

# Ensemble Model

Bùi Tiến Lên

2023



KHOA CÔNG NGHỆ THÔNG TIN  
TRƯỜNG ĐẠI HỌC KHOA HỌC TỰ NHIÊN

# Contents

---



1. Ensemble Model
2. Bagging
3. Boosting
4. Gradient Boosting



# Notation

## Bagging

Bootstrap

Algorithm

Random Forests

## Boosting

AdaBoost

Face Detection

## Gradient Boosting

symbol	meaning
$a, b, c, N \dots$	scalar number
$\mathbf{w}, \mathbf{v}, \mathbf{x}, \mathbf{y} \dots$	column vector
$\mathbf{X}, \mathbf{Y} \dots$	matrix
$\mathbb{R}$	set of real numbers
$\mathbb{Z}$	set of integer numbers
$\mathbb{N}$	set of natural numbers
$\mathbb{R}^D$	set of vectors
$\mathcal{X}, \mathcal{Y}, \dots$	set
$\mathcal{A}$	algorithm

operator	meaning
$\mathbf{w}^T$	transpose
$\mathbf{X}\mathbf{Y}$	matrix multiplication
$\mathbf{X}^{-1}$	inverse

# Big Picture





# Ensemble Model



# Bias vs. Variance

## Bagging

Bootstrap

Algorithm

Random Forests

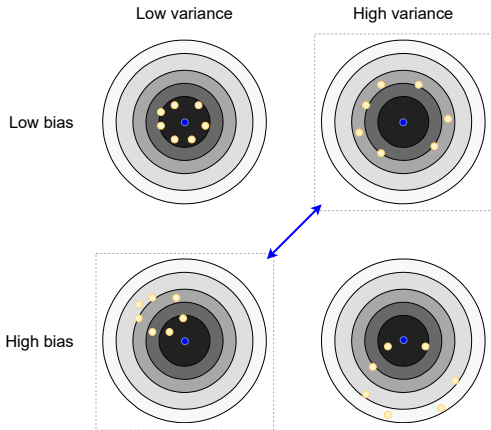
## Boosting

AdaBoost

Face Detection

## Gradient Boosting

- Low-bias models tend to have high variance, and vice versa.





# Basics of Ensembles

## Bagging

Bootstrap

Algorithm

Random Forests

## Boosting

AdaBoost

Face Detection

Gradient  
Boosting

## Concept 1

Instead of providing **one model**, an **ensemble** approach proposes **many models** to the same problem, and **combine** them

- The simplest ensemble  $H$  over models  $\{h_i \in \mathcal{H}, i = 1 \dots T\}$

$$H(\mathbf{x}) = \sum_{i=1}^T \alpha_i h_i(\mathbf{x}) \text{ with } \sum_{i=1}^T \alpha_i = 1 \quad (1)$$

- Why should this be a good idea?
  - combine models  $\rightarrow$  reduce the variance  $\rightarrow$  enhance expected performance.
- However, increase the performance cost

# Basics of Ensembles (cont.)



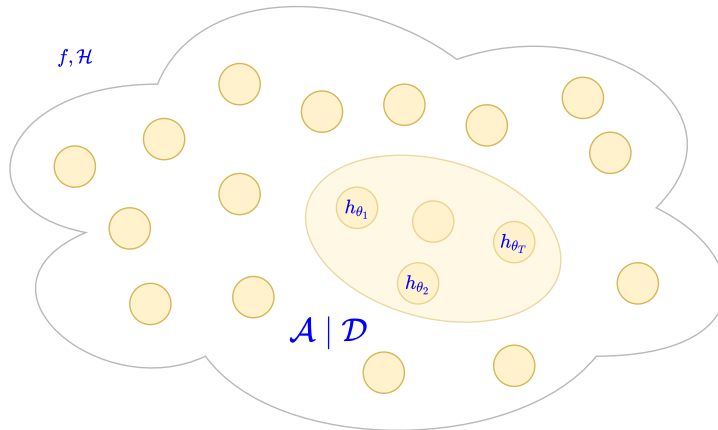
## Bagging

Bootstrap  
Algorithm  
Random Forests

## Boosting

AdaBoost  
Face Detection

## Gradient Boosting







# Why Does it Work?

## Bagging

Bootstrap

Algorithm

Random Forests

## Boosting

AdaBoost

Face Detection

## Gradient Boosting

It has been shown that the expected risk of the average of a set of models is better than the average of the expected risk of these models

