

BAGGING AND BOOSTING

Evgeny Burnaev

Skoltech, Moscow, Russia

- 1 MOTIVATION: CLASSIFICATION PROBLEM
- 2 BAGGING
- 3 BOOSTING
- 4 MARGIN IN CASE OF BOOSTING

1 MOTIVATION: CLASSIFICATION PROBLEM

2 BAGGING

3 BOOSTING

4 MARGIN IN CASE OF BOOSTING

CLASSIFICATION PROBLEM I

- A predictor, feature $\mathbf{x} \in \mathbb{R}^p$ has distribution D
- $h(\mathbf{x})$ is a deterministic function from some concept class
- **Goal:**
 - Based on m training pairs $(\mathbf{x}_i, y_i = h(\mathbf{x}_i))$ drawn from D produce a classifier $\hat{h}(\mathbf{x}) \in \{0, 1\}$
 - Choose \hat{h} to have low generalization error

$$R(\hat{h}) = \mathbb{E}_D \left[1_{\hat{h}(\mathbf{x}) \neq h(\mathbf{x})} \right]$$

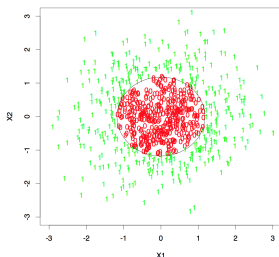
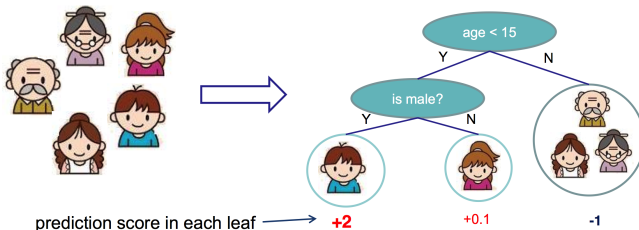


FIGURE : “Sphere” in \mathbb{R}^{10}

CLASSIFICATION AND REGRESSION TREES (CART) I

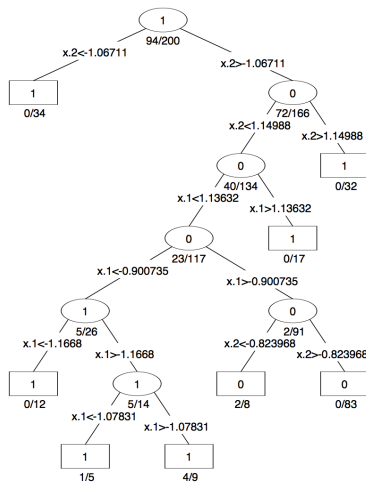
- Classification and Regression Trees:
 - Decision rules
 - Contains one score in each leaf value

Input: age, gender, occupation,... \Rightarrow Does the person like computer games?



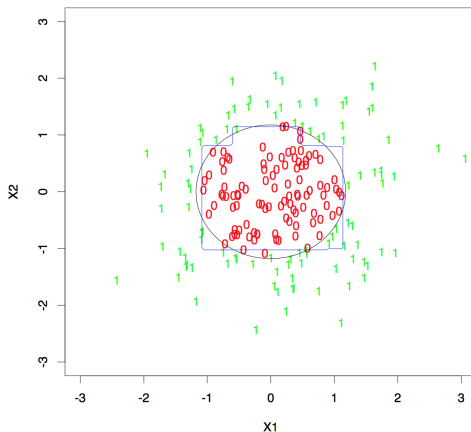
CART II

Sample of size 200



CLASSIFICATION PROBLEM I

Sample of size 200



In case of “Sphere” in \mathbb{R}^{10} CART produces a rather noisy and inaccurate rule $\hat{h}(\mathbf{x})$, with error rates around 30%

- 1 MOTIVATION: CLASSIFICATION PROBLEM
- 2 BAGGING
- 3 BOOSTING
- 4 MARGIN IN CASE OF BOOSTING

MODEL AVERAGING

Classification trees can be simple, but often produce noise (bushy) or weak (stunted) classifiers

- Bagging (Breiman, 1996): Fit many large trees to bootstrap-resampled versions of the training data, and classify by majority vote
- Boosting (Freund & Shapire, 1996): Fit many large or small trees to reweighted versions of the training data. Classify by weighted majority vote

In general

Boosting \succ Bagging \succ Single Tree

STATISTICS: BOOTSTRAP

- Model:
 - we have i.i.d. sample $\{\mathbf{x}_i\}_{i=1}^m \subset \mathbb{R}^1$, generated by some distribution function F
 - We consider some statistics $T_m = g(\mathbf{x}_1, \dots, \mathbf{x}_m)$.
- Problem: estimate variance $\mathbb{V}_F(T_n)$, which depends on some unknown distribution function F

EXAMPLE

Let us consider $T_m = \bar{\mathbf{x}}_m$. Then $\mathbb{V}_F(T_m) = \sigma^2/m$, where $\sigma^2 = \int (\mathbf{x} - \mu)^2 dF(\mathbf{x})$ and $\mu = \int \mathbf{x} dF(\mathbf{x})$. Thus, the variance T_m is a function of F

BOOTSTRAP IDEA

STEP 1. Estimate $\mathbb{V}_F(T_m)$ using $\mathbb{V}_{\hat{F}_m}(T_m)$, where

$$\hat{F}_m(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m 1(\mathbf{x}_i \leq \mathbf{x})$$

STEP 2. Approximate $\mathbb{V}_{\hat{F}_m}(T_m)$ using Monte-Carlo sampling from \hat{F}_m

EXAMPLE

For $T_m = \bar{\mathbf{x}}_m$, $\mathbb{V}_{\hat{F}_m}(T_m) = \hat{\sigma}^2/m$, where $\hat{\sigma}^2 = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i - \bar{\mathbf{x}}_m)^2$. Thus Step 1 is sufficient in this case. However, often we can not provide an explicit expression for $\mathbb{V}_{\hat{F}_m}(T_m)$. Thus we can use Step 2

GENERAL SCHEME

STEP 1. In “Real” World

$$F \Rightarrow \mathbf{x}_1, \dots, \mathbf{x}_m \Rightarrow T_m = g(\mathbf{x}_1, \dots, \mathbf{x}_m)$$

