→ Stagewise Additive method

→ Weak learners (High Bias low variance)

→ Decision stumps



→ ~~1~~/0  1/-1

→ In Random Forest each tree has an equal vote on
   the final prediction but in a Forest of stumps made
   with AdaBoost, some stumps get more say in the
   final prediction than others.

→ In Random Forest each DT is made independently of
   each other. But, in a Forest of stumps made with AdaBoost,
   order is important. The errors that the first stump makes influence
   how the second stump is made and so on....

→ Trees are fit sequentially to improve error of previous trees

---

**Algorithm 10.1** *AdaBoost.M1.*

1. Initialize the observation weights $w_i = 1/N$, $i = 1, 2, \ldots, N$.

2. For $m = 1$ to $M$:

   (a) Fit a classifier $G_m(x)$ to the training data using weights $w_i$.

   (b) Compute
   $$\text{err}_m = \frac{\sum_{i=1}^{N} w_i I(y_i \neq G_m(x_i))}{\sum_{i=1}^{N} w_i}.$$

   (c) Compute $\alpha_m = \log((1 - \text{err}_m)/\text{err}_m)$.

   (d) Set $w_i \leftarrow w_i \cdot \exp[\alpha_m \cdot I(y_i \neq G_m(x_i))]$, $i = 1, 2, \ldots, N$.

3. Output $G(x) = \text{sign}\left[\sum_{m=1}^{M} \alpha_m G_m(x)\right]$.

---

Algorithm

Data

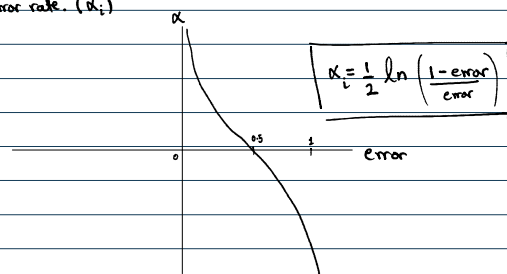| $x_1$ | $x_2$ | $y$ |
|---|---|---|
| 1 | 5 | 1 |
| 2 | 4 | 0 |
| 9 | 2 | 1 |
| 1 | 7 | 0 |
| 2 | 3 | 1 |

1) assign initial weight to each row.
   weight = $\frac{1}{n_c}$ ← total # of rows

2) train a Decision stump
   - check performance on data [calc predictions]
   - calculate model's weight, which depends on the
     error rate. ($\alpha_i$)

| $x_1$ | $x_2$ | $y$ | wt |
|---|---|---|---|
| 1 | 5 | 1 | 1/5 |
| 2 | 4 | 0 | 1/5 |
| 9 | 2 | 1 | 1/5 |
| 1 | 7 | 0 | 1/5 |
| 2 | 3 | 1 | 1/5 |



$$\alpha = \frac{1}{2} \ln\left(\frac{1 - \text{error}}{\text{error}}\right)$$

if error is 1 or 0 equation will panic
add a very small number.

| $x_1$ | $x_2$ | $y$ | wt | $\hat{y}$ |
|---|---|---|---|---|
| 1 | 5 | 1 | 1/5 | 0 |
| 2 | 4 | 0 | 1/5 | 1 |
| 9 | 2 | 1 | 1/5 | 1 |
| 1 | 7 | 0 | 1/5 | 0 |
| 2 | 3 | 1 | 1/5 | 1 |

→ error = sum of weights of misclassified points

error = $\frac{2}{5}$    $\alpha_1 = \frac{1}{2} \ln\left(\frac{1 - \frac{2}{5}}{\frac{2}{5}}\right) \approx 0.2$

