

Ensemble methods

1 / 30

Learning theory

- ▶ Probability distribution P over $\mathcal{X} \times \{0, 1\}$; let $(X, Y) \sim P$.
- ▶ We get $S := \{(x_i, y_i)\}_{i=1}^n$, an iid sample from P .
- ▶ **Goal:** Fix $\epsilon, \delta \in (0, 1)$. With probability at least $1 - \delta$ (over random choice of S), learn a classifier $\hat{f}: \mathcal{X} \rightarrow \{-1, +1\}$ with low error rate

$$\text{err}(\hat{f}) = P(\hat{f}(X) \neq Y) \leq \epsilon.$$

- ▶ **Basic question:** When is this possible?
 - ▶ Suppose I even promise you that there is a perfect classifier from a particular *function class* \mathcal{F} .
(E.g., \mathcal{F} = linear classifiers or \mathcal{F} = decision trees.)
 - ▶ **Default:** Empirical Risk Minimization (i.e., pick classifier from \mathcal{F} with lowest training error rate), but this might be computationally difficult (e.g., for decision trees).
- ▶ **Another question:** Is it easier to learn just non-trivial classifiers in \mathcal{F} (i.e., better than random guessing)?

2 / 30

Boosting

Boosting: Using a learning algorithm that provides “rough rules-of-thumb” to construct a **very accurate predictor**.

Motivation:

Easy to construct classification rules that are **correct more-often-than-not** (e.g., “If $\geq 5\%$ of the e-mail characters are dollar signs, then it’s spam.”), but seems **hard** to find a single rule that is **almost always correct**.

Basic idea:

Input: training data S

For $t = 1, 2, \dots, T$:

1. Choose subset of examples $S_t \subseteq S$ (or a distribution over S).
2. Use “**weak learning**” algorithm to get classifier: $f_t := \text{WL}(S_t)$.

Return an “ensemble classifier” based on f_1, f_2, \dots, f_T .

3 / 30

Boosting: history

- 1984 Valiant and Kearns ask whether “boosting” is theoretically possible (formalized in the PAC learning model).
- 1989 Schapire creates first boosting algorithm, solving the open problem of Valiant and Kearns.
- 1990 Freund creates an optimal boosting algorithm (Boost-by-majority).
- 1992 Drucker, Schapire, and Simard empirically observe practical limitations of early boosting algorithms.
- 1995 **Freund and Schapire** create **AdaBoost**—a boosting algorithm with practical advantages over early boosting algorithms.

Winner of 2004 ACM Paris Kanellakis Award:

For their “seminal work and distinguished contributions [...] to the development of the theory and practice of boosting, a general and provably effective method of producing arbitrarily accurate prediction rules by combining weak learning rules”; specifically, for AdaBoost, which “can be used to significantly reduce the error of algorithms used in statistical analysis, spam filtering, fraud detection, optical character recognition, and market segmentation, among other applications”.

4 / 30

```
input Training data  $\{(x_i, y_i)\}_{i=1}^n$  from  $\mathcal{X} \times \{-1, +1\}$ .
1: initialize  $D_1(i) := 1/n$  for each  $i = 1, 2, \dots, n$  (a probability distribution).
2: for  $t = 1, 2, \dots, T$  do
3:   Give  $D_t$ -weighted examples to WL; get back  $f_t: \mathcal{X} \rightarrow \{-1, +1\}$ .
4:   Update weights:
```

$$z_t := \sum_{i=1}^n D_t(i) \cdot y_i f_t(x_i) \in [-1, +1]$$

$$\alpha_t := \frac{1}{2} \ln \frac{1+z_t}{1-z_t} \in \mathbb{R} \quad (\text{weight of } f_t)$$

$$D_{t+1}(i) := D_t(i) \exp(-\alpha_t \cdot y_i f_t(x_i)) / Z_t \quad \text{for each } i = 1, 2, \dots, n,$$

where $Z_t > 0$ is normalizer that makes D_{t+1} a probability distribution.

```
5: end for
```

```
6: return Final classifier  $\hat{f}(x) := \text{sign}\left(\sum_{t=1}^T \alpha_t \cdot f_t(x)\right)$ .
```

(Let $\text{sign}(z) := 1$ if $z > 0$ and $\text{sign}(z) := -1$ if $z \leq 0$.)