STEP 2. In “Bootstrap” World

$$\hat{F}_m \Rightarrow \{\mathbf{x}_1^*, \dots, \mathbf{x}_m^*\} \Rightarrow T_m^* = g(\mathbf{x}_1^*, \dots, \mathbf{x}_m^*)$$

- **Problem:** how to generate $\mathbf{x}_1^*, \dots, \mathbf{x}_m^*$ from \hat{F}_m ?
- **Solution:** \hat{F}_m has a mass $\frac{1}{m}$ in each of sample point \mathbf{x}_i , $i = 1, \dots, m \Rightarrow$ generating from \hat{F}_m is equivalent to selection with replacement from the initial sample $\{\mathbf{x}_i\}_{i=1}^m$

BOOTSTRAP VARIANCE ESTIMATION

In order to estimate variance of a functional using bootstrap:

1. Select $\mathbf{x}_1^*, \dots, \mathbf{x}_m^* \sim \hat{F}_m$
2. Calculate $T_m^* = g(\mathbf{x}_1^*, \dots, \mathbf{x}_m^*)$
3. Repeat steps 1 and 2 until you get $T_m^{*,1}, \dots, T_m^{*,B}$
4. Set

$$v_{boot} = \frac{1}{B} \sum_{b=1}^B \left(T_m^{*,b} - \frac{1}{B} \sum_{r=1}^B T_m^{*,r} \right)^2$$

Thus we get that

$$\mathbb{V}_F(T_m) \approx \mathbb{V}_{\hat{F}}(T_m) \approx v_{boot}$$

BAGGING

- Bagging or bootstrap averaging averages a given procedure over many samples to reduce its variance
- Let us denote by
 - $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$ a sample of size m
 - $\hat{h}_S(\mathbf{x})$ a classifier, such as a tree, trained using the sample S
- To bag \hat{h} we draw bootstrap samples $S^{*,1}, \dots, S^{*,B}$ each of size m with replacement from the training data
- Then

$$\hat{h}_{\text{bag}}(\mathbf{x}) = \text{MajorityVote} \left\{ \hat{h}_{S^{*,b}}(\mathbf{x}) \right\}_{b=1}^B$$

- Bagging can dramatically reduce the variance of unstable procedures (like trees), leading to improved prediction
- However any simple structure in h (e.g. a tree) is lost

EXAMPLE: BAGGING

Original Tree



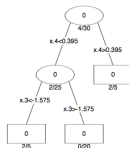
Bootstrap Tree 1



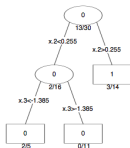
Bootstrap Tree 2



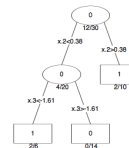
Bootstrap Tree 3



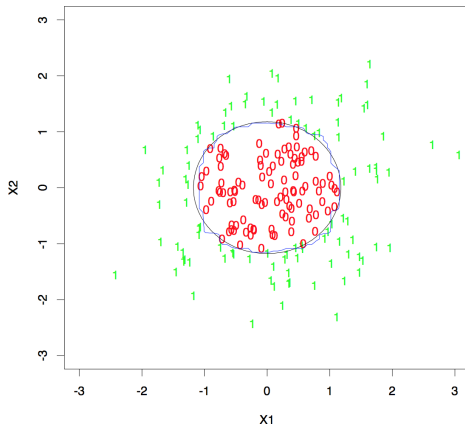
Bootstrap Tree 4



Bootstrap Tree 5



DECISION BOUNDARY: BAGGING



“Sphere” in \mathbb{R}^{10} : Bagging averages many trees, and produces smoother decision boundaries

- 1 MOTIVATION: CLASSIFICATION PROBLEM
- 2 BAGGING
- 3 BOOSTING
- 4 MARGIN IN CASE OF BOOSTING

EXAMPLE: SPAM FILTERING

- **problem:** filter out spam (junk email)
- gather large collection of examples of **spam** and **non-spam**

From: yoav@att.com	Rob, can you review a paper...	non-spam
From: xa412@hotmail.com	Earn money without working!!!! ...	spam
⋮	⋮	⋮

- **goal:** get computer learn from examples to distinguish spam from non-spam
- **main observation:**
 - **easy** to find “rules of thumb” that are “often” correct

if 'v1agr@' occurs in message, then predict “**spam**”
 - **hard** to find single rule that is very highly accurate

THE BOOSTING APPROACH I

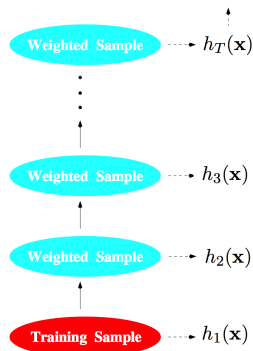
- devise computer program for deriving rough rules of thumb
- apply procedure to subset of emails
- obtain rule of thumb
- apply to 2nd subset of emails
- obtain 2nd rule of thumb
- repeat T times
- aggregate

THE BOOSTING APPROACH II

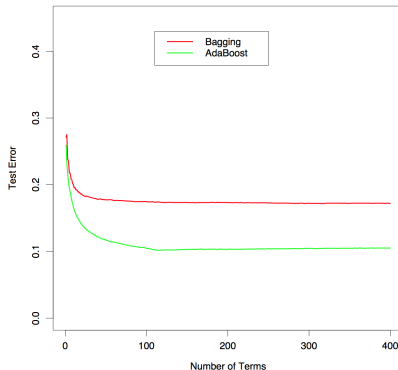
Final Classifier

1. How to choose examples on each round?
 - concentrate on “hardest” examples (those most often misclassified by previous rules of thumb)
2. How to combine rules of thumb into single prediction rule?
 - take (weighted) majority vote of rules of thumb

$$h(\mathbf{x}) = \text{sign} \left[\sum_{t=1}^T \alpha_t h_t(\mathbf{x}) \right]$$



BAGGING AND BOOSTING



- 2000 points, “Sphere” in \mathbb{R}^{10} ; Bayes error rate is 0%
- Trees are grown Best First without pruning
- Leftmost iteration is a single tree

