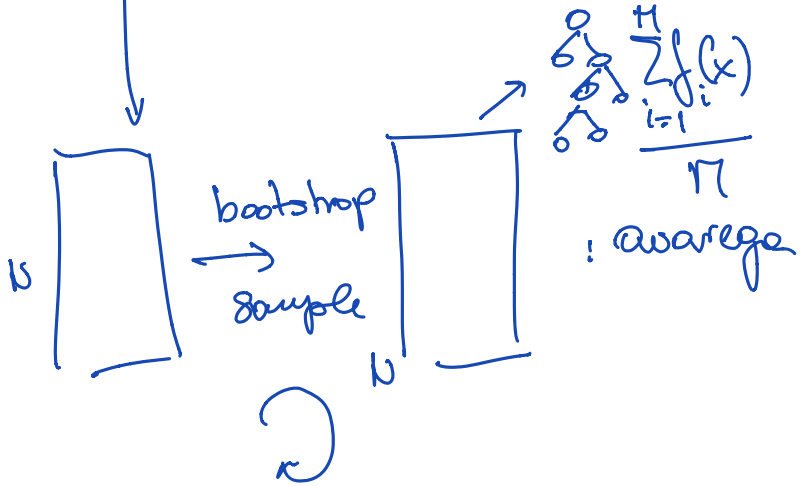


bagging
bootstrap aggregating

Ensembles
reduce the list
on the test set
improve accuracy
interpretability

Boosting

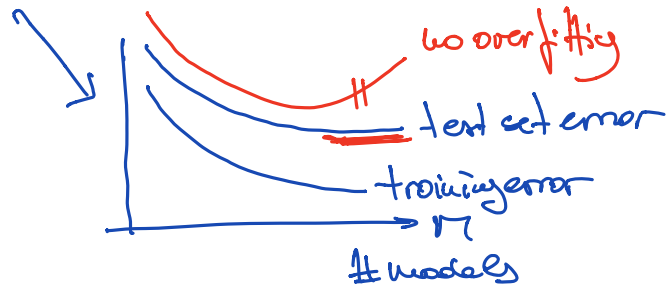
Machine Learning for Data Science 1



no covariance

$$\sigma_{\bar{x}}^2 = \frac{\sum \sigma_i^2 + 2 \sum \sum \text{cov}(x_i, y_j)}{n^2}$$

↓ → 0



Ada Boost.M1 : Shapiro 1997 : classification, binary
 $G_m(x) \in \{-1, 1\}$

1. Initialize training data weights
 $w_i = 1/N \quad i = 1 \dots N$

2. for $m = 1$ to M

fit $G_m(x)$ to the weighted training data

compute $err_m = \frac{\sum w_i I(y_i \neq G_m(x_i))}{\sum w_i}$

compute $\alpha_m = \log \frac{1 - err_m}{err_m}$
→ accuracy → error

set $w_i \leftarrow w_i \cdot e^{\alpha_m I(y_i \neq G_m(x_i))}$

not independent for correct predictions, weights are not changed

3. output

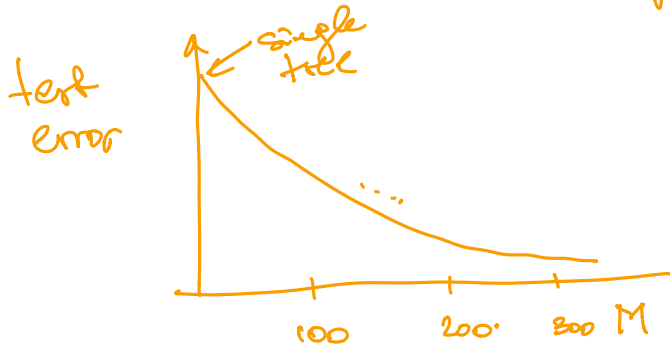
$\rightarrow G(x) = \text{sign} \left[\sum_{i=1}^M \alpha_m G_m(x) \right]$
similar to logistic for incorrect predictions
 $\frac{1 - err}{err}$

$$\alpha_1 G_1(x) + \alpha_2 G_2(x) + \alpha_3 G_3(x) + \alpha_4 G_4(x) = G(x)$$

\uparrow \uparrow \uparrow \uparrow
 noisy weighted weighted ...
 data sample sample ...

