUBSAMPLING

You can expect performance gains when **both** of the following occur:

- There is a small sample size
- The **base learner** is relatively complex

This suggests the usual ‘variance reduction through lowering covariance’ interpretation

The effect is complicated, though as subsampling

- increases the **variance** of each term in the sum
- decreases the **covariance** of each term in the sum

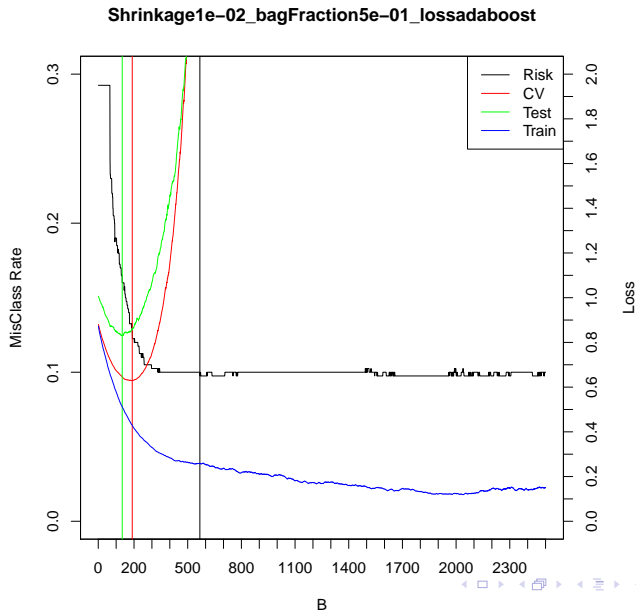
GBM: CHOOSING B

There are three built in methods:

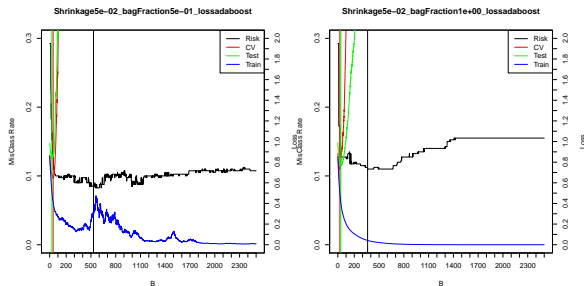
- **INDEPENDENT TEST SET:** using the `nTrain` parameter to say 'use only this amount of data for training'
(Be sure to uniformly permute your data set first.)
- **OUT-OF-BAG (OOB) ESTIMATION:** If `bag.fraction` is > 0 , then `gbm` can use OOB at each iteration to find a good B
(Note: OOB tends to select a too-small B)
- **K -FOLD CROSS VALIDATION (CV):** It will fit `cv.folds+1` models
(The '+1' is the fit on all the data that is reported)

Simulation: Continued

GBM: SOME OBSERVATIONS

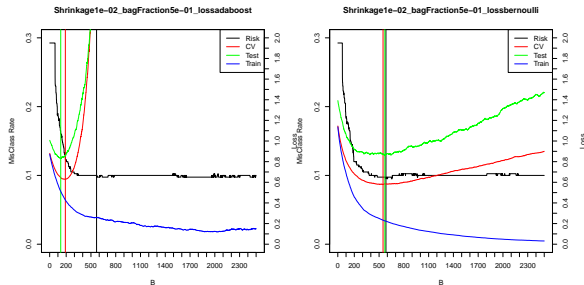


GBM: BAG FRACTION



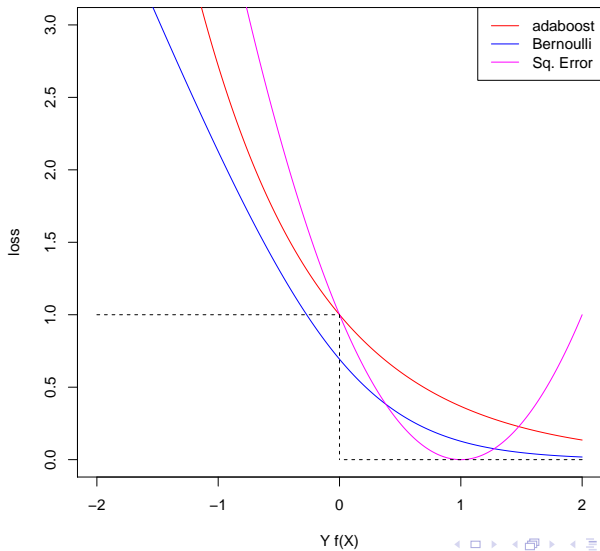
Bag fraction of 0.5 improves the best risk

