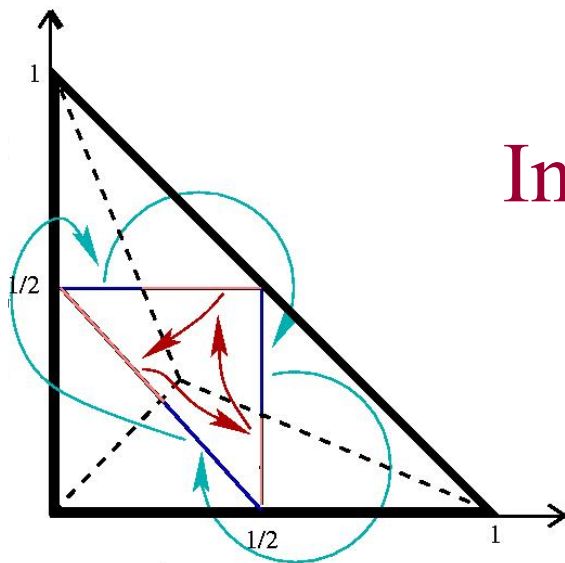# Introduction to Boosting

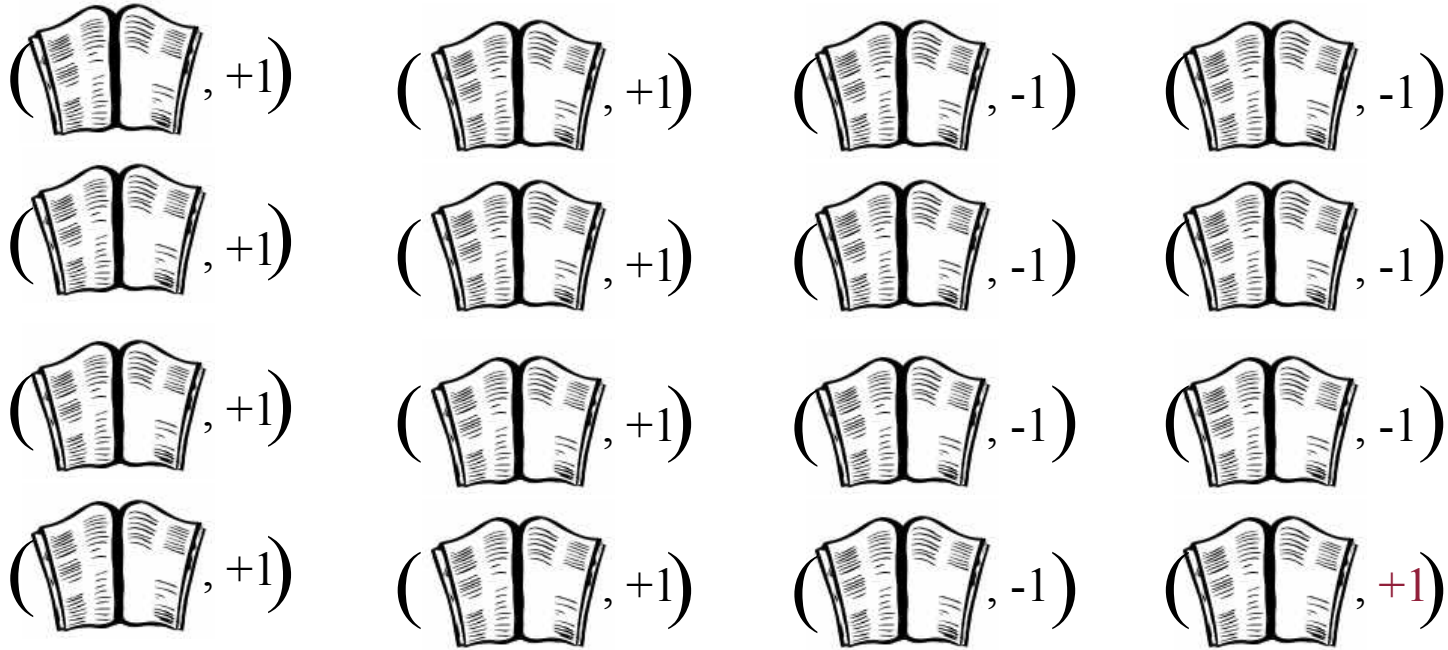## Cynthia Rudin
## PACM, Princeton University

