# Ensemble Methods (Adaboost and GBDT)

Prof.Mingkui Tan

South China University of Technology
Southern Artificial Intelligence Laboratory(SAIL)

November 7, 2017

# Content

# Contents

## Ensemble Learning

- Ensemble learning: Combine numerous weak learners to a strong learner
- Main methods: Boosting, Bagging

## Boosting Method

---

**Algorithm 1:** Adaboost

---

**Input:** $D = \{(\mathbf{x}_1, y_1), ..., (\mathbf{x}_n, y_n)\}$, **where** $\mathbf{x}_i \in X, y_i \in \{-1, 1\}$

**Initialize:** Sample weight distribution $D_1 = \frac{1}{n}$

1 Train a base learner $h_1(\mathbf{x})$ with $D_1$

2 **for** m=2,3,...,M **do**

3     Update the sample distribution $D_m$, to make the wrong predictive samples more important

4     Train a new base learner $h_m(\mathbf{x})$ with $D_m$

5 **end**

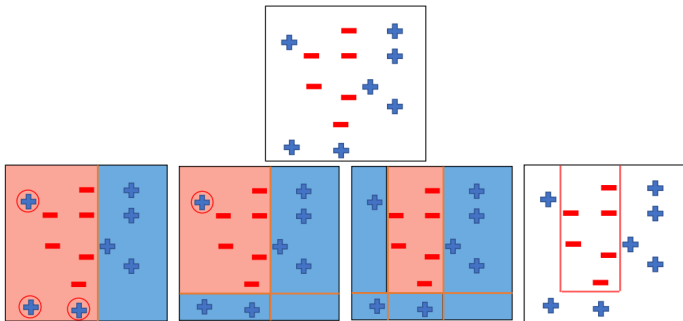**Output:** $H(\mathbf{x}) = \sum_{m=1}^{M} \alpha_m h_m(\mathbf{x})$

---

# Contents

# Adaboost
## How to train the base leaner?

Make the wrong predictive samples more important, and handle it in next round:

# Adaboost
## Sample weight updating formula

$$w_{m+1}(i) = \frac{w_m(i)}{z_m} e^{-\alpha_m y_i h_m(\mathbf{x}_i)}$$

$z_m = \sum_{i=1}^{n} w_m(i) e^{-\alpha_m y_i h_m(\mathbf{x}_i)}$ is normalization term, makes $w_m(i)$ become probability distributions
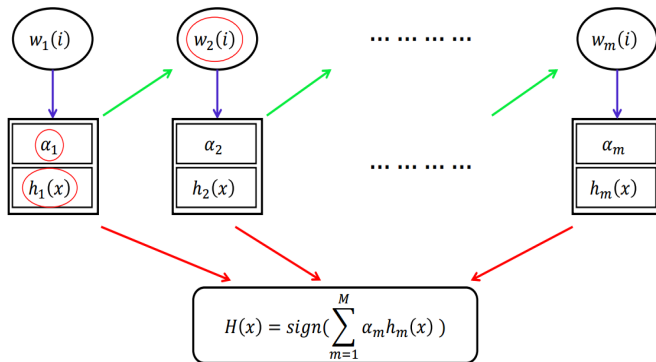
$$w_{m+1}(i) = \begin{cases} \dfrac{w_m(i)}{z_m} e^{-\alpha_m} & \text{for right predictive sample} \\[2ex] \dfrac{w_m(i)}{z_m} e^{\alpha_m} & \text{for wrong predictive sample} \end{cases}$$

so in next round, $\frac{w_{wrong}(i)}{w_{right}(i)} = e^{2\alpha_m} = \frac{1 - \epsilon_m}{\epsilon_m}$ and $\epsilon_m < 0.5$, wrong samples will be more important

## Adaboost
How to combine the base learner?

Every iteration generates a new base learner $h_m(\mathbf{x})$ and its importance score $\alpha_m$



$$H(x) = sign(\sum_{m=1}^{M} \alpha_m h_m(x) )$$

## Adaboost
Evaluate the performance of the base learner

- Base learner

$$h_m(\mathbf{x}) : \mathbf{x} \mapsto \{-1, 1\}$$

- Error rate

$$\epsilon_m = p(h_m(\mathbf{x}_i) \neq y_i) = \sum_{i=1}^{n} w_m(i) \mathbb{I}(h_m(\mathbf{x}_i) \neq y_i)$$

$\epsilon_m$ ¡0.5, or the performance of Adaboost is weaker than random classfication.

## Adaboost
Importance score of base learner

Make the base learner with lower $\epsilon_m$ more important

$$\alpha_m = \frac{1}{2} \log \frac{1 - \epsilon_m}{\epsilon_m}$$

## Adaboost
### Additive model

- Final learner

$$H(\mathbf{x}) = \text{sign}(\sum_{m=1}^{M} \alpha_m h_m(\mathbf{x}))$$

Note: $h_m(\mathbf{x}) = \text{sign}(\mathbf{w}^\top \mathbf{x})$ is a nonlinear function, so the Adaboost can deal with nonlinear problem

## Algorithm

---

**Algorithm 2:** Adaboost

---

**Input:** $D = \{(\mathbf{x}_1, y_1), ..., (\mathbf{x}_n, y_n)\}$, **where** $\mathbf{x}_i \in X, y_i \in \{-1, 1\}$

**Initialize:** Sample distribution $w_m$

**Base learner:** $\mathcal{L}$

1  $w_1(i) = \frac{1}{n}$

2  **for** m=1,2,...,M **do**

3      $h_m(x) = \mathcal{L}(D, w_m)$

4      $\epsilon_m = \sum_{i=1}^{n} w_m(i) \mathbb{I}(h_m(\mathbf{x}_i) \neq y_i)$