Ada Boost. M1: increases the weights of misclassified data instances

dramatic increase of performance → accuracy



- Why?
 - what are we optimizing?
 - are there any other algorithms of this kind?
- ↓

powerful

Boosting & Additive models

G_1, G_2, G_3, \dots

$$G(x) = \text{sign} \left(\sum \alpha_m G_m(x) \right)$$

boosting, we expand the approximation one model at a time

basis function expansion

$$f(x) = \sum_{m=1}^M \beta_m b(x, \gamma_m)$$

\uparrow a set of parameters

* minimize $\sum_{i=1}^N \mathbb{1} \left(y_i, \sum_{m=1}^M \beta_m b(x_i, \gamma_m) \right)$ $M=1000$

α weight

computationally intensive optimization

optimizing all parameters at once
(PROBLETTIC)

ALL AT ONCE

Forward Stagewise Additive Modelling

Approximates the solution to $\min_{\beta, \gamma} L(\sum \dots)$
by adding a new function
to the expansion

w/o adjusting the coefficients
already inferred

1. Initialize $f_0(x) = 0$

2. for $m=1$ to M

compute $\beta_m, \gamma_m = \arg \min_{\beta, \gamma} \sum_{i=1}^n L(y_i, f_{m-1}(x_i) + \beta b(x_i, \gamma))$

set $f_m(x) = f_{m-1}(x) + \beta_m b(x, \gamma_m)$

↳ correct the output
of the previous model
(expansion)

3. output $f_M(x)$

AdaBoost.M1 and Loss Function

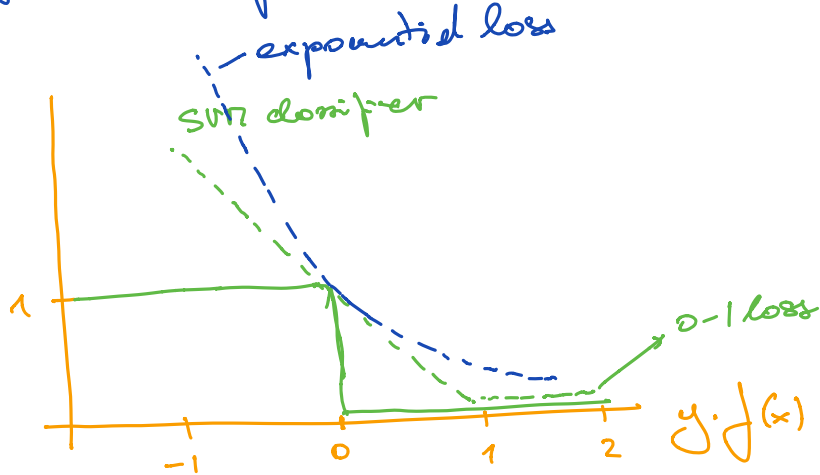
What is the loss function optimized by AdaBoost?

Can similar procedures be used for other loss functions?

AdaBoost \equiv forward stepwise additive model fitting that uses exponential loss

+5 years

$$L(y, f(x)) = \exp(-y f(x)) \leftarrow$$



AdaBoost. II: $G_m(x) = \{-1, 1\}$

For exponential loss, we have to solve:

$$\beta_m, G_m = \arg \min_{\beta, G} \sum_{i=1}^N \exp[-y_i (f_{m-1}(x_i) + \beta G(x_i))]$$

$$= \arg \min_{\beta, G} \sum_{i=1}^N \omega_i^{(m)} \exp(-\beta y_i G(x_i))$$

$e^{-\beta}$ $e^{+\beta}$ $= 1$

$$\omega_i^{(m)} = \exp(-y_i f_{m-1}(x_i))$$

or weight applied to each observation changes with each iteration

$$= \arg \min e^{-\beta} \sum_{y_i = G(x_i)} \omega_i^{(m)} + e^{+\beta} \sum_{y_i \neq G(x_i)} \omega_i^{(m)}$$

$$= \arg \min [(e^{\beta} + e^{-\beta}) \sum_{i=1}^N \omega_i^{(m)} I(y_i \neq G_m(x_i)) + e^{-\beta} \sum_{i=1}^N \omega_i^{(m)}]$$

$$G_m = ? \quad G_m = \arg \min_{G_m} \sum \omega_i^{(m)} I(y_i \neq G_m(x_i))$$

classifier that minimizes weighted error rate

E

What is β ?