### Advisors
Ingrid Daubechies and Robert Schapire

# Say you have a database of news articles…

( [newspaper], +1 )   ( [newspaper], +1 )   ( [newspaper], -1 )   ( [newspaper], -1 )

( [newspaper], +1 )   ( [newspaper], +1 )   ( [newspaper], -1 )   ( [newspaper], -1 )

( [newspaper], +1 )   ( [newspaper], +1 )   ( [newspaper], -1 )   ( [newspaper], -1 )

( [newspaper], +1 )   ( [newspaper], +1 )   ( [newspaper], -1 )   ( [newspaper], +1 )

where articles are labeled '+1' if the category is "entertainment", and '-1' otherwise.

Your goal is: Given a new article [newspaper], find its label.

This is not easy, there are noisy datasets, high dimensions.

# Examples of Statistical Learning Tasks:

- Optical Character Recognition (OCR) (post office, banks), object recognition in images.

- Bioinformatics (analysis of gene array data for tumor detection, protein classification, etc.)

- Webpage classification (search engines), email filtering, document retrieval

- Semantic classification for speech, automatic .mp3 sorting

- Time-series prediction (regression)

Huge number of applications, but all have high dimensional data

# Examples of classification algorithms:

- SVM's (Support Vector Machines – large margin classifiers)
- Neural Networks
- Decision Trees / Decision Stumps (CART)
- RBF Networks
- Nearest Neighbors
- Bayes Net

# Which is the best?
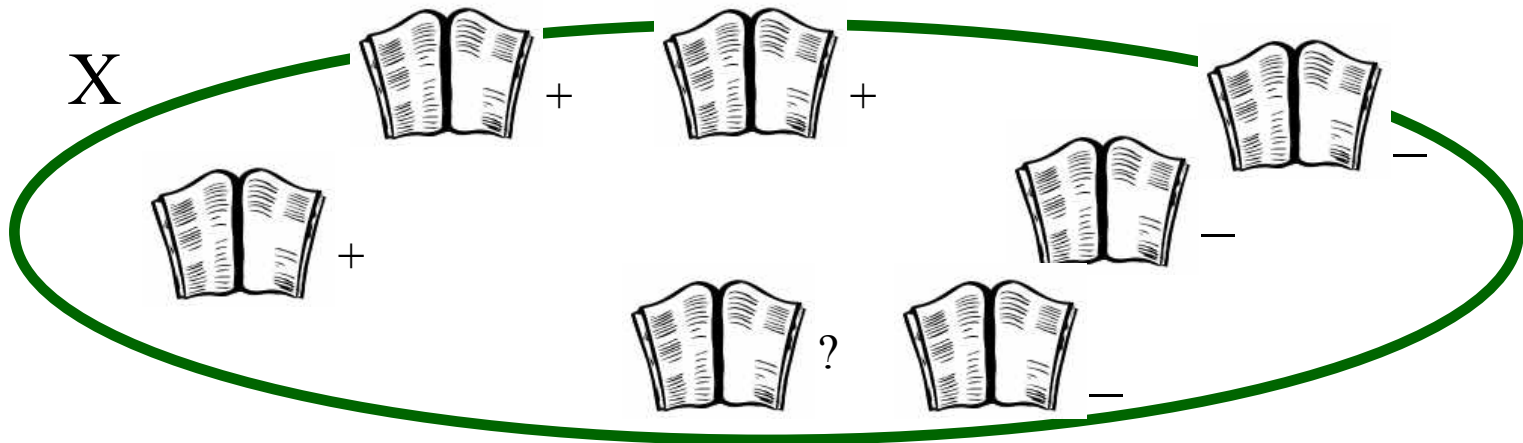
Depends on amount and type of data, and application!

It's a tie between SVM's and **Boosted** Decision Trees/Stumps for general applications.

One can always find a problem where a particular algorithm is the best. Boosted convolutional neural nets are the best for OCR (Yann LeCun et al).

Training Data: $\{(\mathbf{x}_i, y_i)\}_{i=1..m}$ where $(\mathbf{x}_i, y_i)$ is chosen iid from an unknown probability distribution on $X \times \{-1, 1\}$.

