

# CSE 252C: Computer Vision III

Lecturer: Serge Belongie

Scribes: Andrew Rabinovich and Vincent Rabaud

Edited by: Catherine Wah

## LECTURE 16

### Boosting

#### 16.1. Introduction

Since Viola & Jones (2001), boosting – and AdaBoost (Freund & Schapire, 1997) – has taken on a prominent role in computer vision, in particular for object detection. The idea behind boosting is to combine the outputs of many “weak learners” in a weighted fashion to produce a “strong” classifier. We’ll focus on the 2-class case in this lecture; a generalization to multiclass exists.

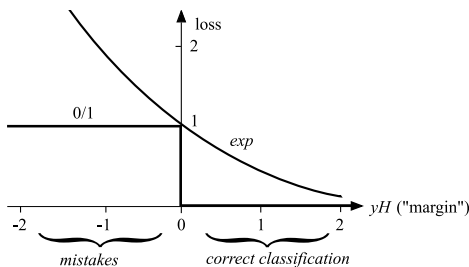
**Definition 16.1.** A *weak learner* is a classifier that is correct at least a bit more than half of the time, *i.e.* it is slightly better than chance.

Formally, suppose we are given a training set:

$$D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\} \text{ with } \mathbf{x}_i \in \mathbb{R}^D, y = \{-1, 1\}.$$

---

<sup>1</sup>Department of Computer Science and Engineering, University of California, San Diego.



**Figure 1.** Exponential loss is the upper bound on the 0/1 loss.

A weak learner is a function  $h_k : \mathbf{x} \mapsto y$ , and the strong classifier has the form:

$$(16.2) \quad H(\mathbf{x}) = \sum_{k=1}^T \alpha_k h_k(\mathbf{x}).$$

The problem is how to set the weights  $\alpha_k$  and train the weak learners  $h_k(\mathbf{x})$ . We'll examine now how this is done with AdaBoost, specifically, discrete AdaBoost (Adaptive Boosting).

## 16.2. AdaBoost

Ideally, one would like to minimize what's known as the 0/1 loss, *i.e.* a penalty of 1 for incorrect classification, and no penalty for correct classification. A 0/1 loss is not differentiable, however, which makes it difficult to work with. AdaBoost uses the so-called exponential loss (see Figure 1):

$$(16.3) \quad \mathcal{L}(H) = \sum_{i=1}^n e^{-y_i H(\mathbf{x}_i)}.$$

Observe that  $y_i \in \{-1, 1\}$ , and  $H(\mathbf{x}_i)$  is multiplied by this; thus, we want it to be the same sign as  $y_i$ .

Now, suppose we are given the first  $t$  weak classifiers, and we want to specify  $h_{t+1}$  and  $\alpha_{t+1}$ . Let  $H_t(\mathbf{x}) = \sum_{k=1}^t \alpha_k h_k(\mathbf{x})$  denote the strong classifier, made up of the first  $t$  weak learners. Iteration  $t + 1$  can then be posed as follows:

$$(16.4) \quad (h_{t+1}, \alpha_{t+1}) = \arg \min_{h, \alpha} \mathcal{L}(H_t + \alpha h)$$

$$(16.5) \quad = \arg \min_{h, \alpha} \sum_{i=1}^n e^{-y_i (H_t(\mathbf{x}_i) + \alpha h(\mathbf{x}_i))}$$

$$(16.6) \quad = \arg \min_{h, \alpha} \sum_{i=1}^n w_t(i) e^{-y_i \alpha h(\mathbf{x}_i)}$$

where  $w_t(i) = \exp(-y_i H_t(\mathbf{x}_i))$  doesn't depend on the min arguments; we can think of it as a per-point weighting at iteration  $t$ .

