# COMS 4771 Lecture 14

1. Boosting

# BOOSTING

# WHAT IS BOOSTING?

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

# WHAT IS 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 hard to find a single rule that is almost always correct.

# WHAT IS 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 hard to find a single rule that is almost always correct.

**Basic idea**:

    **Input**: training data $S$, "weak" learning algorithm $\mathcal{A}$

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

        1. Choose subset of examples $S_t \subseteq S$ or a distribution over $S$).
        2. Call weak learning algorithm to get classifier: $f_t := \mathcal{A}(S_t)$.

    **Return** a weighted majority vote over $f_1, f_2, \ldots, f_T$.

# Boosting: history

1984 Valiant and Kearns ask whether "boosting" is theoretically possible (formalized in the PAC learning model).

# 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.

# 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).

# 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.

# 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.

# 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".*

# ADABOOST

**Input** Training data $S$ from $\mathcal{X} \times \{\pm 1\}$.

      Weak learning algorithm $\mathcal{A}$ (for importance-weighted classification).

1: **initialize** $D_1(x, y) := 1/|S|$ for each $(x, y) \in S$ (a probability distribution).

2: **for** $t = 1, 2, \ldots, T$ **do**

3:     Give $D_t$-weighted examples to $\mathcal{A}$, get back $f_t : \mathcal{X} \to \{\pm 1\}$. <span style="color:red">Dt is the weight for all samples</span>

4:     Update weights:

<span style="color:red">the bigger the Zt, the accurate the</span>

$$z_t := \sum_{(x,y) \in S} D_t(x, y) \cdot y f_t(x) \in [-1, +1]$$

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

$$D_{t+1}(x, y) \propto D_t(x, y) \exp(-\alpha_t \cdot y f_t(x)) \quad \text{for each } (x, y) \in S.$$

5: **end for**

6: **return** Final classifier $f_{\text{final}}(x) := \text{sign}\left( \sum_{t=1}^{T} \alpha_t \cdot f_t(x) \right)$. <span style="color:red">when f(x) is wrong, the larger the exp(-a y f(x))</span>

the probability for a
single point !

the larger, the better

Interpreting $z_t$

If $\Pr_{(X,Y) \sim D_t}[f_t(X) = Y] = \frac{1}{2} + \gamma_t$ for some $\gamma_t \in [-1/2, +1/2]$,

then $z_t = \sum_{(x,y) \in S} D_t(x, y) \cdot y f_t(x) = 2\gamma_t \in [-1, +1]$.
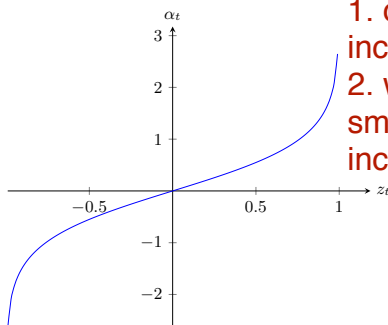
$z_t = 0 \Longleftrightarrow$ random guessing w.r.t. $D_t$.
$z_t > 0 \Longleftrightarrow$ better than random guessing w.r.t. $D_t$.
$z_t < 0 \Longleftrightarrow$ better off using the opposite of $f_t$'s predictions.

Zt is for
classifier

# INTERPRETATION

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



as zt increase:
1. classifier weight increase
2. wrongly labeled smaple's weight increase

Example weights $D_{t+1}(x, y)$

$$D_{t+1}(x, y) \propto D_t(x, y) \cdot \exp(-\alpha_t \cdot y f_t(x))$$
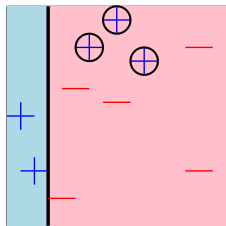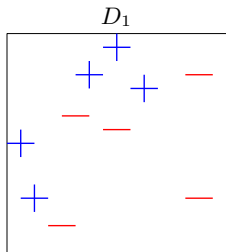
one-level
decision tree

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

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

$D_1$

$f_1$
$z_1 = 0.40,\ \alpha_1 = 0.42$
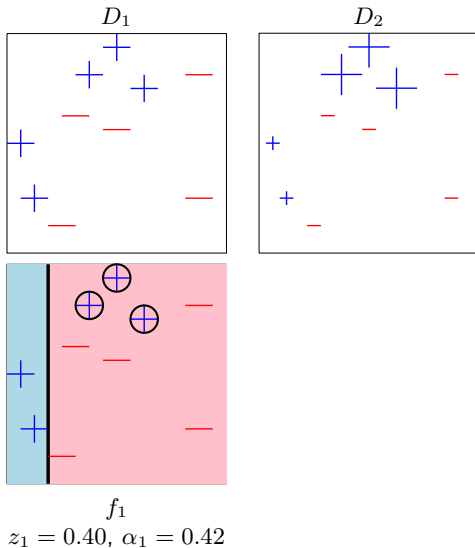
# EXAMPLE: EXECUTION OF ADABOOST