- Let us consider the simplest ensemble  $H$  over models  $h_i$

$$H(\mathbf{x}) = \sum_{i=1}^T \alpha_i h_i(\mathbf{x}) \text{ with } \sum_{i=1}^T \alpha_i = 1 \quad (2)$$

- The MSE risk of  $h_i$  at  $\mathbf{x}$  is

$$e_i(\mathbf{x}) = \mathbb{E}_y[(y - h_i(\mathbf{x}))^2] \quad (3)$$



# Why Does it Work? (cont.)

## Bagging

Bootstrap

Algorithm

Random Forests

## Boosting

AdaBoost

Face Detection

## Gradient Boosting

- The average risk  $\bar{e}(\mathbf{x})$  of a model is

$$\bar{e}(\mathbf{x}) = \sum_i \alpha_i e_i(\mathbf{x}) \quad (4)$$

- The average risk  $e(\mathbf{x})$  of the ensemble is

$$e(\mathbf{x}) = \mathbb{E}_y[(y - H(\mathbf{x}))^2] \quad (5)$$

- Let us define diversity

$$d_i(\mathbf{x}) = (h_i(\mathbf{x}) - H(\mathbf{x}))^2 \quad (6)$$

- The average diversity is

$$\bar{d}(\mathbf{x}) = \sum_i \alpha_i d_i(\mathbf{x}) \quad (7)$$



# Why Does it Work? (cont.)

---

## Bagging

Bootstrap

Algorithm

Random Forests

## Boosting

AdaBoost

Face Detection

## Gradient Boosting

- It can then be shown that

$$e(\mathbf{x}) = \bar{e}(\mathbf{x}) - \bar{d}(\mathbf{x}) \quad (8)$$

$$e(\mathbf{x}) < \bar{e}(\mathbf{x}) \quad (9)$$



# Bagging

- Bootstrap
- Algorithm
- Random Forests



# Bagging

## Bagging

Bootstrap

Algorithm

Random Forests

## Boosting

AdaBoost

Face Detection

## Gradient Boosting

## Underlying idea

A part of the **variance** is due to the specific choice of the training data set

- Let us create many **similar** training data sets,
  - For each of them, let us train a new function  $f_i$
  - The final function will be the *average* of each function outputs.
- 
- How similar? using **bootstrap aggregating**



# Bootstrap

## Bagging

## Bootstrap

## Algorithm

## Random Forests

## Boosting

## AdaBoost

## Face Detection

Gradient  
Boosting

## Concept 2

Given a data set  $\mathcal{D}_n$  with  $n$  examples drawn from  $p(\mathcal{Z}) = p(\mathcal{X}, \mathcal{Y})$ , a bootstrap  $\mathcal{B}_i, i = 1 \dots T$  of  $\mathcal{D}_n$  also contains  $n$  examples:

**for**  $j = 1 \rightarrow n$

the  $j$ -th example of  $\mathcal{B}_i$  is drawn independently with replacement from  $\mathcal{D}_n$

- Some examples from  $\mathcal{D}_n$  are in multiple copies in  $\mathcal{B}_i$
- Some examples from  $\mathcal{D}_n$  are not in  $\mathcal{B}_i$
- The examples were i.i.d. drawn from  $p(\mathcal{Z}) \rightarrow$  the datasets  $\mathcal{B}_i$  are as plausible as  $\mathcal{D}_n$ , but drawn from  $\mathcal{D}_n$  instead of  $p(\mathcal{Z})$ .