GBM: LOSS FUNCTION



Bernoulli loss: CV-minimum $B \approx$ Risk-minimum B

COMPARING LOSS FUNCTIONS



GBM: VARIABLE IMPORTANCE MEASURE

For tree-based methods, there are **variable importance measures**:

RELATIVE.INFLUENCE: For each feature X_j and tree T_b

$$\text{Influence}_j(T_b) = \sum_{\text{Split on } X_j} (\text{Reduction in loss})^2$$

(See equations (43) - (45) in Friedman (2001) for details)

This is aggregated to form

$$\text{Influence}_j = \frac{1}{B} \sum_{b=1}^B \text{Influence}_j(T_b)$$

(There is also **permutation.test.gbm**, but it is currently labeled experimental)

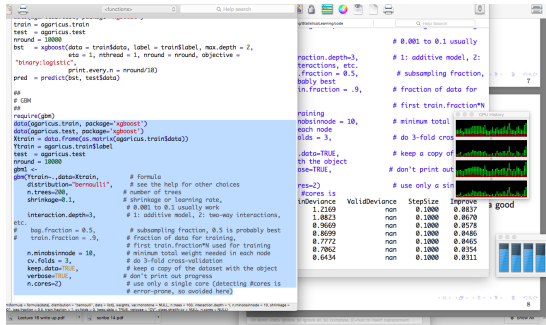
GBM: SAMPLE CODE

```
gbm(Ytrain~.,data=Xtrain,  
    distribution="bernoulli",  
    n.trees=500,  
    shrinkage=0.01,  
    interaction.depth=3,  
    bag.fraction = 0.5,  
    n.minobsinnode = 10,  
    cv.folds = 3,  
    keep.data=TRUE,  
    verbose=TRUE,  
    n.cores=2)
```

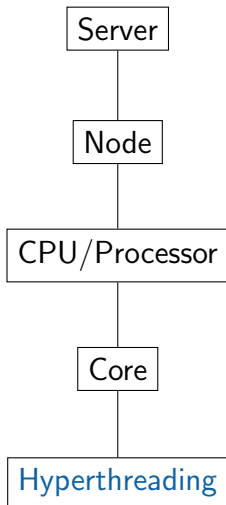
Keep adding trees with **gbm.more**

(If this is taking too long, increase the learning rate, **shrinkage**)

GBM: FIGURES



DISTRIBUTED COMPUTING HIERARCHY



EXAMPLE: A server might have

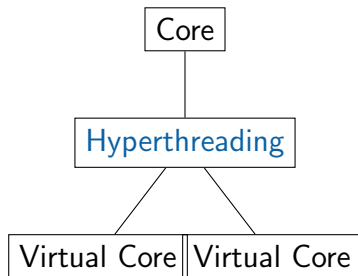
- 64 nodes
- 2 processors per node
- 16 cores per processor
- **hyper threading**

The goal is to somehow allocate a **job** so that these resources are used efficiently

Jobs are composed of **threads**, which are specific computations

HYPERTHREADING

Developed by Intel, Hyperthreading allows for each core to pretend to be two cores



This works by trading off computation and read-time for each core

XGBoost

XGBOOST

This stands for:

EXTREME GRADIENT BOOSTING

It has some advances over **gbm**

XGBOOST: ADVANCES

- **SPARSE MATRICES:** Can use sparse matrices as inputs
(In fact, it has its own matrix-like data structure that is recommended)
- **OPENMP:** Incorporates OpenMP on Windows/Linux
(OpenMP is a **message passing** parallelization paradigm for shared memory parallel programming)
- **LOSS FUNCTIONS:** You can specify your own loss/evaluation functions
(You need to use **xgb.train** for this)