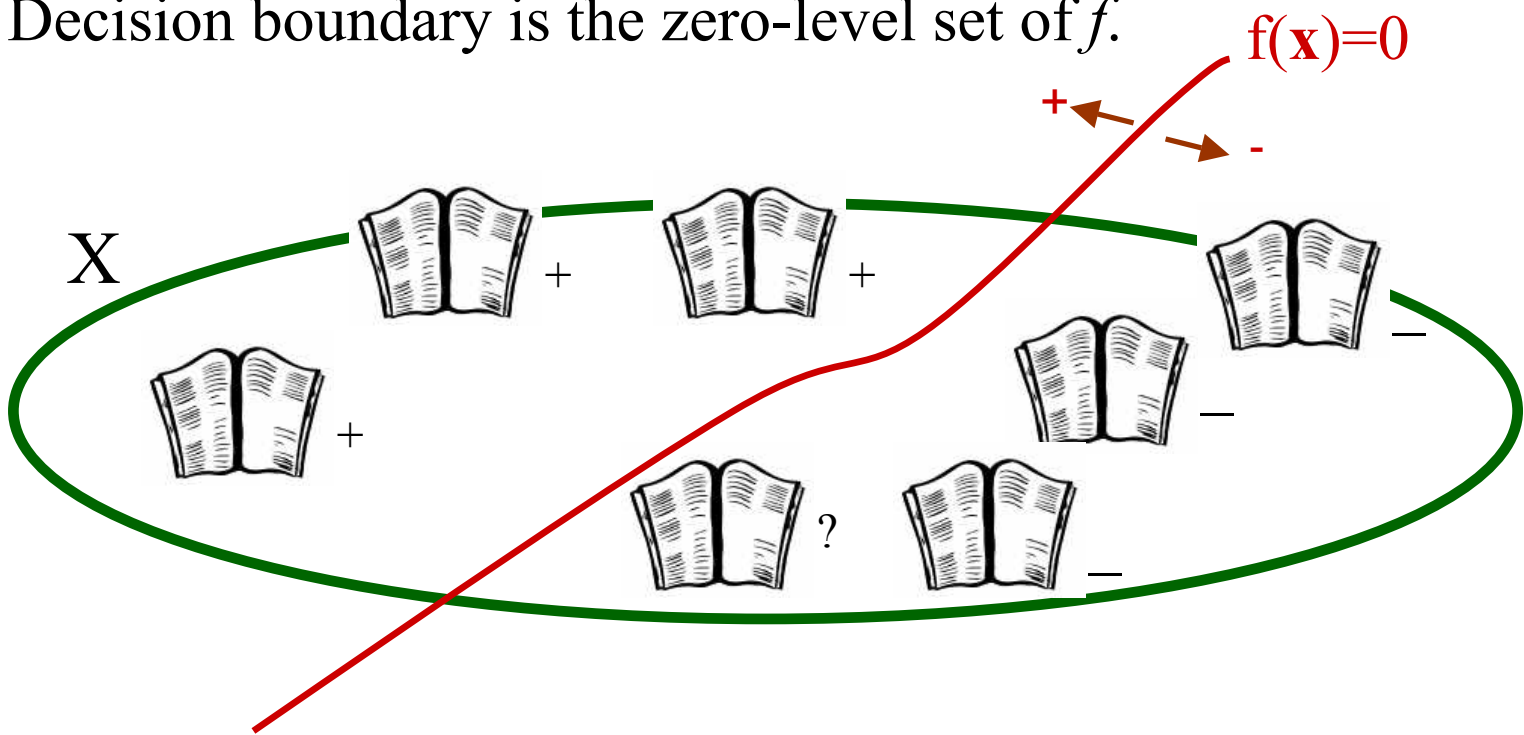"space of all possible articles"    "labels"

Huge Question: Given a new random example $\mathbf{x}$, can we predict its correct label with high probability? That is, can we _generalize_ from our training data?

Huge Question: Given a new random example **x**, can we predict its correct label with high probability? That is, can we _generalize_ from our training data?

Yes!!! That's what the field of statistical learning is all about.

The goal of statistical learning is to characterize points from an unknown probability distribution when given a representative sample from that distribution.

# How do we construct a classifier?

• Divide the space $X$ into two sections, based on the sign of a function $f : X \rightarrow R$.

• Decision boundary is the zero-level set of $f$.