# Example

## Bagging

### Bootstrap

#### Algorithm

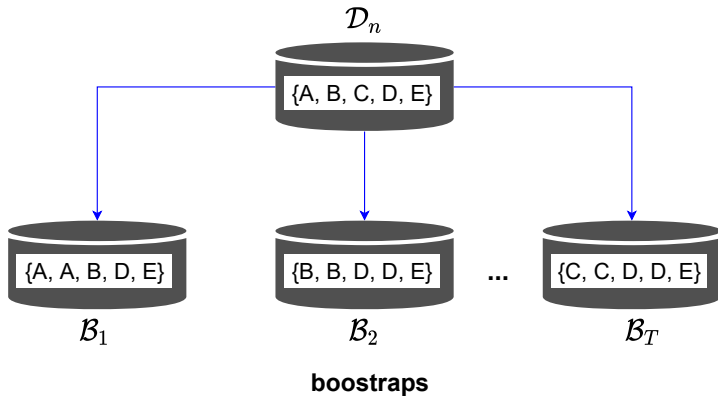
#### Random Forests

## Boosting

### AdaBoost

#### Face Detection

## Gradient Boosting





# Diagram

## Bagging

Bootstrap

Algorithm

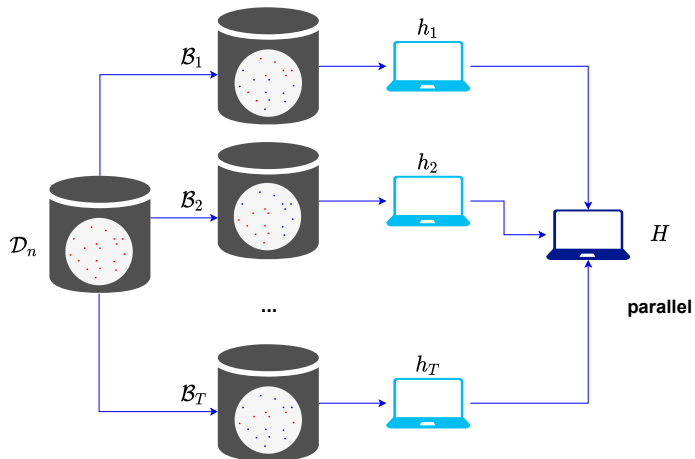
Random Forests

## Boosting

AdaBoost

Face Detection

## Gradient Boosting







# Algorithm

## Bagging

Bootstrap

Algorithm

Random Forests

## Boosting

AdaBoost

Face Detection

Gradient  
Boosting**Training:**

- Given a training set  $\mathcal{D}_n$ , create  $T$  bootstraps  $\mathcal{B}_i$  of  $\mathcal{D}_n$
- For each bootstrap  $\mathcal{B}_i$ , select

$$h_i = \arg \min_{h \in \mathcal{H}} E(h \mid \mathcal{B}_i) \quad (10)$$

**Running:**

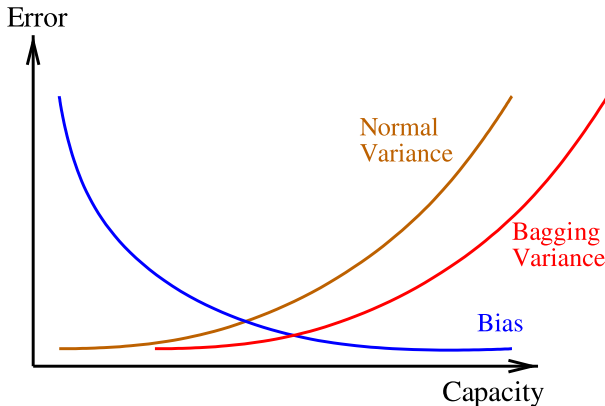
- Given an input  $\mathbf{x}$ , the corresponding output  $\hat{y}$  is:

$$\hat{y} = H(\mathbf{x}) = \frac{1}{T} \sum_{i=1}^T h_i(\mathbf{x}) \quad (11)$$



# Bias + Variance

- **Analysis:** if generalization error is decomposed into **bias** and **variance** terms then **bagging reduces variance**.





# Random Forests

## Concept 3

A **random forest** is an ensemble of decision trees.





# Random Forests (cont.)

Each decision tree  $h_i$  is trained as follows:

- Create a **bootstrap** of the training set
- Select a subset  $m \ll d$  input variables as **potential** split nodes ( $m$  is constant over all trees)
- No pruning of the trees

A decision is taken by **voting** amongst the trees

- Somehow,  $m$  controls the capacity.



# Boosting

- AdaBoost
- Face Detection



# Big Picture

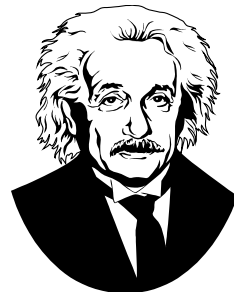
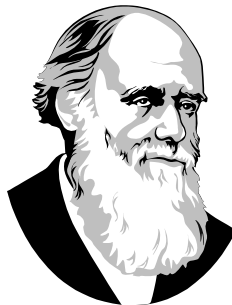
## Bagging