Furthermore, we can assume the weights have been normalized to sum to 1, w.l.o.g.:  $\sum_{i=1}^n w_t(i) = 1$ . (The arg min is not affected.) Now, let's do some more manipulations:

$$(16.7) \quad (h_{t+1}, \alpha_{t+1}) = \arg \min_{h, \alpha} \sum_{i=1}^n w_t(i) \exp(-y_i \alpha h(\mathbf{x}_i))$$

and break this into two cases:

$$\begin{aligned} (h_{t+1}, \alpha_{t+1}) &= \arg \min_{h, \alpha} e^{-\alpha} \sum_{i|y_i=h(\mathbf{x}_i)} w_t(i) + e^{\alpha} \sum_{i|y_i \neq h(\mathbf{x}_i)} w_t(i) \\ &= \arg \min_{h, \alpha} (e^{\alpha} - e^{-\alpha}) \sum_{i=1}^n w_t(i) \mathbf{1}[y_i \neq h(\mathbf{x}_i)] + e^{-\alpha} \underbrace{\sum_{i=1}^n w_t(i)}_1 \\ &= \arg \min_{h, \alpha} \underbrace{(e^{\alpha} - e^{-\alpha})}_{\text{pos. constant}} \sum_{i=1}^n w_t(i) \mathbf{1}[y_i \neq h(\mathbf{x}_i)] + \underbrace{e^{-\alpha}}_{\text{pos. constant}}. \end{aligned}$$

Therefore:

$$(16.8) \quad h_{t+1} = \arg \min_h \sum_{i=1}^n w_t(i) \mathbf{1}[y_i \neq h(\mathbf{x}_i)].$$

Note how the errors are weighted differently per sample.

How about  $\alpha_{t+1}$ ? We define

$$(16.9) \quad \epsilon_t = \sum_{i=1}^n w_t(i) \mathbf{1}[y_i \neq h(\mathbf{x}_i)],$$

one weighted error of the new weak classifier. We need to solve  $\alpha_{t+1} = \arg \min_{\alpha} (e^{\alpha} - e^{-\alpha}) \epsilon_t + e^{-\alpha}$ . Taking the derivative and setting it to zero we get:

$$\begin{aligned} \frac{\partial}{\partial \alpha} ((e^{\alpha} - e^{-\alpha}) \epsilon_t + e^{-\alpha}) &= e^{\alpha} \epsilon_t + e^{-\alpha} \epsilon_t - e^{-\alpha} = 0 \\ e^{2\alpha} \epsilon_t + \epsilon_t - 1 &= 0 \\ e^{2\alpha} &= \frac{1 - \epsilon_t}{\epsilon_t}. \end{aligned}$$

Hence:

$$(16.10) \quad \alpha_{t+1} = \frac{1}{2} \log \frac{1 - \epsilon_t}{\epsilon_t}.$$

Now we're close to being able to state the AdaBoost algorithm, or more accurately, meta-algorithm. Recall we normalized the weights, which we could denote:

$$w_t(i) = \frac{1}{Z_t} \exp(-y_i H_t(\mathbf{x}_i)),$$

where  $Z_t$  is the normalization factor. On the next iteration, we could set the weights as:

$$w_{t+1}(i) = \frac{1}{Z_{t+1}} \exp(-y_i(H_t(\mathbf{x}_i) + \alpha_{t+1}h_{t+1}(\mathbf{x}_i)))$$

and normalize again, but we can do more simplification. Notice the redundancy:

$$\begin{aligned} w_{t+1}(i) &= \frac{1}{Z_{t+1}} \exp(-y_i(H_t(\mathbf{x}_i) + \alpha_{t+1}h_{t+1}(\mathbf{x}_i))) \\ &= w_t(i) \frac{Z_t}{Z_{t+1}} \exp(-y_i \alpha_{t+1} h_{t+1}(\mathbf{x}_i)). \end{aligned}$$