Interpreting z_t

Suppose $(X, Y) \sim D_t$. If

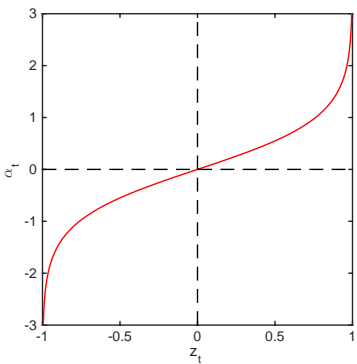
$$P(f(X) = Y) = \frac{1}{2} + \gamma_t,$$

then

$$z_t = \sum_{i=1}^n D_t(i) \cdot y_i f_t(x_i) = 2\gamma_t \in [-1, +1].$$

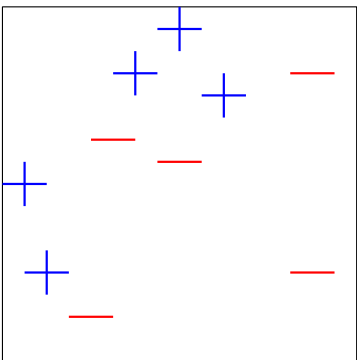
- ▶ $z_t = 0 \iff$ random guessing w.r.t. D_t .
- ▶ $z_t > 0 \iff$ better than random guessing w.r.t. D_t .
- ▶ $z_t < 0 \iff$ better off using the opposite of f 's predictions.

Classifier weights $\alpha_t = \frac{1}{2} \ln \frac{1+z_t}{1-z_t}$



Example weights $D_{t+1}(i)$

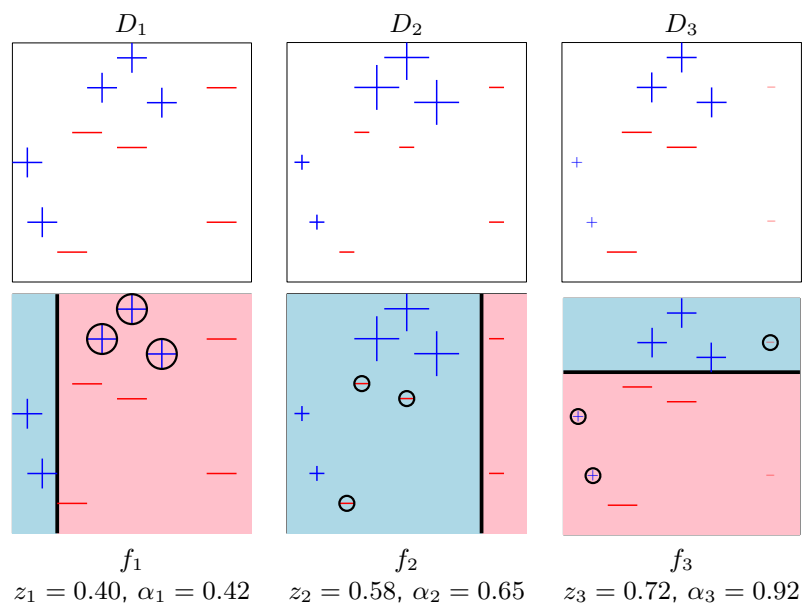
$$D_{t+1}(i) \propto D_t(i) \cdot \exp(-\alpha_t \cdot y_i f_t(x_i)).$$