- Bootstrap
- Algorithm
- Random Forests

## Boosting

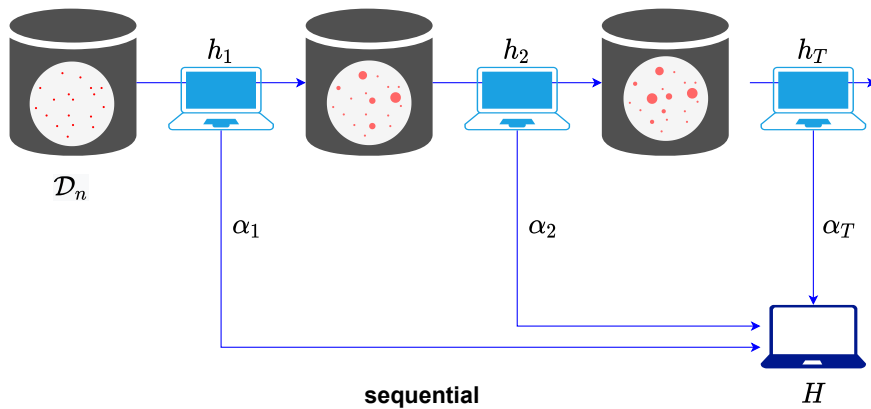
- AdaBoost
- Face Detection

## Gradient Boosting





# Diagram





# Weak vs. Strong Learning Model

## Concept 4

A learning model is **strong** iff every hypothesis  $h$  has low error

## Concept 5

A learning model is **weak** iff every hypothesis  $h$  has high error

Examples of weak classifiers:

- Simple decision trees such as **stumps**
- Simple neural networks such as **perceptrons**
- Haar-like features





# Boosting

## Bagging

Bootstrap  
Algorithm  
Random Forests

## Boosting

AdaBoost  
Face Detection

## Gradient Boosting

Boosting involves three elements:

- A loss function to be optimized  $\ell(.,.)$
- A set of weak learners  $\{h_t(x)\}$
- An additive model  $H(x)$  to add weak learners to minimize the loss function

$$H(x) = \sum_{t=1}^T \alpha_t h_t(x) \quad (12)$$



# Loss function

## Bagging

Bootstrap  
Algorithm  
Random Forests

## Boosting

AdaBoost  
Face Detection

## Gradient Boosting

- Square loss

$$\ell(\hat{y}, y) = (\hat{y} - y)^2 \quad (13)$$

- Absolute loss

$$\ell(\hat{y}, y) = |\hat{y} - y| \quad (14)$$

- Huber loss

$$\ell(\hat{y}, y) = \begin{cases} \frac{1}{2}(\hat{y} - y)^2 & (\hat{y} - y) \leq \delta \\ \delta(|\hat{y} - y| - \delta/2) & (\hat{y} - y) > \delta \end{cases} \quad (15)$$

- Exponential loss

$$\ell(\hat{y}, y) = e^{-\hat{y}y} \quad (16)$$



# Algorithm

## Bagging

Bootstrap

Algorithm

Random Forests

## Boosting

AdaBoost

Face Detection

## Gradient Boosting

- **Initialize**  $H_0 \leftarrow \emptyset$  or  $\alpha_0$
- At each time step  $t$ ,
  - **Choose**  $h_t$  given the performance obtained by previous  $H_{t-1}$ .
  - **Train** a new weak classifier
  - **Find** the new weight  $\alpha_t$  by minimizing the loss function

$$H_t \leftarrow H_{t-1} + \alpha_t h_t \quad (17)$$



# AdaBoost

## Concept 6

**AdaBoost**, short for Adaptive Boosting, is the most popular algorithm in the family of **boosting** algorithms

- **Simplest framework:** binary classification  $H(x)$
- **Simplest requirement:** each weak classifier  $y = h_t(x)$ ,  $y \in \{-1, +1\}$  should perform better than chance
- **Loss function:**

$$\ell((H(x), y)) = e^{-yH(x)} \quad (18)$$



# Algorithm

## Bagging

Bootstrap

Algorithm

Random Forests

## Boosting

AdaBoost

Face Detection

Gradient  
Boosting

**Inputs:** training data  $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$  and a set of weak binary classifiers  $\{h_i \in \mathcal{H}\}$

**Initialize** the weights' distribution of training data

$$(w_1^{(1)}, w_2^{(1)}, \dots, w_N^{(1)}) = \left( \frac{1}{N}, \frac{1}{N}, \dots, \frac{1}{N} \right) \quad (19)$$