Now we use a simple trick:

$$\begin{aligned} -yh(\mathbf{x}_i) &= 2 \cdot \mathbf{1}[y_i \neq h(\mathbf{x}_i)] - 1 \\ \Rightarrow w_{t+1}(i) &= w_t(i) \frac{Z_t}{Z_{t+1}} \exp(2\alpha_{t+1} \mathbf{1}[y_i \neq h_{t+1}(\mathbf{x}_i)]) e^{-\alpha}. \end{aligned}$$

Note that  $Z_t$  and  $e^{-\alpha}$  are constants, and they don't affect the above arg min, so we can just absorb them into  $Z_{t+1}$ :

$$\begin{aligned} w_{t+1}(i) &= w_t(i) \frac{1}{Z_{t+1}} \exp(2\alpha_{t+1} \mathbf{1}[y_i \neq h_{t+1}(\mathbf{x}_i)]) \\ &= w_t(i) \frac{1}{Z_{t+1}} \exp\left(\log \frac{1-\epsilon_t}{\epsilon_t} \mathbf{1}[y_i \neq h_{t+1}(\mathbf{x}_i)]\right). \end{aligned}$$

Finally, observe that:

$$\begin{aligned} Z_{t+1} &= \sum_{i|y_i \neq h_{t+1}(\mathbf{x}_i)} w_t(i) \frac{1-\epsilon_t}{\epsilon_t} + \sum_{i|y_i = h_{t+1}(\mathbf{x}_i)} w_t(i) \\ &= \epsilon_t \frac{1-\epsilon_t}{\epsilon_t} + 1 - \epsilon_t \\ &= 2(1 - \epsilon_t), \end{aligned}$$

which leads to a simple expression for setting the value of the new weight:

$$(16.11) \quad w_{t+1}(i) = \begin{cases} \frac{w_t(i)}{2(1-\epsilon_t)} & \text{if } y_i = h_{t+1}(\mathbf{x}_i) \\ \frac{w_t(i)}{2\epsilon_t} & \text{otherwise.} \end{cases}$$

### 16.2.1. Discrete AdaBoost

Now we can state the Discrete AdaBoost algorithm (Schapire, 1997) (Algorithm 16.1). The empirical error is upper-bounded by  $\prod_{t=1}^T Z_t$  (exercise for the reader: why is this the case?).

---

#### Algorithm 16.1 Discrete AdaBoost

---

**Input:** Training data  $\mathcal{D}$ , number of iterations  $T$ , initial distribution  $w_0(i)$  (*e.g.* uniform) over data points  
**Output:** A strong classifier  $H(\mathbf{x})$   
**for**  $t = 1, \dots, T$  **do**  
   train a weak classifier:  $h_t = \arg \min_h \sum_{i=1}^n w_{t-1}(i) \mathbf{1}[y_i \neq h(\mathbf{x}_i)]$   
   calculate error:  $\epsilon_t = \sum_{i=1}^n w_{t-1}(i) \mathbf{1}[y_i \neq h_t(\mathbf{x}_i)]$   
   set:  $\alpha_t = \frac{1}{2} \log \frac{1-\epsilon_t}{\epsilon_t}$   
   set:  $w_{t+1}(i) = \begin{cases} \frac{w_t(i)}{2(1-\epsilon_t)} & \text{if } y_i = h_{t+1}(\mathbf{x}_i) \\ \frac{w_t(i)}{2\epsilon_t} & \text{otherwise} \end{cases}$   
**end for**  
**return**  $H(\mathbf{x}) = \sum_{k=1}^T \alpha_k h_k(\mathbf{x})$

---

### 16.2.2. Discussion

Some of the pros of AdaBoost are that it is very simple to implement, it supports feature selection on very large sets of features, and has fairly good generalization. However, it can overfit in presence of noise and outliers, and yields a suboptimal solution for  $\alpha_t$ 's.

Many variants of AdaBoost have been proposed (*e.g.* LogitBoost, GentleBoost, BrownBoost, ...); they are not covered in this class. These address some of the cons listed above.

A notable example of AdaBoost applied to the problem of object detection is Viola-Jones, which combined AdaBoost, cascade, and Haar-like filters.