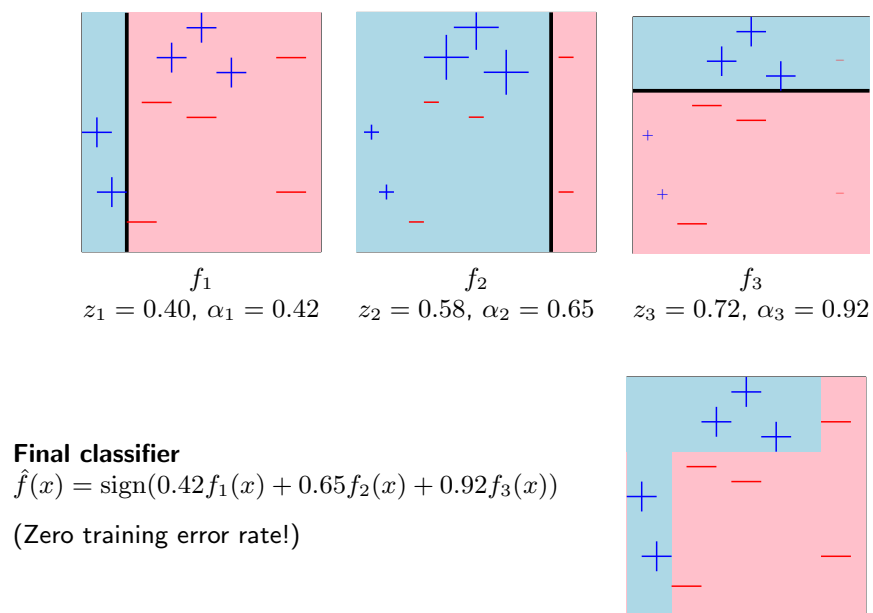
Weak learning algorithm WL: ERM with $\mathcal{F} =$ “decision stumps” on \mathbb{R}^2 (i.e., axis-aligned threshold functions $x \mapsto \text{sign}(vx_i - t)$).
Straightforward to handle importance weights in ERM.

(Example from Figures 1.1 and 1.2 of Schapire & Freund text.)

Example: execution of AdaBoost

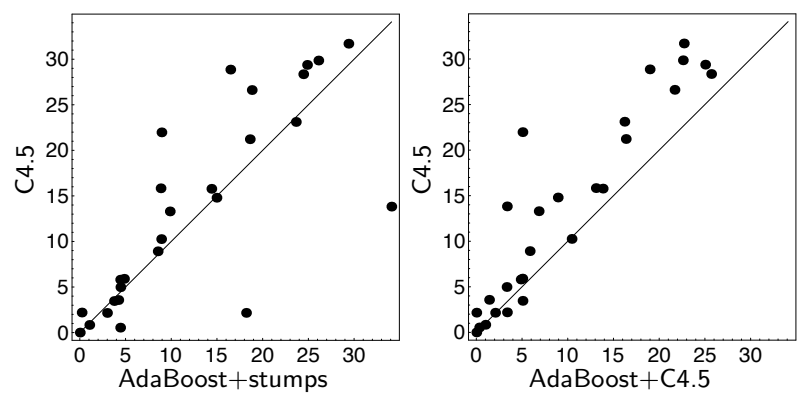


Example: final classifier from AdaBoost



Empirical results

Test error rates of C4.5 and AdaBoost on several classification problems.
Each point represents a single classification problem/dataset from UCI repository.



C4.5 = popular algorithm for learning decision trees.

(Figure 1.3 from Schapire & Freund text.)

Training error rate of final classifier

Recall $\gamma_t := P(f_t(X) = Y) - 1/2 = z_t/2$ when $(X, Y) \sim D_t$.

Training error rate of final classifier from AdaBoost:

$$\text{err}(\hat{f}, \{(x_i, y_i)\}_{i=1}^n) \leq \exp\left(-2 \sum_{t=1}^T \gamma_t^2\right).$$

If average $\bar{\gamma}^2 := \frac{1}{T} \sum_{t=1}^T \gamma_t^2 > 0$, then training error rate is $\leq \exp(-2\bar{\gamma}^2 T)$.