Classifiers divide the space into two pieces for binary classification. Multiclass classification can always be reduced to binary.

# ❧ Overview of Talk ❧

- The Statistical Learning Problem (done)

- Introduction to Boosting and AdaBoost

- AdaBoost as Coordinate Descent

- The Margin Theory and Generalization

# Say we have a "weak" learning algorithm:

- A weak learning algorithm produces weak classifiers.
- (Think of a weak classifier as a "rule of thumb")

Examples of weak classifiers for "entertainment" application:

$$h_1(\;\text{📖}\;) = \begin{cases} +1 \text{ if } \text{📖} \text{ contains the term "movie",} \\ -1 \text{ otherwise} \end{cases}$$

$$h_2(\;\text{📖}\;) = \begin{cases} +1 \text{ if } \text{📖} \text{ contains the term "actor",} \\ -1 \text{ otherwise} \end{cases}$$

$$h_3(\;\text{📖}\;) = \begin{cases} +1 \text{ if } \text{📖} \text{ contains the term "drama",} \\ -1 \text{ otherwise} \end{cases}$$

Wouldn't it be nice to combine the weak classifiers?

# **Boosting algorithms combine weak classifiers in a meaningful way.**

Example:

$$f(\text{📖}) = \text{sign}[.4\ h_1 (\text{📖}) + .3\ h_2 (\text{📖}) + .3\ h_3 (\text{📖})]$$

A boosting algorithm takes as input:
Suppose the article contains the term "movie", and the word "drama", but not the word "actor":

- the weak learning algorithm which produces the weak classifiers
- a large training database
The value of f is sign[.4-.3+.3] = 1, so we label it +1.
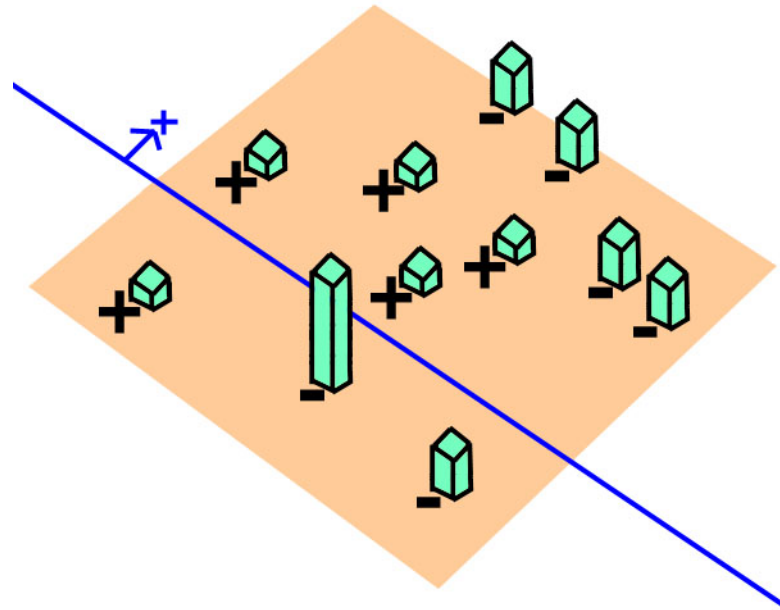
and outputs:

- the coefficients of the weak classifiers to make the combined classifier

# Two ways to use a Boosting Algorithm:

• As a way to increase the performance of already `strong` classifiers.

    • Ex. neural networks, decision trees

• "On their own" with a really basic weak classifier

    • Ex. decision stumps

# AdaBoost
## (Freund and Schapire '95)

-Start with a uniform distribution ("weights") over training examples.
  (The weights tell the weak learning algorithm which examples are important.)
  At the end, make (carefully!) a linear combination of the weak classifiers
  obtained at all iterations.

-Request a weak classifier from the weak learning algorithm, $h_{j_t}:X\rightarrow\{-1,1\}$.