3) Update weights [increasing the importance of misclassified points
                    decreasing the "  " correctly classified "]

| $x_1$ | $x_2$ | $y$ | wt | $\hat{y}$ | new_wt | normalized |
|---|---|---|---|---|---|---|
| 1 | 5 | 1 | 1/5 | 0 | 0.24 | 0.25 |
| 2 | 4 | 0 | 1/5 | 1 | 0.24 | 0.25 |
| 9 | 2 | 1 | 1/5 | 1 | 0.16 | 0.166 |
| 1 | 7 | 0 | 1/5 | 0 | 0.16 | 0.166 |
| 2 | 3 | 1 | 1/5 | 1 | 0.16 | 0.166 |
| | | | | | 0.96 | |

for misclassified

new_wt = curr_wt × $e^{\alpha_i}$

for correctly classified

new_wt = curr_wt × $e^{-\alpha_i}$

| $x_1$ | $x_2$ | $y$ | new.wt | $\hat{y}$ | Range |
|---|---|---|---|---|---|
| 1 | 5 | 1 | 0.25 | 0 | 0 - 0.25 |
| 2 | 4 | 0 | 0.25 | 1 | 0.25 - 0.50 |
| 9 | 2 | 1 | 0.166 | 1 | 0.50 - 0.66 |
| 1 | 7 | 0 | 0.166 | 0 | 0.666 - 0.832 |
| 2 | 3 | 1 | 0.166 | 1 | 0.832 - 1 |

→ **Normalize the weights**

4) Upsampling

| $x_1$ | $x_2$ | $y$ | wt |
|---|---|---|---|
| 9 | 2 | 1 | 1/5 |
| 1 | 5 | 1 | 1/5 |
| 2 | 3 | 1 | 1/5 |
| 1 | 5 | 1 | 1/5 |
| 2 | 4 | 0 | 1/5 |

   - choose n random # b/w 0-1

   e.g. 0.54, 0.2, 0.99, 0.06, 0.31

   - check which range they fall in & add that
     row to new dataset

→ reset weights to 1/n

$\alpha \ l_r × \alpha_i$

↳ n_estimators + learning_rate

5) Redo steps 2-4 for as many decision stumps in model.

6) $\hat{y} = \text{sign}\left(\alpha_1 h_1(x) + \alpha_2 h_2(x) + \alpha_2 h_2(x) + \ldots + \alpha_n h_n(x)\right)$

→ Boosting is a method of gradually strengthening a weak learner model.

→ In estimation stage, the results of multiple models are combined.

→ AdaBoost is an algorithm that increases accuracy by focusing more on misclassified data points during training stage.

---

• AdaBoost : Boosting – Training and Prediction stage

feature  target  initial sampling weights uniformly, with replacement
$x_1$ $x_2$ ... $y$

(weak learner) depth=1, 2 (DT, SVM, ...)

$(\hat{y}_i) = h_1(x_i)$, $y = \{-1, +1\}$

$\epsilon_1 = \sum_{i=1}^{m} D_1(i) \cdot I(y_i \neq \hat{y}_i)$ ← misclassified error = (1-accuracy)

$(\alpha_1) = \frac{1}{2} \ln\left(\frac{1-\epsilon_1}{\epsilon_1}\right)$ ← $\alpha > 0$ $\epsilon < 0.5$

$D_2(i) = \frac{1}{Z} D_1(i) \exp(-\alpha_1 y_i \hat{y}_i)$

* Z: normalization constant

Update sampling weights. The weights of misclassified data points increase, and the weights of well-classified data points decrease.

$(D_1)(i) = \left(\frac{1}{m}\right)$

row sub-sampling → sub set → training → model $h_1$ → predict whole training data

updated sampling weights non-uniform    next sampling weights

$D_2(i)$

row sub-sampling → sub set → training → model $h_2$ → predict whole training data

$\hat{y}_2 = h_2(x_2)$

(repeat) : sequential structure

[ Training Stage ]

test data $x$

model $h_1$   model $h_2$   ...   model $h_7$

$H(x) = sign\left(\sum_{t=1}^{T} \alpha_t h_t(x)\right)$

[ Prediction Stage ]
• The output of h with large ε is unreliable, so it is multiplied by small α. The output of h with small ε is reliable, so it is multiplied by large α.

synthesize results

• A weak learner is a model that performs slightly better than a random learner. For example, for DT, use max_depth 1 or 2, and for SVM, adjust C and γ so that the decision boundary is not detailed.

---

• AdaBoost : Algorithm (binary classification), (y = {-1, +1})

Given: $(x_1, y_1), ..., (x_m, y_m)$ where $x_i \in X$, $y_i \in Y = \{-1, +1\}$    ← Binary classification

Initialize $D_1(i) = 1/m$    ← Initialize the sampling weights to 1/m, uniform distribution.