$$\frac{\partial E}{\partial \beta} = \frac{d \left(\sum_{y=C(x_i)} w_i^{(u)} e^{-\beta} + \sum_{y \neq C(x_i)} w_i^{(u)} e^{\beta} \right)}{d \beta} = 0$$

$$\beta_m = \frac{1}{2} \log \frac{1 - \epsilon}{\epsilon}$$

$$\epsilon = \frac{\sum w_i^{(u)} I(y_i \neq G_m(x_i))}{\sum w_i^{(u)}}$$

looks familiar

$$2\beta_m = \alpha_m$$

The update of approximation

$$f_m(x) = f_{m-1}(x) + \beta_m G_m(x)$$

$$w_i^{(u+1)} = w_i^{(u)} \cdot e^{-\beta_m y_i G_m(x_i)}$$

$$-y_i G_m(x_i) = 2I(y_i \neq G_m(x_i)) - 1$$

$$w_i^{(u+1)} = w_i^{(u)} \cdot e^{2\beta_m I(y_i \neq G_m(x_i))}$$

exactly the
weight
update for
AdaBoost
!!!

Ada Boost Recap

At each iteration

choose classifier that minimizes weighted error

We use this to estimate error rates

We use this to compute the weights

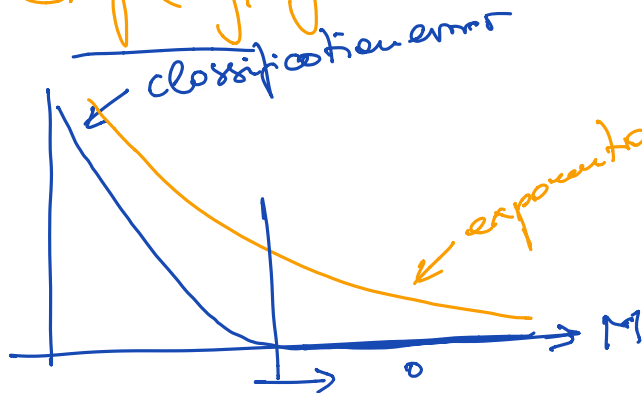
improve the overall classifier with new model

$$f_M(x) = f_{M-1}(x) + \beta_M G_M(x)$$

→ forward stepwise additive model
that minimizes

$$L = \sum_i \exp(y_i - f(x_i))$$

training error



under change
of loss function

↳ what is
the training
error

↳ what
are the
predictions

Boosting Trees ↓ formalism

trees: partition feature space to disjoint regions R_j
where a constant is assigned

$$x \in R_j \Rightarrow f(x) = y_j$$

formally

$$T(x; \Theta) = \sum_{j=1}^J y_j I(x \in R_j)$$

$$\Theta = \arg \min \sum_{j=1}^J \sum_{i=1}^N L(y_i, y_j)$$

all regions all examples

optimization problem

heuristics

fixed y_j given R_j : $\hat{y}_j = \bar{y}_j$ the mean

find R_j : greedy, top-down recursive partitioning

Boosted Tree Model

$$f_{\pi}(x) = \sum_{m=1}^M T(x; \Theta_m)$$

nodes labels
lossing

induced in
forward stepwise manner
at each stage

$$\Theta_m = \arg \min_{\Theta_m} \sum_{i=1}^N L(y_i, f_{m-1}(x_i) + T(x_i; \Theta_m))$$

given the regions

$$f_{jm} = \arg \min_{f_{jm}} \sum_{i=1}^N L(y_i, f_{m-1}(x_i) + f_{jm})$$

finding regions is difficult

but in some cases,
there is heuristic approach

$$|L(y, f(x)) = (y - f(x))^2$$

$$L(y_i, f_{m-1}(x_i) + T_m(x_i)) =$$

$$= \left(\underbrace{(y_i - f_{m-1})}_{\text{residual of the current model}} - \underbrace{T_m(x_i)}_{\epsilon} \right)^2$$

residual of the current model

ϵ

$$f_0(x) = \bar{y}$$

↑
fitting data

→ residual of f_0

↓
fit T_1

$$f_1 = f_0 + T_1$$

↑

→ residual of f_1

↓
fit T_2

$$\frac{f_0 + f_1 + T_2}{4} \dots$$

to minimize this

the tree should be fitted to the residual

Numerical Optimization

$$L(f) = \sum_{i=1}^n L(y_i, f(x_i))$$

Goal: $\min_f L(f)$

$f(x)$: is a sum of trees

$$\hat{f} = \arg \min_f L(f)$$

Numerical optimization

$$\vec{f}_m = \sum_{u=0}^m \vec{f}_u \quad f_u \in \mathbb{R}^D$$

f_0 : initial guess

f_m : induced based on f_{m-1} , which is the sum of previously induced vectors

Steepest descent

$$\vec{f}_m = g_m \cdot \vec{g}_m : \text{gradient}$$

$$g_{im} = \left[\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right] \quad f(x_i) = \underline{f_{m-1}(x)}$$