PAC LEARNING MODEL: NOTATIONS

- X : set of all possible instances or examples, e.g. the set of all men and women characterized by their height and weight
- $c : X \rightarrow \{0, 1\}$: the target concept to learn; can be identified with its support $\{\mathbf{x} \in X : c(\mathbf{x}) = 1\}$
- C : concept class, a set of target concepts c
- D : target distribution, a fixed probability distribution over X . Training and test examples are drawn according to D

PAC LEARNING MODEL: NOTATIONS

- S : training sample
- H : set of concept hypothesis, e.g. the set of all linear classifiers
- The learning algorithm receives sample S and selects a hypothesis h_S from H approximating c

PAC LEARNING MODEL: ERRORS

- **True error or generalization error** of h with respect to the target concept c and distribution D

$$R(h) = \mathbb{P}_{\mathbf{x} \sim D} [h(\mathbf{x}) \neq c(\mathbf{x})] = \mathbb{E}_{\mathbf{x} \sim D} [1_{h(\mathbf{x}) \neq c(\mathbf{x})}]$$

- **Empirical error:** average error of h on the training sample S drawn according to distribution D

$$\begin{aligned}\hat{R}_S(h) &= \mathbb{P}_{\mathbf{x} \sim \hat{D}} [h(\mathbf{x}) \neq c(\mathbf{x})] \\ &= \mathbb{E}_{\mathbf{x} \sim \hat{D}} [1_{h(\mathbf{x}) \neq c(\mathbf{x})}] = \frac{1}{m} \sum_{i=1}^m 1_{h(\mathbf{x}_i) \neq c(\mathbf{x}_i)}\end{aligned}$$

- Note:

$$R(h) = \mathbb{E}_{S \sim D^m} [\hat{R}_S(h)]$$

PAC LEARNING MODEL: DEFINITION

- **PAC Learning:** Probably Approximately Correct Learning (Valiant, 1984)
- **Definition:** concept class C is PAC-learnable if there exists a learning algorithm L such that
 - for all $c \in C$, $\epsilon > 0$, $\delta > 0$, and all distributions D ,

$$\mathbb{P}_{S \sim D^m} [R(h_S) \leq \epsilon] \geq 1 - \delta,$$

- for samples S of size $m = \text{poly}(1/\epsilon, 1/\delta)$ for a fixed polynomial
- Such L is called a strong Learner

PAC LEARNING MODEL: COMMENTS

- Concept class C is known to the algorithm
- Distribution-free model: no assumption on D
- Both training and test examples drawn $\sim D$
- Probably: confidence $1 - \delta$
- Approximately correct: accuracy $1 - \epsilon$
- Efficient PAC-Learning: L runs in time $\text{poly}(1/\epsilon, 1/\delta, N, \text{size}(c))$

WEAK LEARNING

- **Definition** : concept class C is weakly PAC-learnable if there exists a (weak) learning algorithm L and $\gamma > 0$ such that:
 - for all $c \in C$ and $\delta > 0$, and all distributions D ,

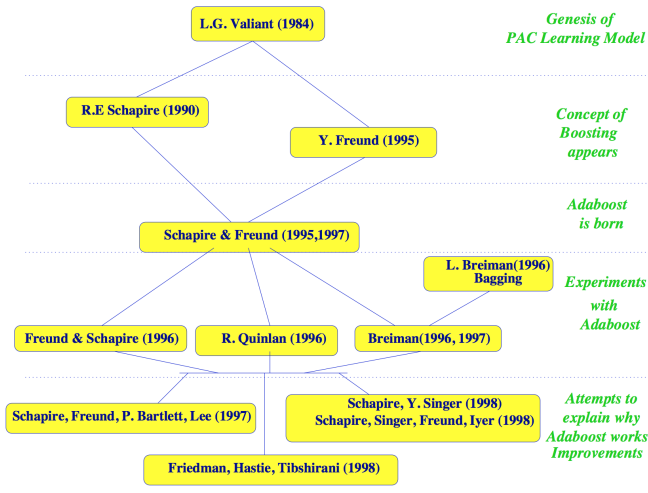
$$\mathbb{P}_{S \sim D^m} \left[R(h_S) \leq \frac{1}{2} - \gamma \right] \geq 1 - \delta,$$

- for samples S of size $m = \text{poly}(1/\delta)$ for a fixed polynomial

THE BOOSTING APPROACH III

- Finding simple relatively accurate base classifiers often not hard \leftarrow weak learner
- Main ideas:
 - use weak learner to create a strong learner
 - combine base classifiers returned by weak learner (ensemble method)
- But how should the base classifiers be combined?

HISTORY OF BOOSTING



BOOSTING A WEAK LEARNER

- Weak learner L produces an h with error rate $\beta = (\frac{1}{2} - \gamma) < \frac{1}{2}$ with $\Pr \geq (1 - \delta)$ for any D
- L has access to continuous stream of training data a class oracle
 - L learns h_1 on first m training points
 - L randomly filters the next batch of training points, extracting $m/2$ points correctly classified by h_1 , $m/2$ incorrectly classified, and produces h_2
 - L builds a third training set of m points for which h_1 and h_2 disagree, and produces h_3
 - L outputs

$$h = \text{MajorityVote}(h_1, h_2, h_3)$$

- **Theorem** (Schapire, 1990): “The Strength of Weak Learnability”

$$R(h) \leq 3\beta^2 - 2\beta^3 < \beta$$

ADABOOST

$$H \subseteq \{-1, +1\}^X$$

$$\text{AdaBoost}(S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\})$$

1. **for** $i \leftarrow 1$ **to** m **do**
2. $D_1(i) \leftarrow \frac{1}{m}$
3. **for** $t \leftarrow 1$ **to** T **do**
4. $h_t \leftarrow$ base classif. with small $\epsilon_t = \Pr_{i \sim D_t} [h_t(\mathbf{x}_i) \neq y_i]$
5. $\alpha_t \leftarrow \frac{1}{2} \log \frac{1-\epsilon_t}{\epsilon_t}$
6. $Z_t \leftarrow 2 [\epsilon_t(1 - \epsilon_t)]^{\frac{1}{2}}$ (normalization factor)
7. **for** $i \leftarrow 1$ **to** m **do**
8. $D_{t+1}(i) \leftarrow \frac{D_t(i) \exp(-\alpha_t y_t h_t(\mathbf{x}_i))}{Z_t}$
9. $f_t \leftarrow \sum_{s=1}^t \alpha_s h_s$
10. **return** $h = \text{sign}(f_T)$

