

# WaldBoost: Learning for Sequential Classification

Jan Šochman

joint work with Jiří Matas

Center for Machine Perception  
Czech Technical University, Prague

<http://cmp.felk.cvut.cz>



## We know that...

- Time-to-decision vs. precision trade-off is inherent in many computer vision problems
- Decision time clearly influences impact of a method:
  - Viola-Jones (2001) – real-time performance [2500 citations](#)
  - Schneiderman-Kanade (1998) – smaller error rates, but 1000x slower [250 citations](#)

## But...

- Machine learning focuses mainly on error minimisation
- Time-to-decision vs. precision trade-off not explicitly stated in the problem formulation
- Speedup usually gained heuristically by code optimisation, hardware implementation, new architecture proposals, ...
- A sequential decision theory is well established in statistics, yet it is not formulated as a *learning task*

## Talk outline

- Learning problem formulation
- WaldBoost algorithm
- Applications
  - Fast face detection
  - Interest point detector speedup by emulation

## Basic notions

$$\mathbf{x} \in \mathcal{X}$$

classified object (e.g. image)

$$y \in \mathcal{Y} = \{-1, +1\}$$

object class

$$(\xi_1, \dots, \xi_M) \text{ where } \xi_t: \mathcal{X} \rightarrow \mathcal{X}_t$$

vector of measurement functions

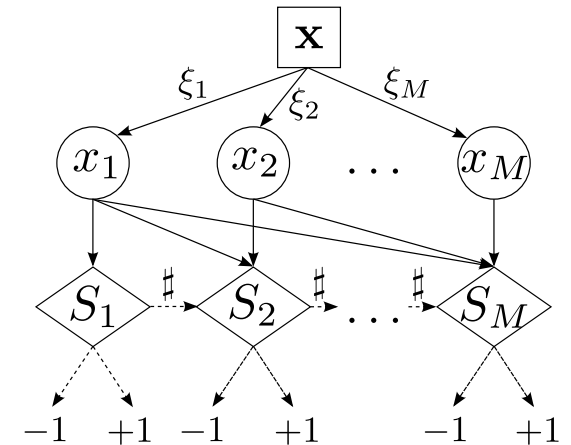
$$(x_1, \dots, x_M) \text{ where } x_t \in \mathcal{X}_t$$

ordered sequence of measurements

$$S = (S_1, \dots, S_M)$$

sequential strategy

$$\text{where } S_t: \mathcal{X}_1 \times \dots \times \mathcal{X}_t \rightarrow \mathcal{Y} \cup \{\#\}$$



## Any sequential strategy $S$ is associated with

- $\alpha_S = P(S = -1 | y = +1), \quad \beta_S = P(S = +1 | y = -1)$

error of the first and the second kind

- $\bar{T}_S^+ = E_{\mathcal{X}^+}[T_S(\mathbf{x})] \quad \bar{T}_S^- = E_{\mathcal{X}^-}[T_S(\mathbf{x})]$

average evaluation time for each class

where  $T_S(\mathbf{x}) = \arg \min_t (S_t(\xi_1(\mathbf{x}), \dots, \xi_t(\mathbf{x})) \neq \#)$

time-to-decision

## Classification task

For given bounds  $\alpha$  and  $\beta$ , find a strategy  $S^*$  such that

$$\alpha_{S^*} \leq \alpha \text{ and } \beta_{S^*} \leq \beta$$

and for any strategy  $S$

$$\bar{T}_{S^*}^+ \leq \bar{T}_S^+ \text{ and } \bar{T}_{S^*}^- \leq \bar{T}_S^- \quad \text{s.t. } \alpha_S \leq \alpha \text{ and } \beta_S \leq \beta$$

## Sequential Probability Ratio Test (SPRT) [Wald 1947]

SPRT is a sequential strategy  $S^*$

$$S_t^* = \begin{cases} +1, & R_t \leq B \\ -1, & R_t \geq A \\ \#, & B < R_t < A \end{cases} \quad R_t = \frac{p(x_1, \dots, x_t | y = -1)}{p(x_1, \dots, x_t | y = +1)}.$$

Practical (and nearly optimal) setting for the thresholds  $A$  and  $B$

$$A = \frac{1 - \beta}{\alpha}, \quad B = \frac{\beta}{1 - \alpha}$$

## Difficulties

SPRT is an **optimal sequential test** in sense of our optimisation problem **but**:

1. the ordering of the measurements is expected to be given
2. multidimensional pdf's  $p(x_1, \dots, x_t | y = \text{class})$  has to be computed in non i.i.d case