Loss

regression $\frac{1}{2} [y_i - f(x_i)]^2$

regression $|y_i - f(x_i)|$

classif.

distance

$$\sum_{k=1}^K p_{ki} \log p_{ki}$$

⋮

$$-\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)}$$

$$y_i - f(x_i) = \underline{\text{residual}}_i$$

$$\text{sign}[y_i - f(x_i)]$$

k-th component

$$I(y_i = G_k) - p_k(x_i)$$

We fit the regression trees

Gradient Tree Boosting Algorithm (L)

1. Initialize $f_0(x) = \underset{f}{\operatorname{argmin}} \sum_{i=1}^n L(y_i, f)$

2. for $m=1$ to M

compute $r_m = - \left[\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f=f_{m-1}}$

fit regression tree to r_m

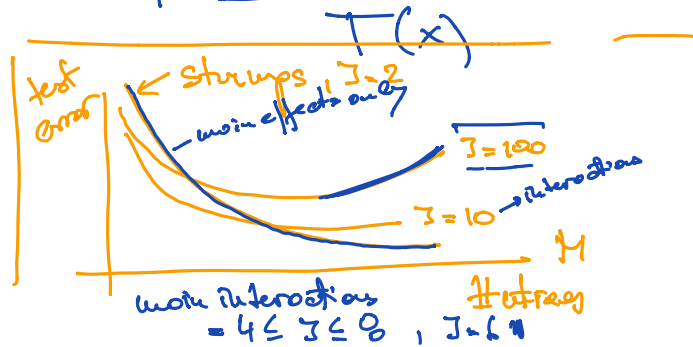
for $j=1 \dots J$ (regions) compute

$f_{m,j} = \underset{f}{\operatorname{argmin}} \sum_{x \in R_{m,j}} L(y_i, f_{m-1}(x_i) + f)$

update $f_m(x) = f_{m-1}(x) + \sum_{j=1}^J \beta_{m,j} I(x \in R_j)$

3. output $\hat{f}(x) = f_M(x)$

choice of the size of the trees



AutoML

H2O.ai Experiment tilucege

DATASETS EXPERIMENTS MLI HELP PY_CLIENT MOJO2-RUNTIME MESSAGES[2] LOGOUT ASD

TRAINING DATA

DRIVERLESS AI 1.3.1 - AI TO DO AI
Licensed to Oracle (SN26943 - For evaluation only, not for production use)

ASSISTANT

DATASET
bcwd.csv

ROWS: **699** COLUMNS: **10** DROPPED COLS: **--** VALIDATION DATASET: **--** TEST DATASET: **--**

TARGET COLUMN: **label** FOLD COLUMN: **--**

WEIGHT COLUMN: **--** TIME COLUMN: **[OFF]**

TYPE: **str** COUNT: **699** UNIQUE: **2** TARGET FREQ: **241**

TUNED 54/432 PARAMETER & FEATURE TUNING MODELS. TUNING (XGBOOST)



EXPERIMENT SETTINGS

8 ACCURACY CLASSIFICATION

2 TIME REPRODUCIBLE

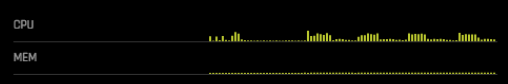
8 INTERPRETABILITY **ENABLE GPU**

EXPERT SETTINGS

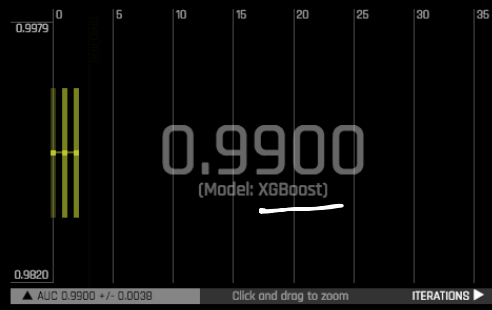
SCORER

BINI
MCC
F06
F1
F2
ACCURACY
LOGLOSS
AUC
AUCPR

CPU / MEMORY



ITERATION DATA - VALIDATION



VARIABLE IMPORTANCE

0_uniformity_cell_size	1.00
7_uniformity_cell_shape	0.69
0_bare_nuclei	0.45
_bland_chromatin	0.27
2_clump_thickness	0.21
5_normal_nucleoli	0.13
5_single_epithelial_cell_size	0.11
3_marginal_adhesion	0.05
4_mitoses	0.01