5      **if** $\epsilon_m > 0.5$ **then**

6          |   **break**

7      **end**

8      $\alpha_m = \frac{1}{2} \log \frac{1 - \epsilon_m}{\epsilon_m}$

9      $w_{m+1}(i) = \frac{w_m(i)}{z_m} e^{-\alpha_m y_i h_m(\mathbf{x}_i)}$,**where** $i = 1, 2, ..., n$ **and**

        $z_m = \sum_{i=1}^{n} w_m(i) e^{-\alpha_m y_i h_m(\mathbf{x}_i)}$

10  **end**

**Output:** $H(\mathbf{x}) = \sum_{m=1}^{M} \alpha_m h_m(\mathbf{x})$

---

# Contents

# Gradient Boosting Decision Trees
GBDT is a decision tree algorithm with iteration

Example: What is the difference between regression tree and GBDT?

Suppose: There are 4 peoples A, B, C and D, whose age are 14, 16, 24, 26 respectively.
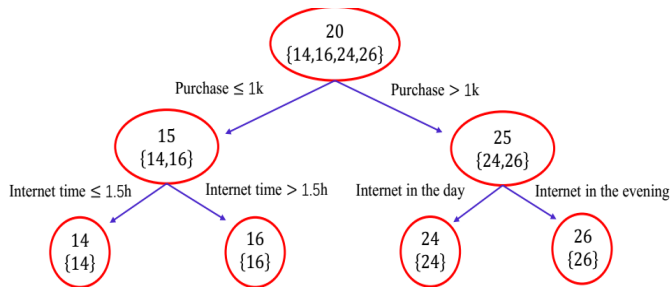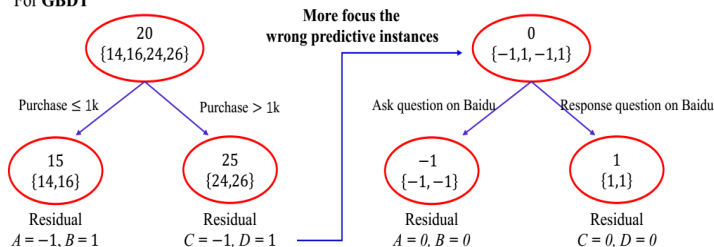


Figure: Single regression tree

# Gradient Boosting Decision Trees

- The key of GBDT is that trees learn all the results and residuals of all trees before.
- The residual is the difference of predictive value and real value, so the predictive value is the sum of all results of trees.



- So,
  A=15+(-1)=14 B=15+1=16 C=25+(-1)=24 D=25+1=26

# Gradient Boosting Decision Trees
Question

Q1:when results of these two algorithms are same, Why do we choose GBDT?

- The motivation of this algorithm is that every calculation of residual is to increase the weight of wrong predictive samples,and the residual of right predictive sample is zero.

- So in the next iteration, model can concentratively address these wrong predictive samples. Another function is to prevent over fitting.

# Gradient Boosting Decision Trees
### Question

Q2: Where does this algorithm reflect gradient boosting?

- In algorithm, residual is the gradient descent direction, which is the derivation of mean square error(MSE). Actually, MSE is the loss function of CART regression tree.

| Setting | Loss Function | $-\partial L(y_i, f(x_i))/\partial f(x_i)$ |
|---|---|---|
| Regression | $\frac{1}{2}[y_i - f(x_i)]^2$ | $y_i - f(x_i)$ |
| Regression | $\lvert y_i - f(x_i) \rvert$ | $\mathrm{sign}[y_i - f(x_i)]$ |
| Regression | Huber | $y_i - f(x_i)$ for $\lvert y_i - f(x_i) \rvert \leq \delta_m$ |
| | | $\delta_m \mathrm{sign}[y_i - f(x_i)]$ for $\lvert y_i - f(x_i) \rvert > \delta_m$ |
| | | where $\delta_m = \alpha$th-quantile$\{\lvert y_i - f(x_i) \rvert\}$ |
| Classification | Deviance | $k$th component: $I(y_i = \mathcal{G}_k) - p_k(x_i)$ |

## Algorithm

### Algorithm 3: GBDT

**Input:** $D = \{(\mathbf{x}_1, y_1), ..., (\mathbf{x}_n, y_n)\}$, **where** $\mathbf{x}_i \in X, y_i \in \{-1, 1\}$
**Initialize:** $f_0(x) = \arg\min_\mu \sum_{i=1}^n L(y_i, \mu)$

1 **for** m=1,2,...,M **do**
2      **for** i=1,2,...,n **do**
3          $r_{im} = - \left[ \frac{\partial L(y_i, f_{m-1}(\mathbf{x}_i))}{\partial f_{m-1}(\mathbf{x}_i)} \right]$
4          Fit a regression tree to targets $r_{im}$ giving terminal regions
            $R_{jm}, j = 1, 2, ..., J_m$
5      **end**
6      **for** j=1,2,...,$J_m$ **do**
7          $\mu_{jm} = \arg\min_\mu \sum_{\mathbf{x}_i \in R_{jm}} L(y_i, f_{m-1}(\mathbf{x}_i) + \mu), j = 1, 2, ..., J_m$
8          Update $f_m(\mathbf{x}) = f_{m-1}(\mathbf{x}) + \sum_{j=1}^{J_m} \mu_{jm} \mathbb{I}(\mathbf{x} \in R_{jm})$
9      **end**
10 **end**
**Output:** $\hat{f}(\mathbf{x}) = f_M(\mathbf{x})$

# THANK YOU!