**Iterate** over  $t = 1, 2, \dots, T$ , use training data with current weights' distribution

1. **Find** a weak classifier  $h_t(\mathbf{x})$  that minimizes the error rate  $e_t$  of over the training data

$$e_t = P(h_t(\mathbf{x}_i) \neq y_i) = \sum_{i=1}^N w_i^{(t)} \mathbb{I}(h_t(\mathbf{x}_i) \neq y_i) \quad (20)$$



# Algorithm (cont.)

## Bagging

Bootstrap

Algorithm

Random Forests

## Boosting

AdaBoost

Face Detection

Gradient  
Boosting

2. **Compute** the weight of classifier  $h_t(\mathbf{x})$

$$\alpha_t = \frac{1}{2} \log \frac{1 - e_t}{e_t} \quad (21)$$

3. **Update** the weights' distribution of training data

$$w_i^{(t+1)} = w_i^{(t)} \exp(-\alpha_t y_i h_t(\mathbf{x}_i)) \quad (22)$$

4. **Normalize** the weights of data points

$$w_i = \frac{w_i}{\sum_i w_i} \quad (23)$$

Ensemble  $T$  weak classifiers

$$\text{sign}[H(\mathbf{x})] = \text{sign} \left[ \sum_{t=1}^T \alpha_t h_t(\mathbf{x}) \right] \quad (24)$$



# Example

## Bagging

Bootstrap

Algorithm

Random Forests

## Boosting

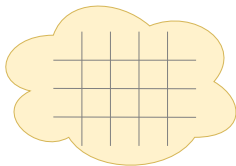
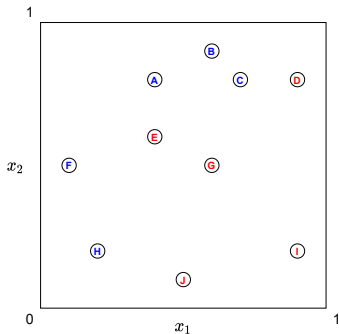
AdaBoost

Face Detection

## Gradient Boosting

- Given a training data set  $\mathcal{D} = \{A, B, C, D, E, F, G, H, I, J\}$ , find a strong classifier from weak classifiers (vertical or horizontal lines)

#	$x_1$	$x_2$	label
A	0.4	0.8	1
B	0.6	0.9	1
C	0.7	0.8	1
D	0.9	0.8	-1
E	0.4	0.6	-1
F	0.1	0.5	1
G	0.6	0.5	-1
H	0.2	0.2	1
I	0.9	0.2	-1
J	0.5	0.1	-1