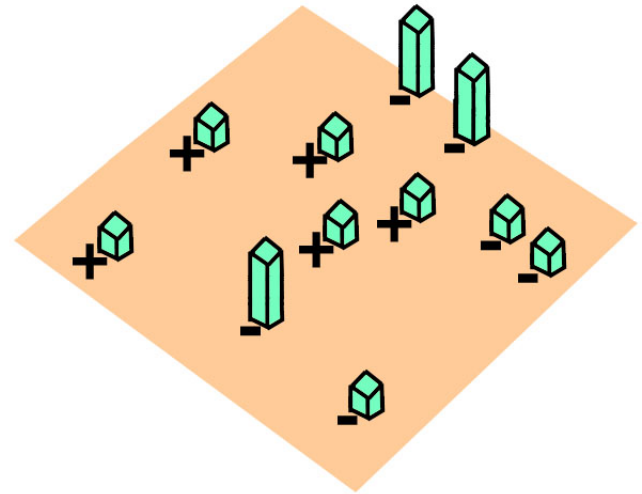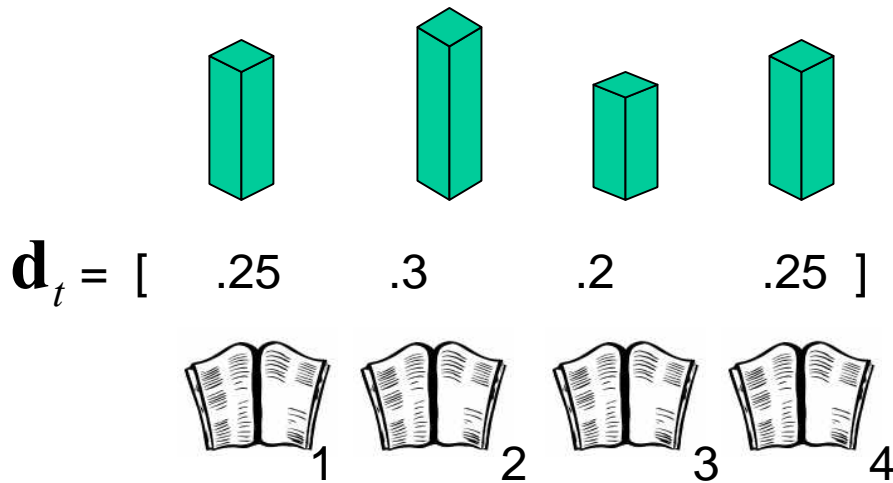$$f_{final}(\mathbf{x}) = sign(\lambda_1 h_1(\mathbf{x}) + \ldots + \lambda_n h_n(\mathbf{x}))$$

-Increase the weights on the training examples that were misclassified.

-(Repeat)

# AdaBoost

Define three important things:

$$\mathbf{d}_t \in \mathrm{R}^m := \text{distribution (``weights'') over examples at time t}$$