For $t = 1, ..., T$
  • Train weak learner using distribution $D_t$    ← Perform sampling of training data according to weights $D_t$ to create a subset. Use the subset to train a model, weak learner ($h_t$).
  • Get weak hypothesis $h_t(x) \to \{-1, +1\}$ with error    ← Use the model ($h_t$) to predict the target class for the entire training data.
    weights : probability by density
    $\epsilon_t = Pr_{i \sim D_t}[h_t(x_i) \neq y_i]$    $\epsilon_t = \sum_{i=1}^{m} D_t(i) \cdot I(y_i \neq \hat{y}_i)$ ← calculate the error = (1-accuracy)
  • Choose $\alpha_t = \frac{1}{2}\log\left(\frac{1-\epsilon_t}{\epsilon_t}\right)$    For $\alpha > 0$, $\epsilon < 0.5$. That is, h should be slightly better than a random prediction.
  • Update:
    $D_{t+1}(i) = \frac{D_t(i)}{Z_t} \times \begin{cases} \exp(-\alpha_t) & \text{if } h_t(x_i) = y_i \\ \exp(\alpha_t) & \text{if } h_t(x_i) \neq y_i \end{cases}$    ← Update the sampling weight ($D_t$) of each data point for the next sampling.
    
    $= \frac{D_t(i)\exp(-\alpha_t y_i h_t(x_i))}{Z_t}$

    $h_t(x_i) = y_i$ : $D_{t+1}(i)$ decreases. Well-matched data points are less likely to be selected for the next round.
    $h_t(x_i) \neq y_i$ : $D_{t+1}(i)$ increases. Mismatched data points are more likely to be selected for the next round.

    where $Z_t$ is a normalization factor (chosen so that $D_{t+1}$ will be a distribution)
  Output
    $H(x) = sign\left(\sum_{t=1}^{T} \alpha_t h_t(x)\right)$    ← Input test data into the trained model ($h_t$) and synthesize the results for each h. α: the weight multiplied by each $h_t(x)$, inversely proportional to ε.

* Reference: Yoav Freund et, al., 1999, A Short Introduction to Boosting. Figure 1: The boosting algorithm AdaBoost

---

• Derivation of α and ε formula
  • Reference : Raul Rojas, 2009, AdaBoost and the Super Bowl of Classifiers A Tutorial Introduction to Adaptive Boosting

Output: $H(x) = sign\left(\sum_{t=1}^{T} \alpha_t h_t(x)\right)$    ← 1999, Yoav Freund

$C_{m-1}(x_i) = \alpha_1 k_1(x_i) + \alpha_2 k_2(x_i) + ... + \alpha_{m-1} k_{m-1}(x_i)$    ← Linear combination of weak learners.

weights for model    The (m-1)th weak learner

$C_m(x_i) = C_{m-1}(x_i) + \alpha_m k_m(x_i)$

$E = \sum_{i=1}^{N} \exp(-y_i \hat{y}_i)$    ← total loss (exponential loss)

$y_i = \{-1, +1\}$
$y_i = \hat{y}_i \to$ E decreases
$y_i \neq \hat{y}_i \to$ E increases

$E = \sum_{i=1}^{N} \exp\{-y_i(C_{m-1}(x_i) + \alpha_m k_m(x_i))\}$

$E = \sum_{i=1}^{N} w_i^{(m)} \exp\{-y_i \alpha_m k_m(x_i)\}$    ← $w_i^{(m)} = \exp\{-y_i C_{m-1}(x_i)\}$

$E = \sum_{y_i = k_m(x_i)} w_i^{(m)} e^{-\alpha_m} + \sum_{y_i \neq k_m(x_i)} w_i^{(m)} e^{\alpha_m}$

* the total cost is the weighted cost of all hits plus the weighted cost of all misses.