ROC PREC-RECALL LIFT GAINS GPU USAGE



H2O.ai Experiment 91e471

1.0.4

Show Experiments

TRAINING DATA

DATASET

BNPParibas-train.csv

ROWS: 114K COLUMNS: 133 DROPPED COLS: 0 TEST DATASET: Yes

TARGET COLUMN

target

TYPE: Int COUNT: 114321 UNIQUE: 2 FREQ: 27300

SCORED 393/426 MODELS ON 4318 FEATURES



EXPERIMENT SETTINGS



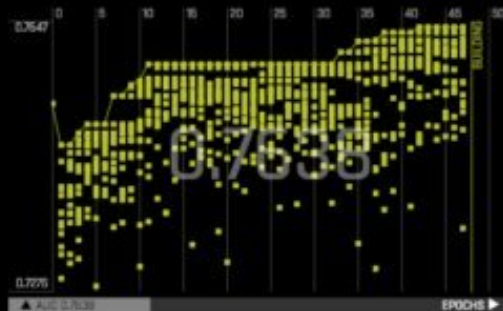
SCORER

IQ
MSE
RMSE
RMSE_L
MAE
SIM
ALL
LOGLOSS

CPU / MEMORY



ITERATION SCORES - INTERNAL VALIDATION



VARIABLE IMPORTANCE

89_v80	1.00
189_CV_TE_v129_v04_v86_0	0.74
187_WoL_v121_v02_v04_v03_v88_v62_v86_0	0.62
14_CV_TE_v86_0	0.30
8_WoL_v02_0	0.22
24_v80	0.19
72_NumToCotTE_v04_0	0.15
182_NumCotTE_v90_v81_v48_v90_v04_v88_0	0.15
181_WoL_v02_v04_v86_v86_0	0.11
149_NumToCotWoL_v1_v81_v84_0	0.11
7_CV_TE_v84_0	0.10
18_CV_TE_v84_0	0.09
146_CV_CatNumInc_v88_v80_median	0.09
185_NumCotTE_v106_v10_v04_v8_v85_v62_v86_v17_0	0.08

GPU USAGE

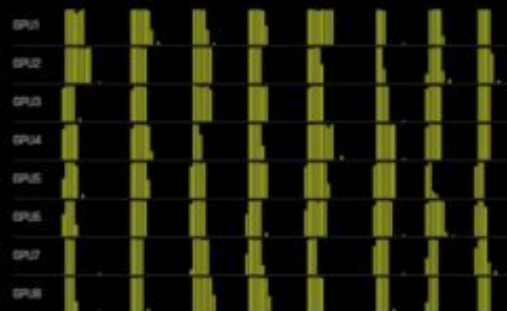


TABLE OF CONTENTS

- Installation Guide
- Get Started with XGBoost
- XGBoost Tutorials
- Frequently Asked Questions
- XGBoost User Forum
- GPU support
- XGBoost Parameters
- Python package
- R package
- JVM package
- Ruby package
- Julia package
- C Package
- C++ Interface
- CLI interface
- Contribute to XGBoost

Get Started with XGBoost

This is a quick start tutorial showing snippets for you to quickly try out XGBoost on the demo dataset on a binary classification task.

Links to Other Helpful Resources

AGGLE

- See [Installation Guide](#) on how to install XGBoost.
- See [Text Input Format](#) on using text format for specifying training/testing data.
- See [Tutorials](#) for tips and tutorials.
- See [Learning to use XGBoost by Examples](#) for more code examples.

Python

```
import xgboost as xgb
# read in data
dtrain = xgb.DMatrix('demo/data/agaricus.txt.train')
dtest = xgb.DMatrix('demo/data/agaricus.txt.test')
# specify parameters via map
param = {'max_depth':2, 'eta':1, 'objective':'binary:logistic' }
num_round = 2
bst = xgb.train(param, dtrain, num_round)
# make prediction
preds = bst.predict(dtest)
```