

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 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` 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)

GBM: VARIABLE IMPORTANCE MEASURE

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

- **relative.influence**
- **permutation.test.gbm**

(This is currently labeled experimental)

These have similar definition relative to **bagging**, however they use **all** of the data instead of the OOB

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)
```

GBM: FIGURES

```

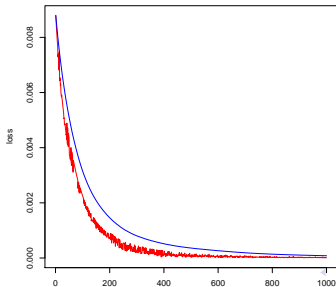
train = agricus.train
test = agricus.test
nround = 10000

bst = xgboost(data = traindata, label = trainlabel, max_depth = 2,
             eta = 1, minleaf = 1, nround = nround, objective =
             "binary:logistic",
             print_every_n = nround/10)
pred = predict(bst, testdata)

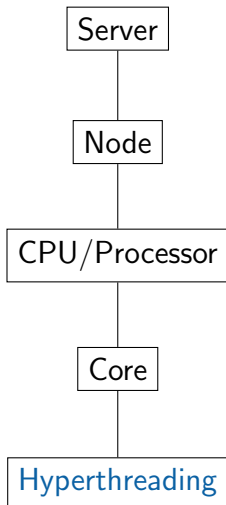
##
## GBM
##
require(gbm)
data(agriscus.train, package="xgboost")
data(agriscus.test, package="xgboost")
Xtrain = data.frame(as.matrix(agriscus.traindata))
Ytrain = agriscus.trainlabel
test = agriscus.test
nround = 10000
gbml <-
gbm(Xtrain~., data=Xtrain,          # formula
     distribution="bernoulli",     # see the help for other choices
     n.trees=200,                 # number of trees
     shrinkage=0.1,              # shrinkage or learning rate,
                               # 0.001 to 0.1 usually work
     interaction.depth=3,        # 1: additive model, 2: two-way interactions,
                               # 3: three-way interactions, etc.
     bag.fraction = 0.5,         # subsampling fraction, 0.5 is probably best
     train.fraction = .3,       # fraction of data for training,
                               # first train.fraction% used for each node
     n.minobsinnode = 10,        # minimum total weights needed in each node
     cv.folds = 3,               # do 3-fold cross-validation
     keep.data=TRUE,            # keep a copy of the dataset with the object
     verbose=TRUE,              # don't print out progress
     n.cores=2)                 # use only a single core (detecting #cores is
                               # error-prone, so avoided here)

```

#cores is	ValidDeviance	StepSize	Improve
1.2169	nan	0.1000	0.0837
1.0623	nan	0.1000	0.0670
0.2669	nan	0.1000	0.0578
0.8699	nan	0.1000	0.0486
0.7772	nan	0.1000	0.0465
0.7062	nan	0.1000	0.0354
0.6434	nan	0.1000	0.0311



DISTRIBUTED COMPUTING HIERARCHY



EXAMPLE: A server might have

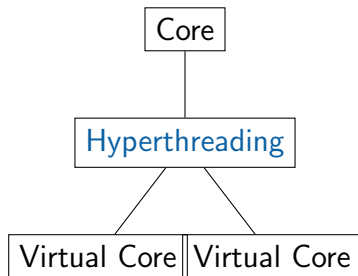
- 64 nodes
- 2 processors per node
- 16 cores per processor
- **hyper threading**

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

BOOSTING: LEARNING SLOW

It is best to set the **learning rate** at a small number.

This is usually calibrated by the computational demands of the problem.

A good strategy is to pick a number, say .001

Run with **n.trees** relatively small and see how long it takes

Keep adding trees with **gbm.more**. If this is taking too long, increase the learning rate

XGBoost

XGBOOST

This stands for:

EXTREME GRADIENT BOOSTING

It has some advances related to **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)