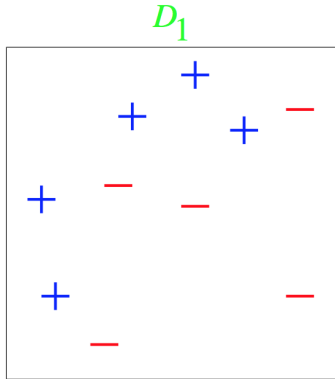
ADABOOST: COMMENTS

- Distribution D_t over training sample:
 - originally uniform
 - at each round, the weight of a misclassified example is increased
 - observation: $D_{t+1}(i) = \frac{e^{-y_i f_t(\mathbf{x}_i)}}{m \prod_{s=1}^t Z_s}$, since

$$\begin{aligned}
 D_{t+1}(i) &= \frac{D_t(i) e^{-\alpha_t y_i h_t(\mathbf{x}_i)}}{Z_t} \\
 &= \frac{D_{t-1}(i) e^{-\alpha_{t-1} y_i h_{t-1}(\mathbf{x}_i)} e^{-\alpha_t y_i h_t(\mathbf{x}_i)}}{Z_{t-1} Z_t} \\
 &= \frac{1}{m} \frac{e^{-y_i \sum_{s=1}^t \alpha_s h_s(\mathbf{x}_i)}}{\prod_{s=1}^t Z_s}
 \end{aligned}$$

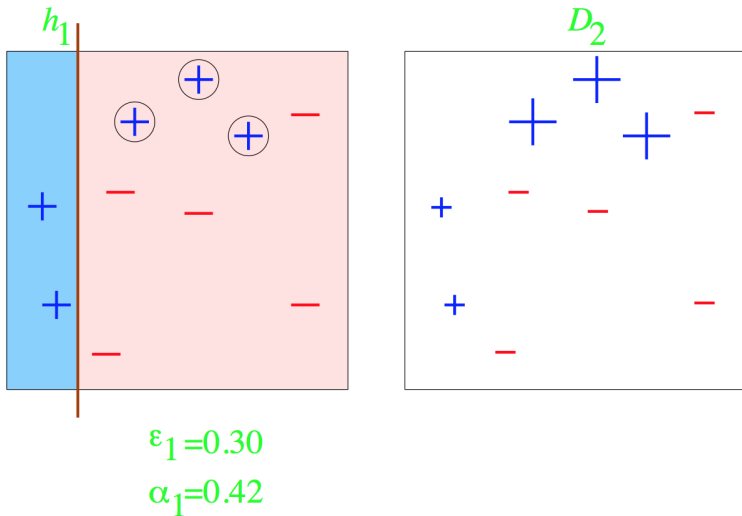
- Weight assigned to base classifier h_t : α_t directly depends on the accuracy of h_t at round t