$E = W_c e^{-\alpha_m} + W_e e^{\alpha_m}$    $(\text{argmin } E)$

$\frac{dE}{d\alpha_m} = -W_c e^{-\alpha_m} + W_e e^{\alpha_m} = 0$

$W_c = w_1^{(m)} + w_3^{(m)} + ...$
$W_e = w_2^{(m)} + w_4^{(m)} + ...$

$\alpha_m = \frac{1}{2}\log\left(\frac{W_c}{W_e}\right)$

$\alpha = \frac{1}{2}\log\left(\frac{W - W_e}{W_e}\right)$    ← $W_c = W - W_e$, $W = W_c + W_e$

$\alpha_m = \frac{1}{2}\log\left(\frac{1-\epsilon_m}{\epsilon_m}\right)$

$W_c = W - W_e$
$W = W_c + W_e$
$\epsilon_m = \frac{W_e}{W} = \frac{\sum w_i^{(m)}}{W}$ for $y_i \neq k_m$

$\epsilon_m \leq 0.5$ ← by assumption
$\epsilon_m = 0.5 \to \alpha = 0$
$\epsilon_m = 0 \to \alpha = \infty$

$\epsilon_m = \sum_{i=1}^{N} w_i^{(m)} \cdot I(y_i \neq \hat{y}_i)$

* α and ε are inversely proportional. If the m-th learner's error is large, α is applied as a small value to reduce the influence of the m-th learner, and vice versa.

---

• Implementation of AdaBoost using DecisionTreeClassifier, (y = {-1, +1})

plot_train(x, y, w=weights[0])   plot_train(x, y, w=weights[1])   plot_train(x, y, w=weights[2])   plot_train(x, y, w=weights[3])

As t increases, the weights around the decision boundary increase. ← Misclassified data points are sampled more and trained more.

plot_train(x, y, w=weights[4])   plot_train(x, y, w=weights[5])

Changes ε, α

Multiply h with small ε by large α, and multiply h with large ε by small α.

Prediction result

$H(x) = sign\left(\sum_{t=1}^{T} \alpha_t h_t(x)\right)$

Training frequently on misclassified samples can lead to overfitting, but as frequency increases (as weights increase) alpha decreases, which reduces the likelihood of overfitting.

• AdaBoost : Algorithm (binary classification), (y = {0, 1})

• Use $y$ = {0, 1} instead of $y$ = {-1, +1}. It can be easily extended to multiclass classification.

**Algorithm 1**: AdaBoost (Freund & Schapire 1997)

1. Initialize the observation weights $w_i = 1/n$, $i = 1, 2, ..., n$
2. **For** $m$ = 1 **to** $M$:
   a) Fit a classifier $T^{(m)}(x)$ to the training data using weights $w_i$
   b) Compute
   $$err^{(m)} = \sum_{i=1}^{n} w_i \cdot I(c_i \neq T^{(m)}(x_i)) / \sum_{i=1}^{n} w_i$$
   c) Compute
   $$\alpha^{(m)} = \log \frac{1 - err^{(m)}}{err^{(m)}}$$
   d) Set
   $$w_i \leftarrow w_i \cdot \exp(\alpha^{(m)} \cdot I(c_i \neq T^{(m)}(x_i))), \quad i = 1, 2, ..., n$$
   e) Re-normalize $w_i$
3. Output
   $$C(x) = \underset{k}{argmax} \sum_{m=1}^{M} \alpha^{(m)} \cdot I(T^{(m)}(x) = k)$$

* Reference: Ji Zhu, et al., 2006, Multi-class AdaBoost, Algorithm 1

In this form, multiple class classification is also possible, but $\alpha < 0$ problem occurs if $err^{(m)} > 0.5$. In binary classification, $\alpha < 0.5$ is possible even in weak learners, but in three classes $\alpha < 0$ occurs easily because the accuracy of random prediction is $1/3$ and $err^{(m)} = 0.67$ → Need to be improved → SAMME algorithm in the next video.

It is always 1 because it is normalized.

Not multiplied by 1/2. The reason was not explained.

Use argmax instead of the sign function.

(example)

| | | | |
|---|---|---|---|
| m=1: | 0 | [1 0] · 0.5 = [0.5  0] |
| m=2: | 0 | [1 0] · 0.7 = [0.7  0] |
| m=3: | 1 | [0 1] · 0.2 = [0   0.2] |
| m=4: | 0 | [1 0] · 0.6 = [0.6  0] |

sum = [1.8  0.2]  argmax = 0 → C(x)

---

• Multiclass Classification – SAMME Algorithm (Stagewise Additive Modeling using a Multi-class Exponential loss function)

• In 2006, Ji Zhu et al. proposed the following algorithm to solve the multiclass problem of Freund's (1995) AdaBoost algorithm.
• The $err(m) < 0.5$ condition in Freund's model was modified to $err(m) < 1 - 1/K$ (K: the number of classes). If K=3, then $err(m) < 0.67$.