## Given

- a training set  $\mathcal{T} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$
- an unordered set (class) of measurement functions  
 $\Xi_u = \{\xi_1, \xi_2, \dots, \xi_M\}$
- two values  $\alpha, \beta$  such that  $0 \leq \alpha, \beta \leq 1$

$$\mathbf{x}_i \in \mathcal{X}, y_i \in \{-1, +1\}$$

$$\xi_t: \mathcal{X} \rightarrow \mathcal{X}_t$$

## Find

- an ordering  $\pi$  of  $\Xi_u$
- a sequential strategy  $S^L$  using  $(\xi_{\pi(1)}, \dots, \xi_{\pi(M)})$

such that

$$\alpha_{S^L} \leq \alpha \text{ and } \beta_{S^L} \leq \beta$$

and for any sequential strategy  $S$  using  $\Xi_u$

$$\bar{T}_{S^L}^+ \leq \bar{T}_S^+ \text{ and } \bar{T}_{S^L}^- \leq \bar{T}_S^- \quad \text{s.t. } \alpha_S \leq \alpha \text{ and } \beta_S \leq \beta$$

---

## How to keep the optimality of SPRT?

- **ordering:** What is a good measurement ordering?
- **decision thresholds:** How to avoid multi-dimensional density estimation?

## Why to select and order the measurements by AdaBoost

- AdaBoost can be seen as a feature (measurement) selector with a principled strategy (minimisation of upper bound on empirical error)
- Is close to sequential decision making, as it produces a sequence of gradually more complex classifiers
- Its output converges to the logarithm of likelihood ratio

## Real AdaBoost [Schapire & Singer, ML 1999]

- Selects iteratively weak classifiers  $h_t: \mathcal{X} \rightarrow \mathbb{R}$  and combines their outputs into a strong classifier

$$f_T(\mathbf{x}) = \sum_{t=1}^T h_t(\mathbf{x})$$

- Classification is given by  $\text{sign}(f_T(\mathbf{x}))$

## Measurements

Each weak classifier is associated with one measurement  $\xi_t \equiv h_t$

$$\mathcal{X}_t = \mathbb{R}$$

$\Rightarrow h_1, h_2, \dots$  **ordered sequence of measurements**

## Ratio needed for SPRT

$$R_t(\mathbf{x}) = \frac{p(h_1(\mathbf{x}), \dots, h_t(\mathbf{x}) | y = -1)}{p(h_1(\mathbf{x}), \dots, h_t(\mathbf{x}) | y = +1)}$$

- Weak classifiers not independent  $\rightarrow$  multi-dimensional  $\rightarrow$  intractable
- We would like to have a test working with  $f_t(\mathbf{x})$  directly

## Asymptotic convergence of AdaBoost [Friedman, TR 1998]

$$\lim_{t \rightarrow \infty} f_t(\mathbf{x}) = -\frac{1}{2} \log R(\mathbf{x}) + \frac{1}{2} \log \frac{P(y = +1)}{P(y = -1)} \qquad R(x) = \frac{p(\mathbf{x}|y=-1)}{p(\mathbf{x}|y=+1)}$$

- AdaBoost response  $f_t(\mathbf{x})$  is (asymptotically) a monotonic function of  $R(\mathbf{x})$
- Thresholding  $R_t(\mathbf{x})$  on  $A \Leftrightarrow$  thresholding  $f_t(\mathbf{x})$  on some  $\theta_t^A$  (similarly for  $B$  and  $\theta_t^B$ )  
for sufficiently large  $t$
- Knowing the thresholds  $\theta_t^A$  and  $\theta_t^B$ , the sequential decision strategy would become

$$S_t = \begin{cases} +1, & f_t(\mathbf{x}) \geq \theta_t^B \\ -1, & f_t(\mathbf{x}) \leq \theta_t^A \\ \#, & \theta_t^A < f_t(\mathbf{x}) < \theta_t^B \end{cases} \qquad \begin{array}{l} \text{inequalities inverted} \\ \text{because } f(\mathbf{x}) \approx -R(x) \end{array}$$