TOY EXAMPLE

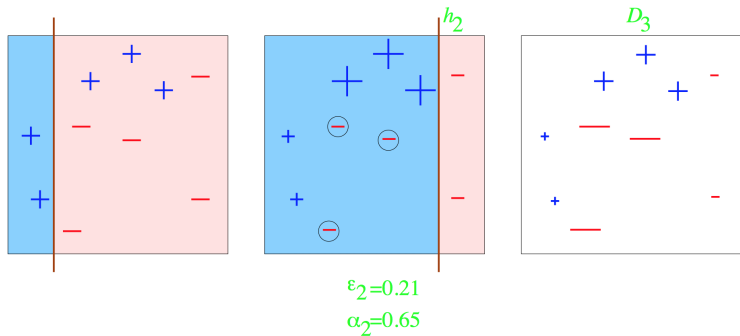


Weak classifiers = vertical or horizontal half-planes

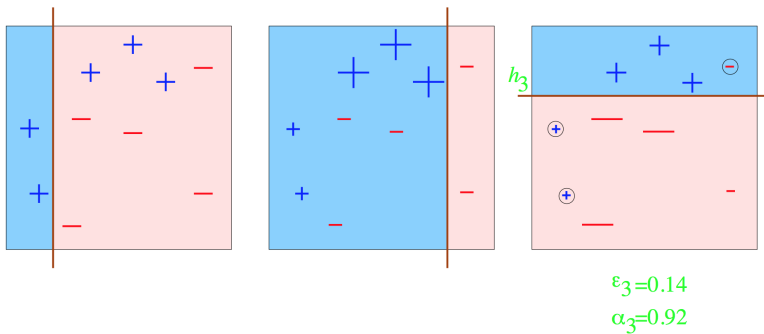
TOY EXAMPLE: ROUND 1



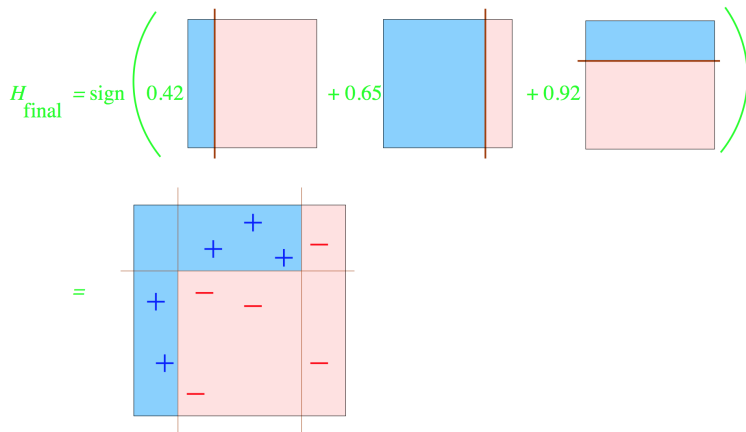
TOY EXAMPLE: ROUND 2

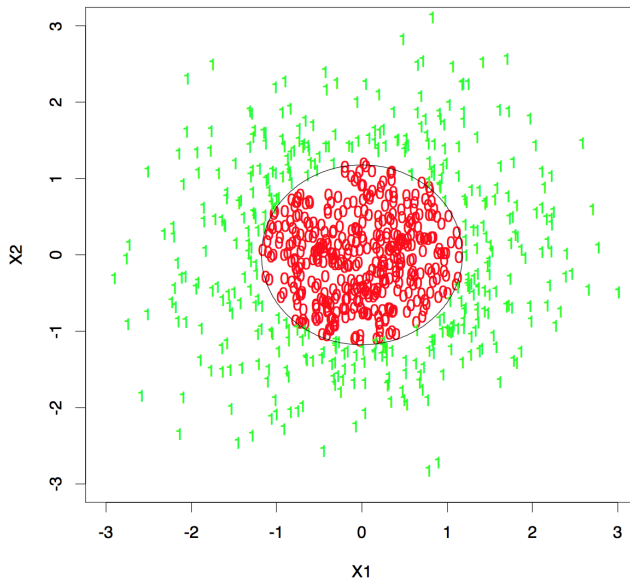


TOY EXAMPLE: ROUND 3

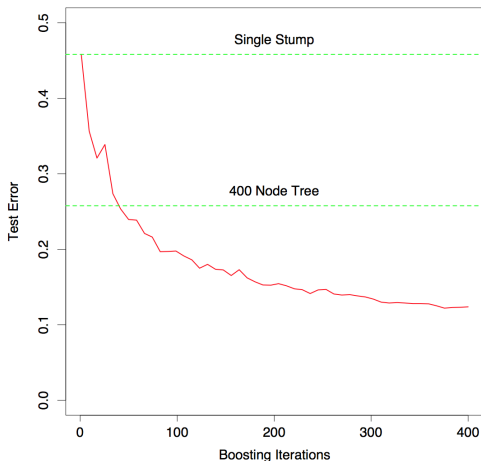


TOY EXAMPLE: FINAL CLASSIFIER



EXAMPLE: “SPHERE” IN \mathbb{R}^{10} 

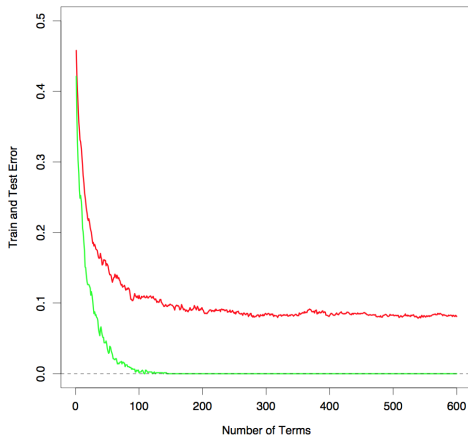
BOOSTING STUMPS



“Sphere” in \mathbb{R}^{10} : A stump is a two-node tree, after a single split.
Boosting stumps works remarkably well on this problem

TRAINING & TEST ERROR

Stumps



“Sphere” in \mathbb{R}^{10} : Boosting drives the training error to zero. Further iterations continue to improve test error in many examples

BOOSTING NOISY PROBLEMS I

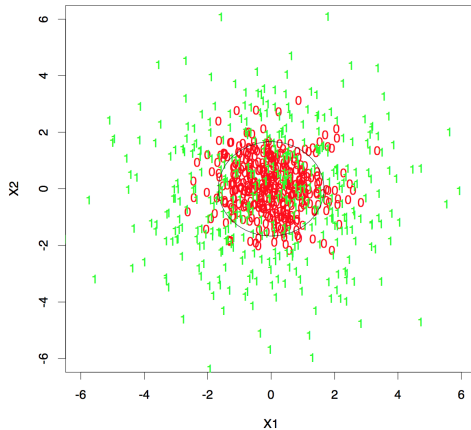
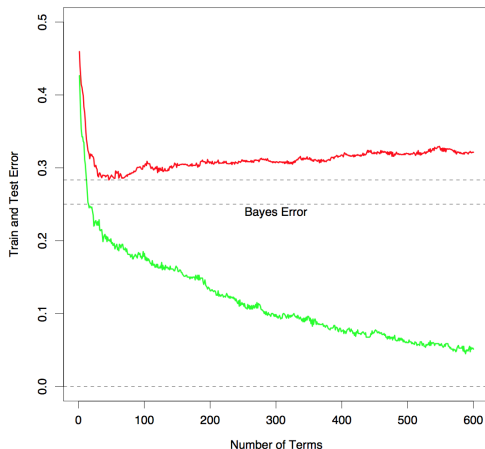


FIGURE : “Gaussians” in \mathbb{R}^{10} . Bayes error is 25%

BOOSTING NOISY PROBLEMS II

Stumps



“Gaussians” in \mathbb{R}^{10} . Bayes error is 25%. Here the test error does increase, but quite slowly

BOUND ON EMPIRICAL ERROR

- **Theorem:** The empirical error of the classifier output by AdaBoost verifies:

$$\hat{R}(h) \leq \exp \left[-2 \sum_{t=1}^T \left(\frac{1}{2} - \epsilon_t \right)^2 \right]$$

- if further for all $t \in [1, T]$, $\gamma \leq \left(\frac{1}{2} - \epsilon_t \right)$, then

$$\hat{R}(h) \leq \exp(-2\gamma^2 \cdot T)$$

- $\gamma > 0$ does not need to be known in advance:
adaptive boosting

- **Proof:** Since, as we saw, $D_{t+1}(i) = \frac{e^{-y_i f_t(\mathbf{x}_i)}}{m \prod_{s=1}^t Z_s}$,

$$\begin{aligned}\hat{R}(h) &= \frac{1}{m} \sum_{i=1}^m 1_{y_i f_T(\mathbf{x}_i) \leq 0} \leq \frac{1}{m} \sum_{i=1}^m \exp(-y_i f_T(\mathbf{x}_i)) \\ &\leq \frac{1}{m} \sum_{i=1}^m \left[m \prod_{t=1}^T Z_t \right] D_{T+1}(i) = \prod_{t=1}^T Z_t\end{aligned}$$