### Multi-class AdaBoost

Ji Zhu, Department of Statistics University of Michigan Ann Arbor, MI 48109
Saharon Rosset, Data Analytics Group IBM Research Center Yorktown Heights, NY 10598
Hui Zou, School of Statistics University of Minnesota Minneapolis, MN 55455
Trevor Hastie, Department of Statistics Stanford University Stanford, CA 94305
January 12, 2006

**Abstract**

Boosting has been a very successful technique for solving the two-class classification problem. In going from two-class to multi-class classification, most algorithms have been restricted to reducing the multi-class classification problem to multiple two-class problems. In this paper, we propose a new algorithm that naturally extends the original AdaBoost algorithm to the multiclass case without reducing it to multiple two-class problems. Similar to AdaBoost in the two-class case, this new algorithm combines weak classifiers and only requires the performance of each weak classifier be better than random guessing (rather than 1/2). We further provide a statistical justification for the new algorithm using a novel multi-class exponential loss function and forward stage-wise additive modeling. As shown in the paper, the new algorithm is extremely easy to implement and is highly competitive with the best currently available multi-class classification methods.

→ **Algorithm 2: SAMME**    $y = \{0, 1\} \rightarrow y = \{0, 1, 2, ... K-1\}$

1. Initialize the observation weights $w_i = 1/n$, $i = 1, 2, ..., n$
2. **For** m = 1 **to** M:
   a) Fit a classifier $T^{(m)}(x)$ to the training data using weights $w_i$    if K=2, $err < 1 - \frac{1}{2} (0.5)$
   b) Compute
   $$err^{(m)} = \sum_{i=1}^{n} w_i \cdot I(c_i \neq T^{(m)}(x_i)) / \sum_{i=1}^{n} w_i$$    if K=3, $err < 1 - \frac{1}{3} (0.67)$
   c) Compute
   $$\alpha^{(m)} = \log \frac{1 - err^{(m)}}{err^{(m)}} + \log(K - 1)$$    This part was added to the Freund (1997) model.
   d) Set
   $$w_i \leftarrow w_i \cdot \exp(\alpha^{(m)} \cdot I(c_i \neq T^{(m)}(x_i))), \quad i = 1, 2, ..., n$$
   For $\alpha^{(m)} > 0$, $err^{(m)} < 1 - \frac{1}{K}$
   e) Re-normalize $w_i$
3. Output
   $$C(x) = \underset{k}{argmax} \sum_{m=1}^{M} \alpha^{(m)} \cdot I(T^{(m)}(x) = k)$$

$$\log \frac{1 - err^{(m)}}{err^{(m)}} + \log(K - 1) > 0$$

proof: $\frac{1 - err^{(m)}}{err^{(m)}} > \frac{1}{K-1}$

$\frac{1}{err^{(m)}} > \frac{K}{K-1}$

---

• AdaBoost Regression - Algorithm

**3. BOOSTING:** In bagging, each training example is equally likely to be picked. In boosting, the probability of a particular example being in the training set of a particular machine depends on the performance of the prior machines on that example. The following is a modification of Adaboost.R [Freund and Schapire (1996a)].

① Initially, to each training pattern we assign a weight $w_i = 1/N_1$, $i = 1, ..., N_1$. Pick $N_1$ samples (with replacement) to form the training set.

Repeat the following while the average loss L defined below is less than 0.5.

2. Construct a regression machine t from that training set. Each machine makes a hypothesis $h_t$: $x \rightarrow y$

③ Pass every member of the training set through this machine to obtain a prediction $y_i^{(p)}(x_i)$   $i = 1, ..., N_1$

④ Calculate a loss for each training sample $(L) = L[|y_i^{(p)} - y_i|]$. The loss L may be of any functional form as long as $L \in [0, 1]$. If we let
$$D = \sup |y_i^{(p)}(x_i) - y_i| \quad i = 1, ..., N_1$$

then we have three candidate loss functions:

$$L_i = \frac{|y_i^{(p)} - y_i|}{D} \quad \text{(linear)}$$
$$L_i = \frac{|y_i^{(p)} - y_i|^2}{D^2} \quad \text{(square law)}$$
$$L_i = 1 - \exp\left(\frac{-|y_i^{(p)} - y_i|}{D}\right) \quad \text{(exponential)}$$