“AdaBoost” = “Adaptive Boosting”

Some γ_t could be small, even negative—only care about overall average $\bar{\gamma}^2$.

What about true error rate?

Combining classifiers

Let \mathcal{F} be the function class used by the **weak learning algorithm** WL.

The function class used by AdaBoost is

$$\mathcal{F}_T := \left\{ x \mapsto \text{sign} \left(\sum_{t=1}^T \alpha_t f_t(x) \right) : f_1, f_2, \dots, f_T \in \mathcal{F}, \alpha_1, \alpha_2, \dots, \alpha_T \in \mathbb{R} \right\}$$

i.e., linear combinations of T functions from \mathcal{F} .

Complexity of \mathcal{F}_T grows **linearly with T** .

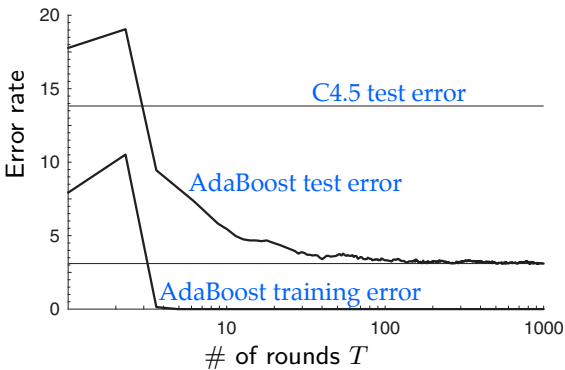
Theoretical guarantee (e.g., when \mathcal{F} = decision stumps in \mathbb{R}^d):
With high probability (over random choice of training sample),

$$\text{err}(\hat{f}) \leq \underbrace{\exp(-2\gamma^2 T)}_{\text{Training error rate}} + \underbrace{O\left(\sqrt{\frac{T \log d}{n}}\right)}_{\text{Error due to finite sample}}.$$

Theory suggests danger of overfitting when T is very large.
Indeed, this does happen sometimes ... **but often not!**

A typical run of boosting

AdaBoost+C4.5 on “letters” dataset.



(# nodes across all decision trees in \hat{f} is $> 2 \times 10^6$)

Training error rate is zero after just five rounds,
but test error rate continues to decrease, even up to 1000 rounds!

(Figure 1.7 from Schapire & Freund text)

Boosting the margin

Final classifier from AdaBoost:

$$\hat{f}(x) = \text{sign} \left(\underbrace{\frac{\sum_{t=1}^T \alpha_t f_t(x)}{\sum_{t=1}^T |\alpha_t|}}_{g(x) \in [-1, +1]} \right).$$

Call $y \cdot g(x) \in [-1, +1]$ the **margin** achieved on example (x, y) .

New theory [Schapire, Freund, Bartlett, and Lee, 1998]:

- ▶ **Larger margins** \Rightarrow **better resistance to overfitting**, independent of T .
- ▶ AdaBoost **tends to increase margins** on training examples.

(Similar but not the same as SVM margins.)

On “letters” dataset:

	$T = 5$	$T = 100$	$T = 1000$
training error rate	0.0%	0.0%	0.0%
test error rate	8.4%	3.3%	3.1%
% margins ≤ 0.5	7.7%	0.0%	0.0%
min. margin	0.14	0.52	0.55

Linear classifiers

Regard function class \mathcal{F} used by weak learning algorithm as “feature functions”:

$$x \mapsto \phi(x) := (f(x) : f \in \mathcal{F}) \in \{-1, +1\}^{\mathcal{F}}$$

(possibly infinite dimensional!).