## Search procedure (for $\theta_t^A$ only)

- We want to find  $\theta_t^A$  corresponding to the condition

$$R_t(\mathbf{x}) = \frac{p(h_1(\mathbf{x}), \dots, h_t(\mathbf{x})|y = -1)}{p(h_1(\mathbf{x}), \dots, h_t(\mathbf{x})|y = +1)} \geq A$$

classification to -1 class

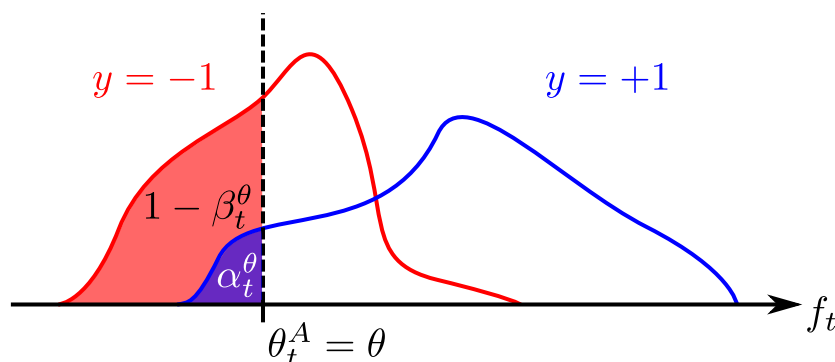
- This can be rewritten as

$$\underbrace{\sum_{\mathbf{x} \in D_t^-} p(h_1, \dots, h_t|y = -1)}_{1 - \beta_t \text{ (true negatives)}} \geq A \underbrace{\sum_{\mathbf{x} \in D_t^-} p(h_1, \dots, h_t|y = +1)}_{\alpha_t \text{ (false negatives)}}$$

$$D_t^- = \{\mathbf{x}; R_t(\mathbf{x}) \geq A\}$$

- Using the asymptotic property

$$D_t^- = \{\mathbf{x}; R_t(\mathbf{x}) \geq A\} = \{\mathbf{x}; f_t(\mathbf{x}) \leq \theta_t^A\}$$



Search for  $\theta$  such that

$$1 - \beta_t^\theta \geq A\alpha_t^\theta$$

## Input:

- A training set  $\mathcal{T} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$   $\mathbf{x}_i \in \mathcal{X}, y_i \in \{-1, +1\}$
- A set (class) of weak classifiers  $\mathcal{H} = \{h_i\}_{i=1}^M$  each one associated to one measurement
- Desired final false negative rate  $\alpha$  and false positive rate  $\beta$
- The number of iterations  $T$

Set  $A = (1 - \beta)/\alpha$  and  $B = \beta/(1 - \alpha)$

**for**  $t = 1, \dots, T$

1. Find  $h_t \in \mathcal{H}$  by AdaBoost using  $\mathcal{T}$  and add it to the strong classifier  $f_t = f_{t-1} + h_t$
2. Find the decision thresholds  $\theta_t^A$  and  $\theta_t^B$  for  $f_t$
3. Remove from training samples  $\mathbf{x} \in \mathcal{T}$  for which  $f_t(\mathbf{x}) \geq \theta_t^B$  or  $f_t(\mathbf{x}) \leq \theta_t^A$
4. Sample new data into the training set  $\mathcal{T}$

**end**

## Output:

- Ordered set of weak classifiers  $\{h_t\}_{t=1}^T$
- The decision thresholds  $\{\theta_t^A\}_{t=1}^T$  and  $\{\theta_t^B\}_{t=1}^T$

**Given:**  $h_t, \theta_t^A, \theta_t^B, \gamma$  ( $t = 1, \dots, T$ )

**Input:** an object  $\mathbf{x}$

**For**  $t = 1, \dots, T$  (SPRT execution)

    If  $f_t(\mathbf{x}) \geq \theta_t^B$ , classify  $\mathbf{x}$  to the class +1 and terminate

    If  $f_t(\mathbf{x}) \leq \theta_t^A$ , classify  $\mathbf{x}$  to the class -1 and terminate

**end**

If  $f_T(\mathbf{x}) > \gamma$ , classify  $\mathbf{x}$  as +1. Classify  $\mathbf{x}$  as -1 otherwise.

## Non-symmetric decision problems (e.g. face detection)

- Positive (face) and negative (background) class complexities are unbalanced
- Positive samples difficult to collect
- Positive samples appear rarely during classification
- Missed detection rate more serious than false positive rate

Set  $\beta = 0$  in WaldBoost learning. Then

SPRT thresholds

$$A = \frac{1 - 0}{\alpha} = \frac{1}{\alpha}$$

$$B = \frac{0}{1 - \alpha} = 0$$

SPRT strategy

$$S_t^* = \begin{cases} +1, & R_t \leq 0 \\ -1, & R_t \geq 1/\alpha \\ \#, & 0 < R_t < 1/\alpha \end{cases}$$

WaldBoost strategy

$$S_t = \begin{cases} -1, & f_t(\mathbf{x}) \leq \theta_t^A \\ \#, & \theta_t^A < f_t(\mathbf{x}) \end{cases}$$