* classification
$$err^{(m)} = \sum_{i=1}^{n} w_i \cdot I(c_i \neq T^{(m)}(x_i)) / \sum_{i=1}^{n} w_i$$
$$\alpha^{(m)} = \log \frac{1 - err^{(m)}}{err^{(m)}}$$
$$w_i \leftarrow w_i \cdot \exp(\alpha^{(m)} \cdot I(c_i \neq T^{(m)}(x_i)))$$

⑤ Calculate an average loss:  $\bar{L} = \sum_{i=1}^{n} L_i w_i$

⑥ Form $\beta = \frac{\bar{L}}{1 - \bar{L}}$   ← (Proportional to average loss)

$\beta$ is a measure of confidence in the predictor. Low $\beta$ means high confidence in the prediction.

⑦ Update the weights: $w_i \rightarrow w_i \beta^{**} (1 - L_i)$, where ** indicates exponentiation. The smaller the loss, the more the weight is reduced, making the probability smaller that this pattern will be picked as a member of the training set for the next machine in the ensemble.

⑧ For a particular input $x_i$, each of the T machines makes a prediction $h_t$, t=1, ..., T. Obtain the cumulative prediction $h_f$ using the T predictors:
$$h_f = \inf\{y \in Y: \sum_{t: h_t \leq y} \log(\frac{1}{\beta_t}) \geq \frac{1}{2} \sum_t \log(\frac{1}{\beta_t})\}$$

This is the weighted median

Weight of each machine ($\alpha$ role in classification). It is inversely proportional to the average loss. The weight of machines with large losses is reduced.

---

• AdaBoost Regression - Algorithm

• In algorithm step 8, the weighted median is used for prediction.
• Here we will try both weighted average and weighted median.

model weights
$$w_t = \log(\frac{1}{\beta_t}) / \sum_t \log(\frac{1}{\beta_t}) \leftarrow \text{Weight for the value estimated by each model}$$

① Using weighted average
$$h_f = \sum \hat{y}_t w_t \leftarrow \text{weighted average}$$

# Method-1: Using weighted average
★ wavg_pred = np.sum(y_pred * w, axis=1)

② Using weighted median as in the paper
$$h_f = \inf\{y \in Y: \sum_{t: h_t \leq y} \log(\frac{1}{\beta_t}) \geq \frac{1}{2} \sum_t \log(\frac{1}{\beta_t})\}$$
$$h_f = \inf\{y \in Y: \sum_{t: h_t \leq y} w_t \geq \frac{1}{2} \sum_t w_t\}$$

If $y$ is sorted in ascending order, the $y$ at the position where the sum of the lower w is half of the total sum of w becomes the final output.

The y value estimated by the first model.
the number of models (ex. T = 50)

ex)
y_pred: [0 1 2 3 4 5 6 7 8 9 10 11 ...  49]
         $\hat{y}_0$ $\hat{y}_1$ $\hat{y}_2$ ...

argsort, ascending

i_pred: [3 0 9 11 5 4 8 6 1 10 2 7 ... 34]

w: [$w_3$ $w_0$ $w_9$ $w_{11}$ $w_5$ ... ]  models

accumulated w: [$a_0$ $a_1$ $a_2$ $a_3$ $a_4$ ... $a_{49}$]  (w-acc)

is ($a_i \geq 0.5 \times a_{49}$)?

False or True: [F F F ...  F T T T ...]  the position of weighted median (ex. 23)

weighted median of y_pred = $\hat{y}_{34}$ ← final estimate

# weighted median: (sum of the lower w ≥ half of the total sum of w)
i_pred = np.argsort(y_pred, axis=1)
w_acc = np.cumsum(w[i_pred], axis=1)     # accumulated w
is_med = w_acc >= 0.5 * w_acc[:, -1][:, np.newaxis]
i_med = is_med.argmax(axis=1)      # 23
y_med = i_pred[np.arange(n_test), i_med]  # 34
wmed_pred = np.array(y_pred[np.arange(n_test), y_med]) # final estimate