- Now, since Z_t is a normalization factor,

$$\begin{aligned}Z_t &= \sum_{i=1}^m D_t(i) e^{-\alpha_t y_i h_t(\mathbf{x}_i)} \\ &= \sum_{i: y_i h_t(\mathbf{x}_i) \geq 0} D_t(i) e^{-\alpha_t} + \sum_{i: y_i h_t(\mathbf{x}_i) < 0} D_t(i) e^{+\alpha_t} \\ &= (1 - \epsilon_t) e^{-\alpha_t} + \epsilon_t e^{\alpha_t} \\ &= (1 - \epsilon_t) \sqrt{\frac{\epsilon_t}{1 - \epsilon_t}} + \epsilon_t \sqrt{\frac{1 - \epsilon_t}{\epsilon_t}} = 2\sqrt{\epsilon_t(1 - \epsilon_t)}\end{aligned}$$

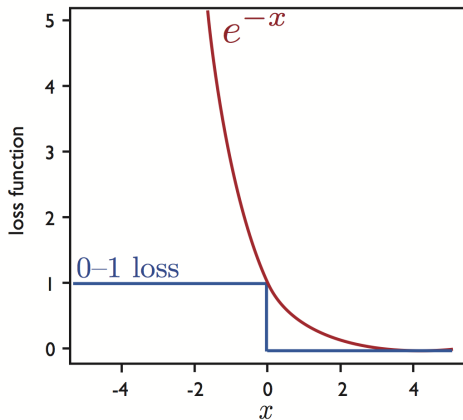
- Thus

$$\begin{aligned}\prod_{t=1}^T Z_t &= \prod_{t=1}^T 2\sqrt{\epsilon_t(1-\epsilon_t)} = \prod_{t=1}^T \sqrt{1-4\left(\frac{1}{2}-\epsilon_t\right)^2} \\ &\leq \prod_{t=1}^T \exp\left[-2\left(\frac{1}{2}-\epsilon_t\right)^2\right] = \exp\left[-2\sum_{t=1}^T \left(\frac{1}{2}-\epsilon_t\right)^2\right]\end{aligned}$$

- Comments:
 - α_t is a minimizer of $\alpha \rightarrow (1-\epsilon_t)e^{-\alpha} + \epsilon_t e^{\alpha}$
 - since $(1-\epsilon_t)e^{-\alpha_t} = \epsilon_t e^{\alpha_t}$, at each round Ada Boost assigns the same probability mass to correctly classified and misclassified instances

ADABOOST = COORDINATE DESCENT

- Objective Function: convex and differentiable



- **Direction** unit vector \mathbf{e}_k with best directional derivative

$$F'(\bar{\alpha}_{t-1}, \mathbf{e}_k) = \lim_{\eta \rightarrow 0} \frac{F(\bar{\alpha}_{t-1} + \eta \mathbf{e}_k) - F(\bar{\alpha}_{t-1})}{\eta}$$

- Since $F(\bar{\alpha}_{t-1} + \eta \mathbf{e}_k) = \sum_{i=1}^m e^{-y_i \sum_{j=1}^K \bar{\alpha}_{t-1,j} h_j(\mathbf{x}_i) - \eta y_i h_k(\mathbf{x}_i)}$,

$$\begin{aligned} F'(\bar{\alpha}_{t-1}, \mathbf{e}_k) &= -\frac{1}{m} \sum_{i=1}^m y_i h_k(\mathbf{x}_i) e^{-y_i \sum_{j=1}^K \bar{\alpha}_{t-1,j} h_j(\mathbf{x}_i)} \\ &= -\frac{1}{m} \sum_{i=1}^m y_i h_k(\mathbf{x}_i) \bar{D}_t(i) \bar{Z}_t \\ &= -\left[\sum_{i=1}^m \bar{D}_t(i) 1_{y_i h_k(\mathbf{x}_i)=+1} - \sum_{i=1}^m \bar{D}_t(i) 1_{y_i h_k(\mathbf{x}_i)=-1} \right] \frac{\bar{Z}_t}{m} \\ &= -[(1 - \bar{\epsilon}_{t,k}) - \bar{\epsilon}_{t,k}] \frac{\bar{Z}_t}{m} = [2\bar{\epsilon}_{t,k} - 1] \frac{\bar{Z}_t}{m} \end{aligned}$$

Here $[2\bar{\epsilon}_{t,k} - 1]$ is a direction corresponding to the base classifier with the smallest error

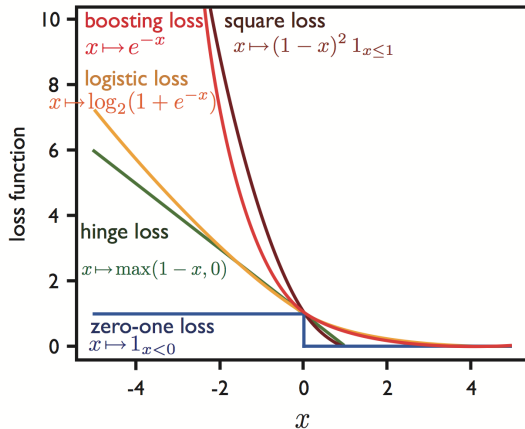
- **Step size:** η is chosen to minimize $F(\bar{\alpha}_{t-1} + \eta \mathbf{e}_k)$

$$\begin{aligned}
 \frac{dF(\bar{\alpha}_{t-1} + \eta \mathbf{e}_k)}{d\eta} = 0 &\Leftrightarrow - \sum_{i=1}^m y_i h_k(\mathbf{x}_i) e^{-y_i \sum_{j=1}^K \bar{\alpha}_{t-1,j} h_j(\mathbf{x}_i)} e^{-\eta y_i h_k(\mathbf{x}_i)} \\
 &\Leftrightarrow - \sum_{i=1}^m y_i h_k(\mathbf{x}_i) \bar{D}_t(i) \bar{Z}_t e^{-\eta y_i h_k(\mathbf{x}_i)} = 0 \\
 &\Leftrightarrow - \sum_{i=1}^m y_i h_k(\mathbf{x}_i) \bar{D}_t(i) e^{-\eta y_i h_k(\mathbf{x}_i)} = 0 \\
 &\Leftrightarrow - \left[(1 - \bar{\epsilon}_{t,k}) e^{-\eta} - \bar{\epsilon}_{t,k} e^{\eta} \right] = 0 \\
 &\Leftrightarrow \eta = \frac{1}{2} \log \frac{1 - \bar{\epsilon}_{t,k}}{\bar{\epsilon}_{t,k}}
 \end{aligned}$$