$D_1$

$D_2$

$f_1$
$z_1 = 0.40, \ \alpha_1 = 0.42$

$D_1$

$D_2$

$f_1$
$z_1 = 0.40, \ \alpha_1 = 0.42$

$f_2$
$z_2 = 0.58, \ \alpha_2 = 0.65$

$f_1$
$z_1 = 0.40, \ \alpha_1 = 0.42$

$f_2$
$z_2 = 0.58, \ \alpha_2 = 0.65$

# EXAMPLE: EXECUTION OF ADABOOST



$D_1$      $D_2$      $D_3$

$f_1$
$z_1 = 0.40, \; \alpha_1 = 0.42$

$f_2$
$z_2 = 0.58, \; \alpha_2 = 0.65$

$f_3$
$z_3 = 0.72, \; \alpha_3 = 0.92$

# EXAMPLE: FINAL CLASSIFIER FROM ADABOOST



$f_1$
$z_1 = 0.40, \ \alpha_1 = 0.42$
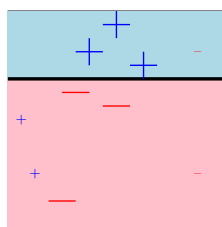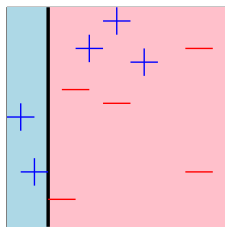
$f_2$
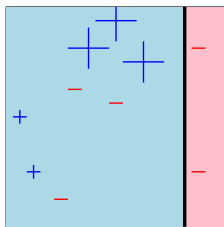$z_2 = 0.58, \ \alpha_2 = 0.65$

$f_3$
$z_3 = 0.72, \ \alpha_3 = 0.92$

# EXAMPLE: FINAL CLASSIFIER FROM ADABOOST



$f_1$
$z_1 = 0.40, \ \alpha_1 = 0.42$

$f_2$
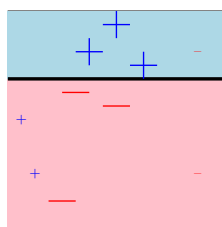$z_2 = 0.58, \ \alpha_2 = 0.65$

$f_3$
$z_3 = 0.72, \ \alpha_3 = 0.92$

**Final classifier**

$f_{\mathsf{final}}(x) = \mathrm{sign}(0.42 f_1(x) + 0.65 f_2(x) + 0.92 f_3(x))$

(Zero training error!)

**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.)

Recall $\gamma_t := \Pr_{(X,Y) \sim D_t}[f_t(X) = Y] - 1/2 = z_t/2$.

**Training error of final classifier from AdaBoost:**

$$\mathrm{err}(f_{\mathsf{final}}, S) \ \leq \ \exp\left(-2\sum_{t=1}^{T} \gamma_t^2\right).$$

Recall $\gamma_t := \Pr_{(X,Y) \sim D_t}[f_t(X) = Y] - 1/2 = z_t/2$.

**Training error of final classifier from AdaBoost**:

$$\mathrm{err}(f_{\mathsf{final}}, S) \ \leq \ \exp\left(-2\sum_{t=1}^{T} \gamma_t^2\right).$$

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

Recall $\gamma_t := \Pr_{(X,Y) \sim D_t}[f_t(X) = Y] - 1/2 = z_t/2$.

**Training error of final classifier from AdaBoost**:

$$\text{err}(f_{\text{final}}, S) \ \leq \ \exp\left(-2\sum_{t=1}^{T} \gamma_t^2\right).$$

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

**"AdaBoost" = "Adaptive Boosting"**
Some $\gamma_t$ could be small (or even negative!)—only care about overall average $\bar{\gamma}^2$.

Recall $\gamma_t := \Pr_{(X,Y) \sim D_t}[f_t(X) = Y] - 1/2 = z_t/2$.

**Training error of final classifier from AdaBoost**:

$$\text{err}(f_{\text{final}}, S) \leq \exp\left(-2\sum_{t=1}^{T} \gamma_t^2\right).$$

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

**"AdaBoost" = "Adaptive Boosting"**
Some $\gamma_t$ could be small (or even negative!)—only care about overall average $\bar{\gamma}^2$.