# Round 1

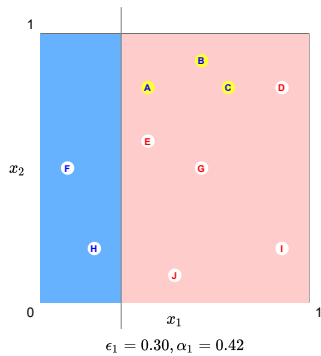
## Bagging

Bootstrap  
Algorithm  
Random Forests

## Boosting

**AdaBoost**  
Face Detection

## Gradient Boosting







# Round 2

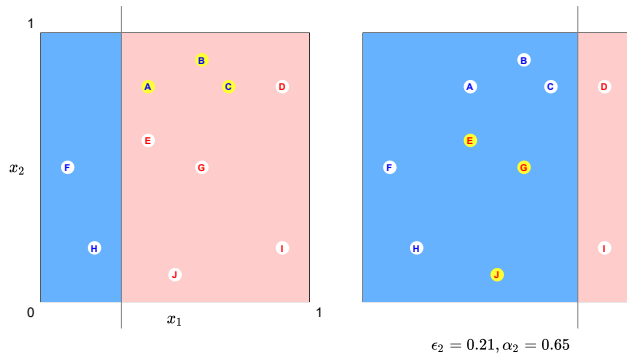
## Bagging

Bootstrap  
Algorithm  
Random Forests

## Boosting

**AdaBoost**  
Face Detection

## Gradient Boosting





# Round 3

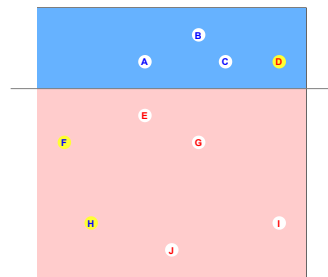
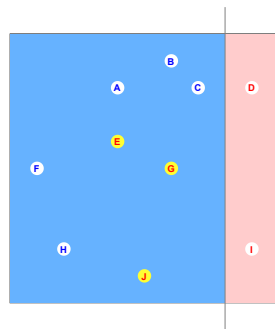
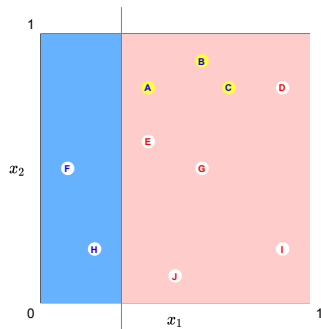
## Bagging

Bootstrap  
Algorithm  
Random Forests

## Boosting

AdaBoost  
Face Detection

## Gradient Boosting



$$\epsilon_3 = 0.14, \alpha_3 = 0.92$$



# The combined classifier

## Bagging

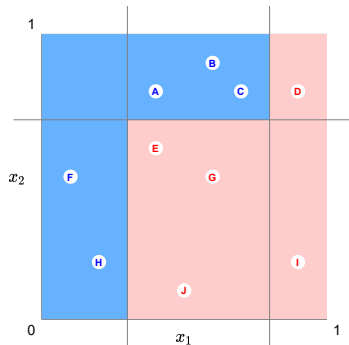
Bootstrap  
Algorithm  
Random Forests

## Boosting

AdaBoost  
Face Detection

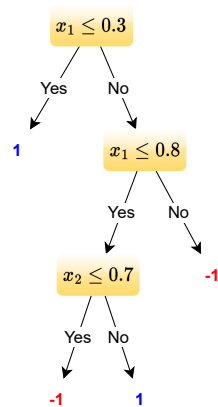
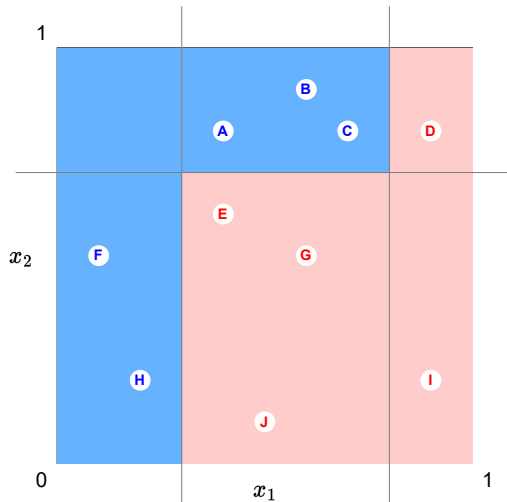
## Gradient Boosting

$$H_{final} = \text{sign} \left[ 0.42 \begin{array}{|c|} \hline \text{Blue} \\ \hline \end{array} + 0.65 \begin{array}{|c|} \hline \text{Blue} \\ \hline \end{array} + 0.92 \begin{array}{|c|} \hline \text{Blue} \\ \hline \end{array} \right]$$





# Tree based classifier





# Analysis

## Bagging

Bootstrap

Algorithm

Random Forests

## Boosting

AdaBoost

Face Detection

Gradient  
Boosting

- Selection of  $\alpha_t$  comes from minimizing

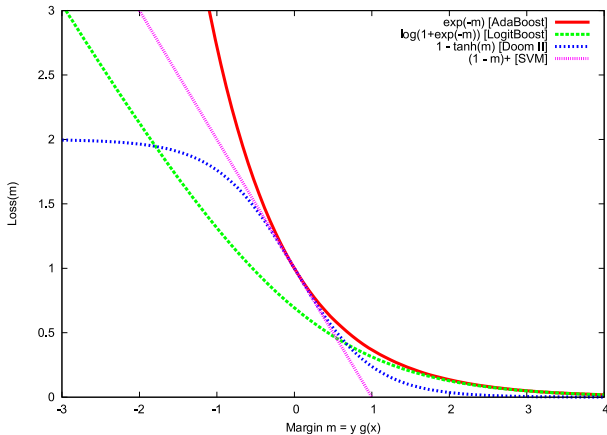
$$\arg \min_{\alpha_t} \sum_{i=1}^N \exp(-y_i [H_{t-1}(\mathbf{x}_i) + \alpha_t h_t(\mathbf{x}_i)]) \quad (25)$$