Thus, step size matches base classifier weight of AdaBoost

ALTERNATIVE LOSS FUNCTIONS

- Examples of several convex upper bounds on the zero-one loss



STANDARD USE IN PRACTICE

- **Base Learners:** decision trees, quite often just decision stumps (trees of depth one)
- Boosting stumps
 - data in \mathbb{R}^N , e.g. $N = 2$ (height(x), weight(x))
 - associate a stump to each component
 - pre-sort each component: $O(Nm \log m)$
 - at each round, find best component and threshold
 - total complexity: $O((m \log m)N + mNT)$
 - stumps are not weak learners (XOR problem)

OVERFITTING?

- Assume that $\text{VCdim} = d$ and for a fixed T , define

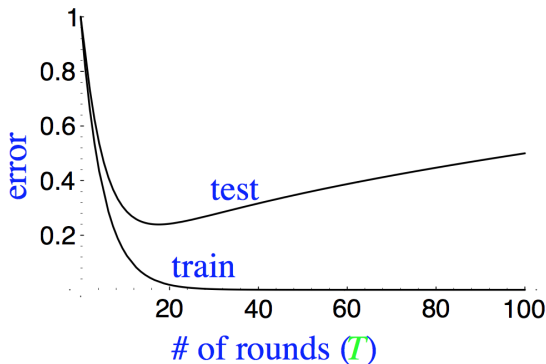
$$\mathcal{F}_T = \left\{ \text{sign} \left(\sum_{t=1}^T \alpha_t h_t - b \right) : \alpha_t, b \in \mathbb{R}, h_t \in H \right\}$$

- \mathcal{F}_T can form a very rich family of classifiers. It can be shown (Freund & Shapire, 1997) that

$$\text{VCdim}(\mathcal{F}_T) \leq 2(d+1)(T+1) \log_2((T+1)e)$$

- This suggests that AdaBoost could overfit for large values of T , and that is in fact observed in some cases, but in various others it is not!

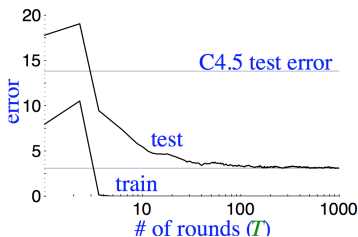
HOW WILL TEST ERROR BEHAVE? (A FIRST GUESS)



Expect:

- training error to continue to drop (or reach zero)
- test error to increase when h_{final} becomes “too complex”
 - “Occams razor”
 - overfitting: hard to know when to stop training

EMPIRICAL OBSERVATIONS



(boosting C4.5 on
"letter" dataset)

Expect:

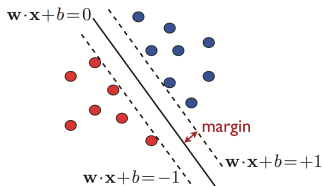
- test error does not increase, even after 1000 rounds
— (total size > 2,000,000 nodes)
- test error continues to drop even after training error is zero!

	# rounds		
	5	100	1000
train error	0.0	0.0	0.0
test error	8.4	3.3	3.1

- Occams razor wrongly predicts "simpler" rule is better

- 1 MOTIVATION: CLASSIFICATION PROBLEM
- 2 BAGGING
- 3 BOOSTING
- 4 MARGIN IN CASE OF BOOSTING**

MARGIN IN A LINEAR SEPARABLE CASE



- Linear separating hyperplane $\mathbf{w} \cdot \mathbf{x} + b = 0$
- The distance of any point $\mathbf{x}_0 \in \mathbb{R}^N$ to a hyperplane is given by

$$\frac{|\mathbf{w} \cdot \mathbf{x}_0 + b|}{\|\mathbf{w}\|}$$

- For a hyperplane that does not pass through any sample point, we can scale \mathbf{w} and b s.t. $\min_{(\mathbf{x}, y) \in S} |\mathbf{w} \cdot \mathbf{x} + b| = 1$
- The margin is given by

$$\rho = \min_{(\mathbf{x}, y) \in S} \frac{|\mathbf{w} \cdot \mathbf{x} + b|}{\|\mathbf{w}\|} = \frac{1}{\|\mathbf{w}\|}$$

SVM OPTIMIZATION PROBLEM

- **Constrained Optimization** in a separable case:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{s.t.} \quad y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, \quad i \in [1, m]$$

- **Dual Constrained Optimization:**

$$\max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$

$$\text{s.t.} \quad 0 \leq \alpha_i \leq C \quad \text{and} \quad \sum_{i=1}^m \alpha_i y_i = 0, \quad i \in [1, m]$$

- **Solution** $h(\mathbf{x}) = \text{sgn}(\sum_{i=1}^m \alpha_i y_i (\mathbf{x}_i \cdot \mathbf{x}) + b)$, with

$$b = y_i - \sum_{j=1}^m \alpha_j y_j (\mathbf{x}_j \cdot \mathbf{x}_i)$$

for any SV \mathbf{x}_i with $0 < \alpha_i < C$

OPTIMAL HYPERPLANE

- Multiplying both sides of $b = y_i - \sum_{j=1}^m \alpha_j y_j (\mathbf{x}_j \cdot \mathbf{x}_i)$ by $\alpha_i y_i$ we get that

$$\sum_{i=1}^m \alpha_i y_i b = \sum_{i=1}^m \alpha_i y_i^2 - \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$

- From KKT conditions we get that $\mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i$
- Thus, since $y_i^2 = 1$, we obtain

$$0 = \sum_{i=1}^m \alpha_i - \|\mathbf{w}\|^2,$$

i.e.

$$\rho^2 = \frac{1}{\|\mathbf{w}\|^2} = \frac{1}{\sum_{i=1}^m \alpha_i} = \frac{1}{\|\boldsymbol{\alpha}\|_1}$$