**What about true error?**

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

# COMBINING CLASSIFIERS

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

The function class used by AdaBoost is

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

i.e., "linear combinations of $T$ functions from $\mathcal{F}$".
Complexity of $\mathcal{F}_T$ grows *linearly* with $T$.

## COMBINING CLASSIFIERS

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

The function class used by AdaBoost is

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

i.e., "linear combinations of $T$ functions from $\mathcal{F}$".
Complexity of $\mathcal{F}_T$ grows *linearly* with $T$.

**Theoretical guarantee**: with high probability over choice of i.i.d. sample $S$,

$$\text{err}(f) \leq \text{err}(f, S) + O\left(\sqrt{\frac{T \log |\mathcal{F}_{|S}|}{|S|}}\right) \quad \forall f \in \mathcal{F}_T.$$

# COMBINING CLASSIFIERS

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

The function class used by AdaBoost is

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

i.e., "linear combinations of $T$ functions from $\mathcal{F}$".
Complexity of $\mathcal{F}_T$ grows *linearly* with $T$.

**Theoretical guarantee**: with high probability over choice of i.i.d. sample $S$,

$$\text{err}(f) \ \leq \ \text{err}(f, S) + O\left(\sqrt{\frac{T \log |\mathcal{F}_{|S|}|}{|S|}}\right) \quad \forall f \in \mathcal{F}_T.$$

**Theory suggests danger of over-fitting when $T$ is very large.**

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

The function class used by AdaBoost is

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

i.e., "linear combinations of $T$ functions from $\mathcal{F}$".
Complexity of $\mathcal{F}_T$ grows *linearly* with $T$.

this part get bigger

**Theoretical guarantee**: with high probability over choice of i.i.d. sample $S$,

$$\text{err}(f) \ \leq \ \text{err}(f, S) + O\left( \sqrt{\frac{T \log |\mathcal{F}_{|S|}}{|S|}} \right) \quad \forall f \in \mathcal{F}_T.$$
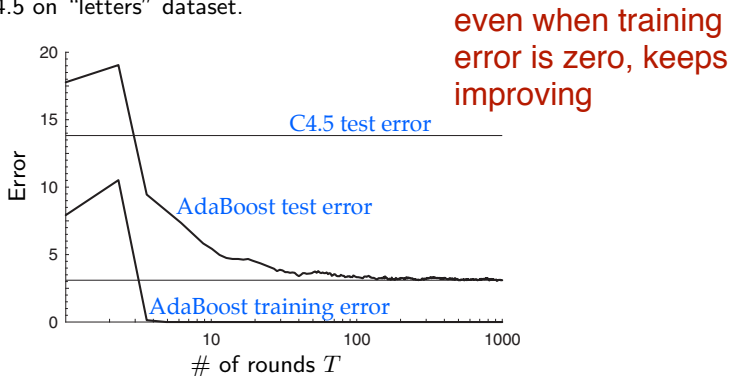
can lead to overfitting

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

As the T increases, the training error goes down!

# A TYPICAL RUN OF BOOSTING

AdaBoost+C4.5 on "letters" dataset.



even when training error is zero, keeps improving

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

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

(Figure 1.7 from Schapire & Freund text)

# Boosting the margin

Final classifier from AdaBoost:

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

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

# Boosting the margin

Final classifier from AdaBoost:

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

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 ⇒ better generalization error**, independent of $T$.
- ▶ AdaBoost tends to increase margins on training examples.

(Similar but not the same as SVM margins.)

Final classifier from AdaBoost:

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

the normlizer has no effect over sign!

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** ⇒ **better generalization error**, independent of $T$.
- AdaBoost tends to increase margins on training examples.

(Similar but not the same as SVM margins.)

actually increase large margin a different kind of margin

**On "letters" dataset:**

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

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

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

(possibly infinite dimensional!).

|F| is the size of feature functions' collection

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

$$x \mapsto \boldsymbol{\phi}(x) := (f(x) : f \in \mathcal{F}) \in \{\pm 1\}^{\mathcal{F}}$$

(possibly infinite dimensional!).

AdaBoost's final classifier is a *linear classifier* in $\{\pm 1\}^{\mathcal{F}}$:

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

where

it's interesting!