$\mathbf{d}_t = [$ .25     .3     .2     .25 $]$
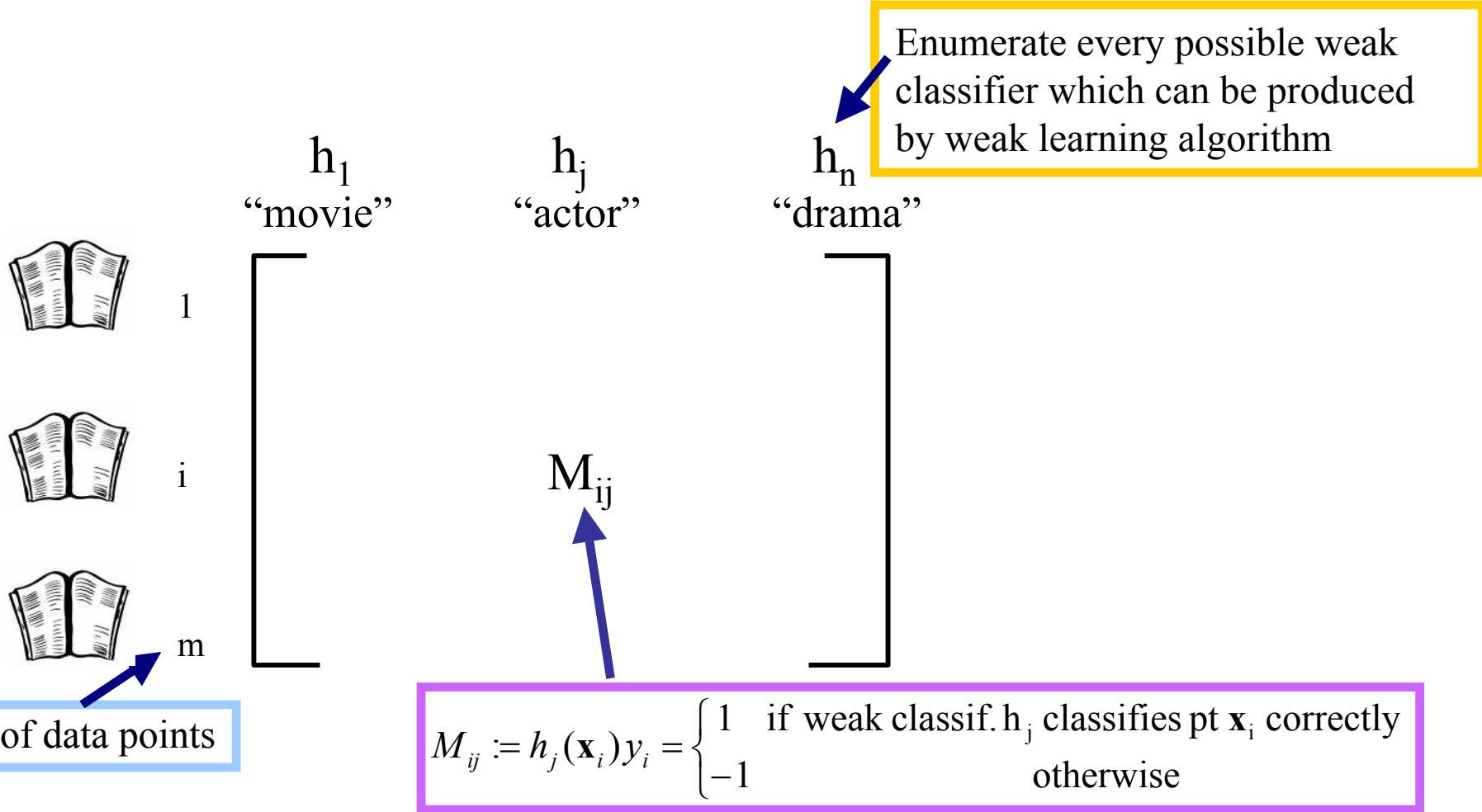
1     2     3     4

# AdaBoost

Define three important things:

$\lambda_t \in R^n :=$ coeffs of weak classifiers for the linear combination

$$f_t(\mathbf{x}) = \text{sign}(\lambda_{t,1} h_1(\mathbf{x}) + \ldots + \lambda_{t,n} h_n(\mathbf{x}))$$
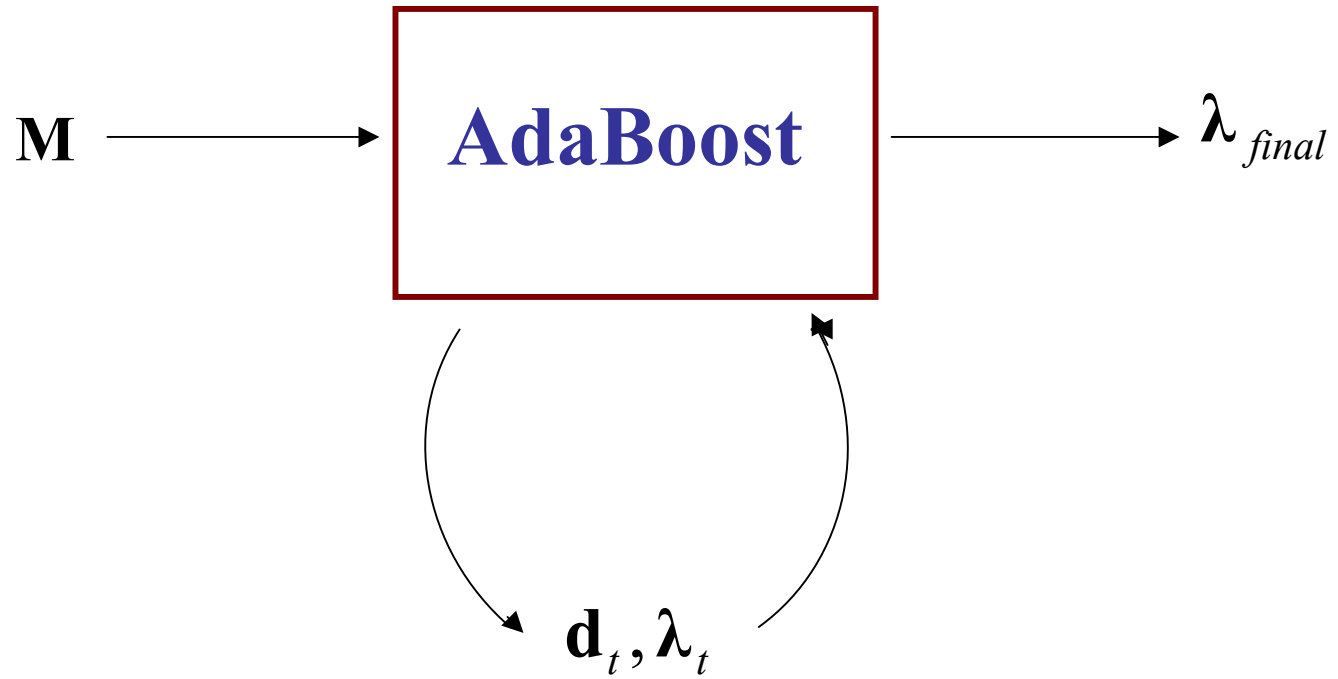
# AdaBoost