- If each weak classifier is always better than chance, then AdaBoost can be proven to **converge to 0 training error**
- Even after training error is 0, generalization error continues to improve: the **margin** continues to grow
- **Sampling** can often be replaced by **weighting**



# Cost Functions

- Comparison of various cost functions related to AdaBoost





# Margin

## Bagging

Bootstrap

Algorithm

Random Forests

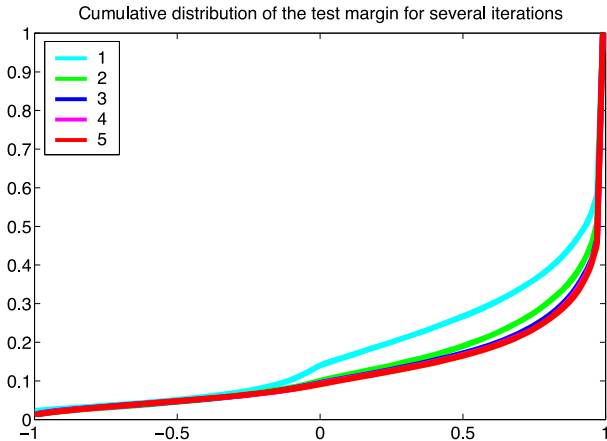
## Boosting

AdaBoost

Face Detection

## Gradient Boosting

- The AdaBoost margin is defined as the distribution of  $y \cdot h(x)$

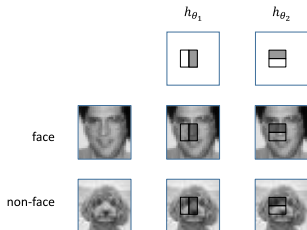




# Face Detection

Face detection framework was proposed in 2001 by Paul Viola and Michael Jones using AdaBoost

- Some hypotheses  $h_\theta$



- Haar-like features for each hypothesis

$$h_\theta = \sum_{(x,y) \in \text{dark area}} \text{image}(x,y) - \sum_{(x,y) \in \text{white area}} \text{image}(x,y) \quad (26)$$





# Gradient Boosting



# What is Gradient Boosting

---

***Gradient Boosting = Gradient Descent + Boosting***

- In Adaboost, “losses” are identified by high-weight data points
- In Gradient Boosting, “losses” are identified by gradients
- Gradient Boosting can be applied for different problems

Regression → Classification → Ranking



# Algorithm

## Bagging

Bootstrap

Algorithm

Random Forests

## Boosting

AdaBoost

Face Detection

## Gradient Boosting

**Input:** training set  $\{(x_i, y_i)\}_{i=1}^n$  and a differentiable loss function  $\ell(y, H(x))$ , number of iterations  $M$ .

1. Initialize model with a constant value:

$$H_0(x) = \gamma_0 = \arg \min_{\gamma} \sum_{i=1}^n \ell(y_i, \gamma). \quad (27)$$

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

2.1 Compute so-called pseudo-residuals:

$$r_{im} = - \left[ \frac{\partial \ell(y_i, H(x_i))}{\partial H(x_i)} \right]_{H(x)=H_{m-1}(x)}, \text{ for } i = 1, \dots, m \quad (28)$$



# Algorithm (cont.)

**2.2** Fit a base learner (or weak learner, e.g. tree) closed under scaling  $h_m(x)$  to pseudo-residuals, i.e. train it using the training set  $\{(x_i, r_{im})\}_{i=1}^n$ .

**2.3** Compute multiplier by solving the following one-dimensional optimization problem:

$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n \ell(y_i, H_{m-1}(x_i) + \gamma h_m(x_i)). \quad (29)$$

**2.4** Update the model:

$$H_m(x) = H_{m-1}(x) + \gamma_m h_m(x) \quad (30)$$

**3.** Output  $H_M(x)$



# Important points to remember

---

- Bagging is predominantly a variance-reduction technique, while boosting is primarily a bias-reduction technique.
- This explains why bagging is often used in combination with high-variance models such as tree models, whereas boosting is typically used with high-bias models such as linear classifiers.

# References

---



Goodfellow, I., Bengio, Y., and Courville, A. (2016).

*Deep learning.*

MIT press.



Lê, B. and Tô, V. (2014).

*Cở sở trí tuệ nhân tạo.*

Nhà xuất bản Khoa học và Kỹ thuật.



Russell, S. and Norvig, P. (2021).

*Artificial intelligence: a modern approach.*

Pearson Education Limited.