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

# Exponential loss

AdaBoost is a particular "coordinate descent" algorithm for

$$\min_{\boldsymbol{w} \in \mathbb{R}^{\mathcal{F}}} \quad \frac{1}{|S|} \sum_{(\boldsymbol{x}, y) \in S} \ell_{\exp}(y \langle \boldsymbol{w}, \boldsymbol{\phi}(\boldsymbol{x}) \rangle)$$
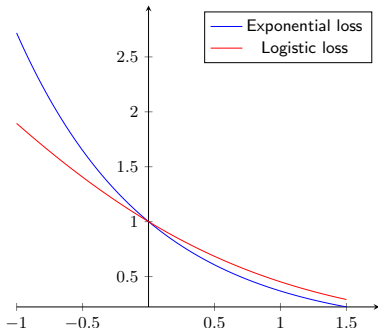
# EXPONENTIAL LOSS

actually an algorithm for solving loss_exp target!!!

AdaBoost is a particular "coordinate descent" algorithm for

great!!!!

$$\min_{\boldsymbol{w} \in \mathbb{R}^{\mathcal{F}}} \quad \frac{1}{|S|} \sum_{(\boldsymbol{x}, y) \in S} \ell_{\exp}(y \langle \boldsymbol{w}, \boldsymbol{\phi}(\boldsymbol{x}) \rangle)$$



**Exponential loss**:

$$\ell_{\exp}(z) = \exp(-z).$$

$\mathbb{E}[\ell_{\exp}(Y \cdot g(x))|X = x]$ is minimized by

$$g(x) = \frac{1}{2} \ln\left(\frac{\eta(x)}{1 - \eta(x)}\right)$$

where $\eta(x) = \Pr[Y = +1|X = x]$.

## Face detection

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

## 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., $24 \times 24$, $48 \times 48$).
- Classify each patch as "face" or "not face".

## 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., $24 \times 24$, $48 \times 48$).
- ▶ Classify each patch as "face" or "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.

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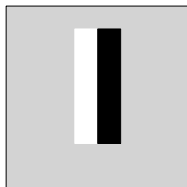
**Face detector architecture by Viola & Jones (2001)**: major achievement in computer vision; **detector actually usable in real-time**.

▶ Think of each image patch ($d \times d$-pixel gray-scale) as a vector $\boldsymbol{x} \in [0, 1]^{d^2}$.

**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 $\boldsymbol{x} \in [0,1]^{d^2}$.

- ▶ Used weak learning algorithm that picks linear classifiers $f_{\boldsymbol{w},\theta}(\boldsymbol{x}) = \text{sign}(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta)$, where $\boldsymbol{w}$ has a very particular form:

**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 $\boldsymbol{x} \in [0,1]^{d^2}$.

▶ Used weak learning algorithm that picks linear classifiers
$f_{\boldsymbol{w},\theta}(\boldsymbol{x}) = \mathrm{sign}(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta)$, where $\boldsymbol{w}$ has a very particular form:



$\langle \boldsymbol{w}, \boldsymbol{x} \rangle =$ average pixel value in black box
$-$ average pixel value in white box

**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 $\boldsymbol{x} \in [0,1]^{d^2}$.

- ▶ Used weak learning algorithm that picks linear classifiers $f_{\boldsymbol{w},\theta}(\boldsymbol{x}) = \mathrm{sign}(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta)$, where $\boldsymbol{w}$ has a very particular form:



$\langle \boldsymbol{w}, \boldsymbol{x} \rangle =$ average pixel value in black box

$-$ average pixel value in white box

**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 $\boldsymbol{x} \in [0,1]^{d^2}$.

▶ Used weak learning algorithm that picks linear classifiers $f_{\boldsymbol{w},\theta}(\boldsymbol{x}) = \mathrm{sign}(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta)$, where $\boldsymbol{w}$ has a very particular form:



$\langle \boldsymbol{w}, \boldsymbol{x} \rangle =$ average pixel value in black box

$-$ average pixel value in white box

**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 $\boldsymbol{x} \in [0,1]^{d^2}$.

▶ Used weak learning algorithm that picks linear classifiers $f_{\boldsymbol{w},\theta}(\boldsymbol{x}) = \mathrm{sign}(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta)$, where $\boldsymbol{w}$ has a very particular form:

$\langle \boldsymbol{w}, \boldsymbol{x} \rangle =$ average pixel value in black box
$- $ average pixel value in white box

**Face detector architecture by Viola & Jones (2001)**: major achievement in computer vision; **detector actually usable in real-time**.
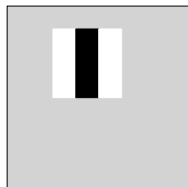
- ▶ Think of each image patch ($d \times d$-pixel gray-scale) as a vector $\boldsymbol{x} \in [0,1]^{d^2}$.