AdaBoost’s final classifier is a *linear classifier* in $\{-1, +1\}^{\mathcal{F}}$:

$$\hat{f}(x) = \text{sign} \left(\sum_{t=1}^T \alpha_t f_t(x) \right) = \text{sign} \left(\sum_{f \in \mathcal{F}} w_f f(x) \right) = \text{sign}(\langle w, \phi(x) \rangle)$$

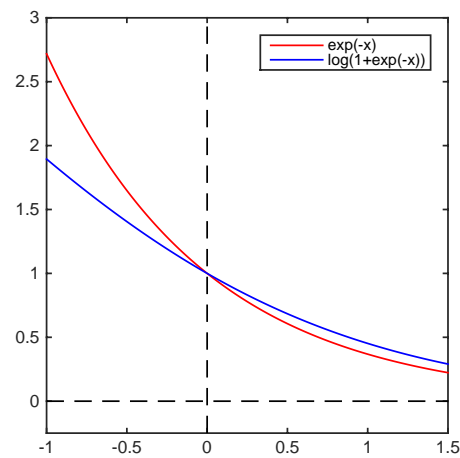
where

$$w_f := \sum_{t=1}^T \alpha_t \cdot \mathbb{1}\{f_t = f\} \quad \forall f \in \mathcal{F}.$$

Exponential loss

AdaBoost is a particular “coordinate descent” algorithm (similar to but not the same as gradient descent) for

$$\min_{\mathbf{w} \in \mathbb{R}^{\mathcal{F}}} \frac{1}{n} \sum_{i=1}^n \exp(-y_i \langle \mathbf{w}, \phi(x_i) \rangle).$$



17 / 30

More on boosting

Many variants of boosting:

- ▶ AdaBoost with different loss functions.
- ▶ Boosted decision trees = boosting + decision trees.
- ▶ Boosting algorithms for *ranking* and *multi-class*.
- ▶ Boosting algorithms that are robust to certain kinds of noise.
- ▶ Boosting for online learning algorithms (very new!).
- ▶ ...

Many connections between boosting and other subjects:

- ▶ Game theory, online learning
- ▶ “Geometry” of information (replace $\|\cdot\|_2^2$ with relative entropy divergence)
- ▶ Computational complexity
- ▶ ...

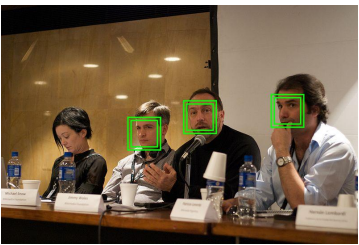
18 / 30

An application of AdaBoost

Application: face detection

Face detection

Problem: Given an image, locate all of the faces.



As a classification problem:

- ▶ Divide up images into patches (at varying scales, e.g., 24×24 , 48×48).
- ▶ Learn classifier $f: \mathcal{X} \rightarrow \mathcal{Y}$, where $\mathcal{Y} = \{\text{face}, \text{not face}\}$.

Many other things built on top of face detectors (e.g., face tracking, face recognizers); now in every digital camera and iPhoto/Picasa-like software.

Main problem: how to make this **very fast**.

20 / 30

Face detectors via AdaBoost [Viola & Jones, 2001]

Face detector architecture by Viola & Jones (2001): major achievement in computer vision; **detector actually usable in real-time.**

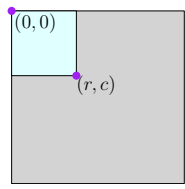
- ▶ Think of each image patch ($d \times d$ -pixel gray-scale) as a vector $x \in [0, 1]^{d^2}$.
- ▶ Used weak learning algorithm that picks linear classifiers $f_{w,t}(x) = \text{sign}(\langle w, x \rangle - t)$, where w has a very particular form:



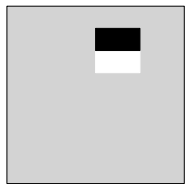
- ▶ AdaBoost combines many “rules-of-thumb” of this form.
 - ▶ Very simple.
 - ▶ **Extremely fast to evaluate** via pre-computation.

Viola & Jones “integral image” trick

“Integral image” trick:



For every image, pre-compute $s(r, c) = \text{sum of pixel values in rectangle from } (0, 0) \text{ to } (r, c)$ (single pass through image).