- Since  $R_t(\mathbf{x}) > 0$ , only decisions to the class  $-1$  (background) allowed
- Only the  $-1$  class pruned during training
- All measurements evaluated for positive samples (does not affect the speed much since they are rare)

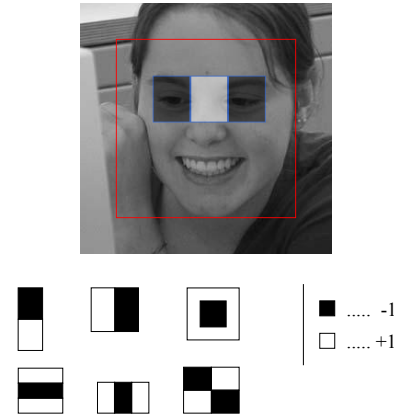
# Application 1: Fast Face Detection



scanning process



used measurements



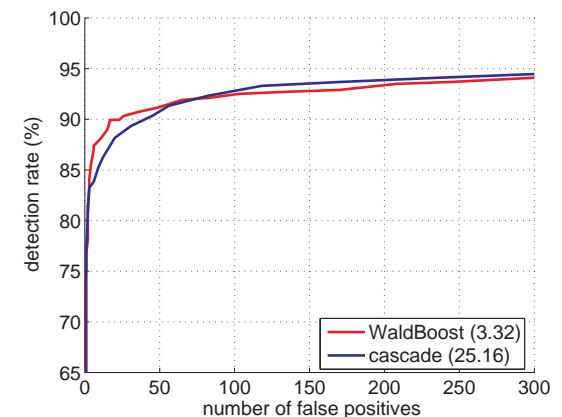
results



## Results

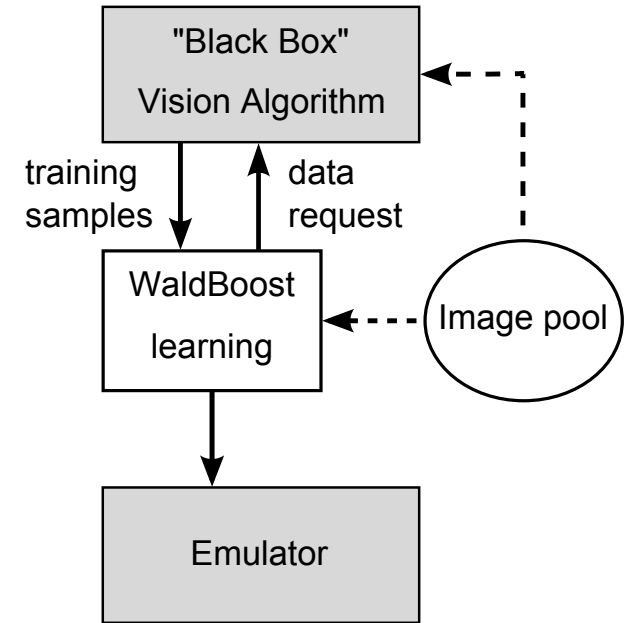
- State-of-the-art detection results
- **Fastest detector** among those with comparable detection rate
- Principle way of learning a cascade-like classifier
  - only two global parameters:  $\alpha$  and  $\beta$
- Web demo

<http://cmp.felk.cvut.cz/demos/FaceDetection/>



## Idea

- Assume: slow binary-valued “black-box” algorithm is available
- Use the black-box algorithm to generate a training set
- WaldBoost optimises:
  - output similarity
  - evaluation speed



## Results

- A general framework for speeding up existing algorithms by a sequential classifier learned by the WaldBoost algorithm
- Demonstrated on Hessian-Laplace and Kadir-Brady saliency detectors

	HL	KB
<b>original</b>	0.9s	1m 48s
<b>SURF</b>	0.09s	—
speed-up	10×	—
<b>WaldBoost</b>	0.10s	0.76s
speed-up	9×	142×



- Time to decision vs. precision trade-off problem was formalised as a constrained optimisation problem
- WaldBoost overcomes limitations of SPRT to a priori ordered measurements and known multi-dimensional pdf's
- Even though WaldBoost is greedy (sub-optimal), it is highly flexible and has interesting applications
  - Fast face detection
  - Fast emulators of interest point detectors
- Other related papers
  - Randomized RANSAC with SPRT [Matas & Chum, ICCV 2005]
  - On-line WaldBoost for tracking [Grabner et al., ICPR 2008]
  - Sequential correspondence selection by cosegmentation [Čech et al., PAMI 2009]

**Thank you for attention!**

The End