- ▶ Used weak learning algorithm that picks linear classifiers $f_{\boldsymbol{w},\theta}(\boldsymbol{x}) = \text{sign}(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta)$, where $\boldsymbol{w}$ has a very particular form:

**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 $\boldsymbol{x} \in [0,1]^{d^2}$.

▶ Used weak learning algorithm that picks linear classifiers
$f_{\boldsymbol{w},\theta}(\boldsymbol{x}) = \mathrm{sign}(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \theta)$, where $\boldsymbol{w}$ has a very particular form:
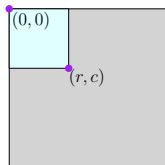


note the two different patches

▶ **Many possible "rules-of-thumb" of this form**.
AdaBoost combines several of them to build an accurate classifier.

# Viola & Jones "integral image" trick
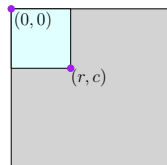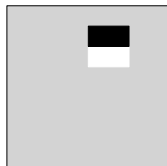
**"Integral image" trick**:



For every image, pre-compute

$s(r, c) = $ sum of pixel values in rectangle from $(0, 0)$ to $(r, c)$

(single pass through image).

**"Integral image" trick**:



For every image, pre-compute

$s(r, c)$ = sum of pixel values in rectangle from $(0,0)$ to $(r, c)$

(single pass through image).



To compute inner product

$\langle \boldsymbol{w}, \boldsymbol{x} \rangle$ = average pixel value in black box
$-$ average pixel value in white box

just need to add and subtract a few $s(r, c)$ values.

**"Integral image" trick**:



For every image, pre-compute

$s(r, c) =$ sum of pixel values in rectangle from $(0, 0)$ to $(r, c)$
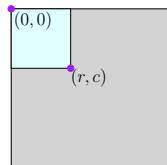
(single pass through image).



To compute inner product

$\langle \boldsymbol{w}, \boldsymbol{x} \rangle =$ average pixel value in black box

$-$ average pixel value in white box

just need to add and subtract a few $s(r, c)$ values.

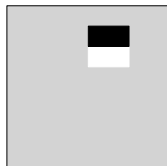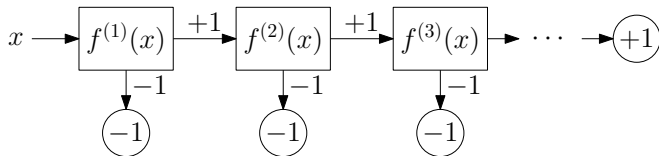$\Rightarrow$ Evaluating "rules-of-thumb" classifiers is **extremely fast**.

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

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

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

$$x \longrightarrow \boxed{f^{(1)}(x)} \xrightarrow{+1} \boxed{f^{(2)}(x)} \xrightarrow{+1} \boxed{f^{(3)}(x)} \longrightarrow \cdots \longrightarrow \bigcirc{+1}$$

with each stage branching $-1$ down to $\bigcirc{-1}$.

- ▶ Each $f^{(\ell)}$ is trained (using AdaBoost), adjust threshold (before passing through $\mathrm{sign}$) 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.

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

**Solution**: Train several classifiers (using AdaBoost), and arrange in a special kind of **decision list** called a **cascade**:
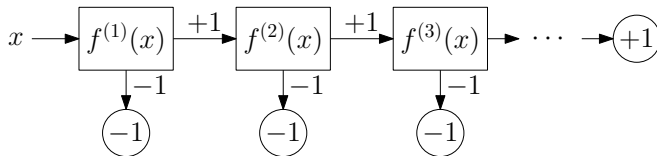
$$x \longrightarrow \boxed{f^{(1)}(x)} \xrightarrow{+1} \boxed{f^{(2)}(x)} \xrightarrow{+1} \boxed{f^{(3)}(x)} \longrightarrow \cdots \longrightarrow \bigcirc\!\!\!\!\!+1$$

with each stage having a $-1$ branch leading to $\bigcirc\!\!\!\!-1$.

- ▶ Each $f^{(\ell)}$ is trained (using AdaBoost), adjust threshold (before passing through $\mathrm{sign}$) 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**.

# More on boosting

**Many variants of boosting**:

- AdaBoost.L and LogitBoost (replaces $\ell_{\exp}$ with $\ell_{\log}$).
- Forward-{step,stage}wise regression (replaces $\ell_{\exp}$ with $\ell_{\mathrm{sq}}$).
- Boosted decision trees $=$ boosting $+$ decision trees, often with $\ell_{\mathrm{sq}}$. (See ESL Chapter 10.)
- Boosting algorithms for *ranking* and *multi-class*.
- Boosting algorithms that are robust to certain kinds of noise.
- . . .

**Many connections between boosting and other subjects**:

- Game theory, online learning
- Information geometry
- Computational complexity
- . . .