Define:

$$\mathbf{M} \in \mathbf{R}^{\,m \times n} := \text{matrix of hypotheses and data}$$

Enumerate every possible weak classifier which can be produced by weak learning algorithm

$$
\begin{array}{ccc}
h_1 & h_j & h_n \\
\text{"movie"} & \text{"actor"} & \text{"drama"}
\end{array}
$$

1

i

$$M_{ij}$$

m

# of data points

$$M_{ij} := h_j(\mathbf{x}_i) y_i = \begin{cases} 1 & \text{if weak classif. } h_j \text{ classifies pt } \mathbf{x}_i \text{ correctly} \\ -1 & \text{otherwise} \end{cases}$$

The matrix M has too many columns to actually be enumerated. M acts as the only input to AdaBoost.

$$\mathbf{M} \longrightarrow \boxed{\textbf{AdaBoost}} \longrightarrow \boldsymbol{\lambda}_{\textit{final}}$$

$$\mathbf{d}_t, \boldsymbol{\lambda}_t$$

# AdaBoost (Freund and Schapire 95)

$$\boldsymbol{\lambda}_1 = 0$$

Initialize coeffs to 0

for $t = 1..T_{final}$

$$d_{t,i} = \frac{1}{\displaystyle\sum_{i'=1}^{m} e^{-(\mathbf{M}\boldsymbol{\lambda}_t)_{i'}}} e^{-(\mathbf{M}\boldsymbol{\lambda}_t)_i} \quad \text{for all } i$$

Calculate (normalized) distribution

$$j_t \in \arg\max_j (\mathbf{d_t^T M})_j$$

Request weak classif. from weak learning algorithm

$$r_t = (\mathbf{d_t^T M})_{j_t}$$

$$\alpha_t = \frac{1}{2} \ln\left(\frac{1 + r_t}{1 - r_t}\right)$$

$$\boldsymbol{\lambda}_{t+1} = \boldsymbol{\lambda}_t + \alpha_t \mathbf{e}_{j_t}$$

Update linear combo of weak classifiers

end for

# AdaBoost (Freund and Schapire 95)

$$\boldsymbol{\lambda}_1 = 0$$

for $t = 1..T_{final}$

$$d_{t,i} = \frac{1}{\sum_{i'=1}^{m} e^{-(\mathbf{M}\boldsymbol{\lambda}_t)_{i'}}} e^{-(\mathbf{M}\boldsymbol{\lambda}_t)_i} \quad \text{for all } i$$

$$j_t \in \arg\max_j (\mathbf{d_t^T M})_j$$

$$r_t = (\mathbf{d_t^T M})_{j_t}$$

$$\alpha_t = \frac{1}{2} \ln\left(\frac{1+r_t}{1-r_t}\right)$$

$$\boldsymbol{\lambda}_{t+1} = \boldsymbol{\lambda}_t + \alpha_t \mathbf{e}_{j_t}$$

end for

"Edge" or "correlation" of weak classifier $j_t$.

$$(\mathbf{d}_t^T \mathbf{M})_{j_t} = \sum_{i=1}^{m} d_{t,i} y_i h_{j_t}(\mathbf{x}_i) = \mathbf{E}_{\mathbf{d}_t}[y_i h_{j_t}]$$

# ❧ AdaBoost as Coordinate Descent ❧

Breiman, Mason et al., Duffy and Helmbold, etc. noticed that AdaBoost is a coordinate descent algorithm.