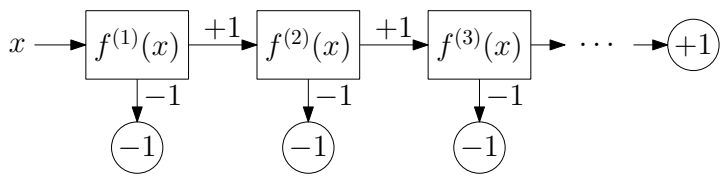
To compute inner product $\langle w, x \rangle = \text{average pixel value in black box} - \text{average pixel value in white box}$ just need to add and subtract a few $s(r, c)$ values.

⇒ Evaluating “rules-of-thumb” classifiers is **extremely fast**.

Viola & Jones cascade architecture

Problem: severe class imbalance (most patches don’t contain a face).

Solution: Train several classifiers (each using AdaBoost), and arrange in a special kind of **decision list** called a **cascade**:



- ▶ Each $f^{(\ell)}$ is trained (using AdaBoost), adjust threshold (before passing through sign) to *minimize false negative rate*.
- ▶ Can make $f^{(\ell)}$ in later stages more complex than in earlier stages, since most examples don’t make it to the end.

⇒ (Cascade) classifier evaluation **extremely fast**.

Viola & Jones detector: example results



Two key points:

- ▶ AdaBoost effectively combines many fast-to-evaluate “weak classifiers”.
- ▶ Cascade structure optimizes speed for common case.

Bagging

Bagging = **B**ootstrap **a**ggregating (Leo Breiman, 1994).

Input: training data $\{(x_i, y_i)\}_{i=1}^n$ from $\mathcal{X} \times \{-1, +1\}$.

For $t = 1, 2, \dots, T$:

1. Randomly pick n examples *with replacement* from training data
→ $\{(x_i^{(t)}, y_i^{(t)})\}_{i=1}^n$ (a *bootstrap* sample).
2. Run learning algorithm on $\{(x_i^{(t)}, y_i^{(t)})\}_{i=1}^n$
→ classifier f_t .

Return a majority vote classifier over f_1, f_2, \dots, f_T .

Question: if n individuals are picked from a population of size n *u.a.r. with replacement*, what is the probability that a given individual is *not* picked?

Answer:

$$\left(1 - \frac{1}{n}\right)^n$$

For large n :

$$\lim_{n \rightarrow \infty} \left(1 - \frac{1}{n}\right)^n = \frac{1}{e} \approx 0.3679.$$

Implications for bagging:

- ▶ Each bootstrap sample contains about 63% of the data set.
- ▶ Remaining 37% can be used to estimate error rate of classifier trained on the bootstrap sample.
- ▶ Can average across bootstrap samples to get estimate of *bagged classifier's* error rate (sort of).

Random Forests (Leo Breiman, 2001).

Input: training data $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ from $\mathbb{R}^d \times \{-1, +1\}$.

For $t = 1, 2, \dots, T$:

1. Randomly pick n examples *with replacement* from training data
→ $\{(\mathbf{x}_i^{(t)}, y_i^{(t)})\}_{i=1}^n$ (a *bootstrap* sample).
2. Run variant of decision tree learning algorithm on $\{(\mathbf{x}_i^{(t)}, y_i^{(t)})\}_{i=1}^n$,
where each split is chosen by only considering a random subset of \sqrt{d} features (rather than all d features)
→ decision tree classifier f_t .

Return a majority vote classifier over f_1, f_2, \dots, f_T .

1. Theoretical concept of weak and strong learning.
2. AdaBoost algorithm; concept of margins in boosting.
3. Interpreting AdaBoost's final classifier as a linear classifier, and interpreting AdaBoost as a coordinate descent algorithm.
4. Structure of decision lists / cascades.
5. Concept of bootstrap samples; bagging and random forests.