- Coordinate descent is a minimization algorithm like gradient descent, except that we only move along coordinates.

- We cannot calculate the gradient because of the high dimensionality of the space!

- "coordinates" = weak classifiers
  "distance to move in that direction" = the update $\alpha_t$

AdaBoost minimizes the following function via coordinate descent:

$$F(\boldsymbol{\lambda}) := \sum_{i=1}^{m} e^{-(\mathbf{M}\boldsymbol{\lambda})_i}$$

Choose a direction:   $j_t \in \arg\max_j (\mathbf{d_t^T M})_j$

Choose a distance to move in that direction:

$$r_t = (\mathbf{d_t^T M})_{j_t}$$

$$\alpha_t = \frac{1}{2} \ln\left(\frac{1 + r_t}{1 - r_t}\right)$$

$$\boldsymbol{\lambda}_{t+1} = \boldsymbol{\lambda}_t + \alpha_t \mathbf{e}_{j_t}$$

The function $F(\lambda) := \sum_{i=1}^{m} e^{-(\mathbf{M}\lambda)_i}$ is convex:

1) If the data is non-separable by the weak classifiers, the minimizer of $F$ occurs when the size of $\lambda$ is finite.

(This case is ok. AdaBoost converges to something we understand.)

2) If the data is separable, the minimum of $F$ is 0

(This case is confusing!)

The original paper suggested that AdaBoost would probably overfit…

But it didn't in practice!

Why not?

The margin theory!

# ❧ Boosting and Margins ☙

- We want the boosted classifier (defined via $\lambda$) to generalize well, i.e., we want it to perform well on data that is not in the training set.

- The margin theory: The margin of a boosted classifier indicates whether it will generalize well. (Schapire et al. '98)

- Large margin classifiers work well in practice, (but there's more to this story).

Think of the margin as the confidence of a prediction.

# Generalization Ability of Boosted Classifiers

Can we guess whether a boosted classifier $f$ generalizes well?

 • Can not calculate $\mathrm{Pr}_{error}(f)$

Minimize the rhs of a (loose) inequality such as this one (Schapire et al.)

When there are no training errors, with probability at least 1-$\delta$,

$$\mathrm{Pr}_{error}(f) \leq O\left( \frac{1}{\sqrt{m}} \left( \frac{d \log^2(m/d)}{(\mu(f))^2} + \log(1/\delta) \right)^{1/2} \right).$$

Probability that classifier $f$ makes an error on a random position $\mathbf{x} \in X$

# of training examples

d=VC dim. of hyp. space, d$\leq$m

margin of $f$

# The margin theory:

When there are no training errors, with high probability:
(Schapire et al, '98)

$$\Pr_{error}(f) \leq \tilde{O}\left( \frac{\sqrt{d/m}}{\mu(f)} \right)$$

d=VC dim. of hyp. space, d≤m

# of training examples

margin of $f$

Probability that classifier $f$ makes an error on a random position $\mathbf{x} \in X$

Large margin = better generalization = smaller probability of error

For Boosting, the margin of combined classifier $f_\lambda$
(where $f_\lambda := \text{sign}(\lambda_1 h_1 + \ldots + \lambda_n h_n)$) is defined by

$$\text{margin} := \mu(f_\lambda) := \min_i \frac{(\mathbf{M}\lambda)_i}{\|\lambda\|_1}.$$

# Does AdaBoost produce maximum margin classifiers?

(AdaBoost was invented before the margin theory…)

(Grove and Schuurmans '98)
- yes, empirically.

(Schapire, et al. '98)
- proved AdaBoost achieves at least half the maximum possible margin.

(Rätsch and Warmuth '03)
- yes, empirically.
- improved the bound.

(R, Daubechies, Schapire '04)
- no, it doesn't.

AdaBoost performs mysteriously well!

AdaBoost performs better than algorithms which are designed to maximize the margin

Still open: Why does AdaBoost work so well?
Does AdaBoost converge?

Better / more predictable boosting algorithms!