GEOMETRIC MARGIN: GENERAL LINEAR CASE

- The geometric margin $\rho(\mathbf{x})$ of a point \mathbf{x} with label y w.r.t. a linear classifier $h : \mathbf{x} \rightarrow \mathbf{w} \cdot \mathbf{x} + b$ is its distance to the hyperplane $\mathbf{w} \cdot \mathbf{x} + b = 0$:

$$\rho(\mathbf{x}) = \frac{y(\mathbf{w} \cdot \mathbf{x} + b)}{\|\mathbf{w}\|}$$

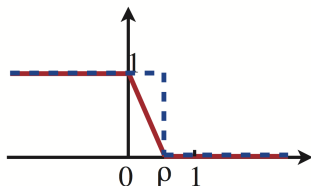
- The margin of a linear classifier h for a sample $S = (\mathbf{x}_1, \dots, \mathbf{x}_m)$ is the minimum margin over the points in the sample

$$\rho = \min_{1 \leq i \leq m} \frac{y_i(\mathbf{w} \cdot \mathbf{x}_i + b)}{\|\mathbf{w}\|}$$

MARGIN LOSS FUNCTION

- Here, the training data is not assumed to be separable. The quantity $\rho > 0$ should thus be interpreted as the margin we wish to achieve
- For any $\rho > 0$, the ρ -margin loss
 $L_\rho(y, y') = \Phi_\rho(yy') : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}_+$ with

$$\Phi_\rho(t) = \begin{cases} 0, & \text{if } \rho \leq t, \\ 1 - \frac{t}{\rho}, & \text{if } 0 \leq t \leq \rho \\ 1, & \text{if } t \leq 0 \end{cases}$$



EMPIRICAL MARGIN LOSS

- Given a sample $S = (\mathbf{x}_1, \dots, \mathbf{x}_m)$ and a hypothesis h , the empirical margin loss is defined by

$$\hat{R}_\rho(h) = \frac{1}{m} \sum_{i=1}^m \Phi_\rho(y_i h(\mathbf{x}_i))$$

- Note that for any $i \in [1, m]$, $\Phi_\rho(y_i h(\mathbf{x}_i)) \leq 1_{y_i h(\mathbf{x}_i) \leq \rho}$. Thus, the empirical margin loss can be upper-bounded as follows:

$$\hat{R}_\rho(h) \leq \frac{1}{m} \sum_{i=1}^m 1_{y_i h(\mathbf{x}_i) \leq \rho}$$

- This upper bound is the fraction of the points in the training sample S that have been misclassified or classified with confidence less than ρ

MARGIN FOR BOOSTING I

- Kernel classifier:

$$f(\mathbf{x}) = \mathbf{w} \cdot \Phi(\mathbf{x})$$

- Let us denote by $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_T)^T$,
 $\mathbf{h}(\mathbf{x}) = (h_1(\mathbf{x}), \dots, h_T(\mathbf{x}))$
- Boosting: linear combination of base classifiers

$$f(\mathbf{x}) = \sum_{t=1}^T \alpha_t h_t = \boldsymbol{\alpha} \cdot \mathbf{h}(\mathbf{x})$$

MARGIN FOR BOOSTING II

- The L_1 -margin $\rho(\mathbf{x})$ of a point $\mathbf{x} \in X$ with label $y \in \{-1, +1\}$ for a linear combination of base classifiers $f(\mathbf{x}) = \sum_{t=1}^T \alpha_t h_t$ with $\alpha \neq 0$ and $h_t \in H$ for all $t \in [1, T]$ is defined as

$$\rho(\mathbf{x}) = \frac{yg(\mathbf{x})}{\sum_{t=1}^T |\alpha_t|} = \frac{y \sum_{t=1}^T \alpha_t h_t(\mathbf{x})}{\|\alpha\|_1} = y \frac{\alpha \cdot \mathbf{h}(\mathbf{x})}{\|\alpha\|_1}$$

- The L_1 -margin of a linear combination classifier f with respect to a sample $S = (\mathbf{x}_1, \dots, \mathbf{x}_m)$ is the minimum margin of the points within the sample:

$$\rho = \min_{i \in [1, m]} y_i \frac{\alpha \cdot \mathbf{h}(\mathbf{x}_i)}{\|\alpha\|_1}$$

MARGIN FOR BOOSTING III

- If $\alpha_t \geq 0$, $\rho(\mathbf{x}) \in [-1, 1]$. Thus $|\rho(\mathbf{x})|$ can be interpreted as the confidence of the classifier $f(\mathbf{x})$ in that label
- It is well-known that $\frac{|\alpha \cdot \mathbf{x}|}{\|\alpha\|_p}$ is the L_q distance of \mathbf{x} to the hyperplane of equation $\alpha \cdot \mathbf{x} = 0$, where $1/p + 1/q = 1$, $p, q \geq 1$
- Thus

$$\rho_1(\mathbf{x}) = \frac{|\alpha \cdot \mathbf{h}(\mathbf{x})|}{\|\alpha\|_1} \quad \text{and} \quad \rho_2(\mathbf{x}) = \frac{|\alpha \cdot \mathbf{h}(\mathbf{x})|}{\|\alpha\|_2}$$

measures $\|\cdot\|_\infty$ and $\|\cdot\|_2$ distances of the feature vector $\mathbf{h}(\mathbf{x})$ to the hyperplane $\alpha \cdot \mathbf{h} = 0$

MARGIN OPTIMIZATION

- By definition of the L_1 -margin, the maximum margin for a sample S is given by

$$\rho = \max_{\alpha} \min_{i \in [1, m]} y_i \frac{\alpha \cdot \mathbf{h}(\mathbf{x}_i)}{\|\alpha\|_1}$$

- The optimization problem (if we do not need to tune $\mathbf{h}(\mathbf{x})$):

$$\begin{aligned} & \max_{\alpha} \rho \\ & \text{s.t. : } y_i \frac{\alpha \cdot \mathbf{h}(\mathbf{x}_i)}{\|\alpha\|_1} \geq \rho, \quad \forall i \in [1, m] \end{aligned}$$

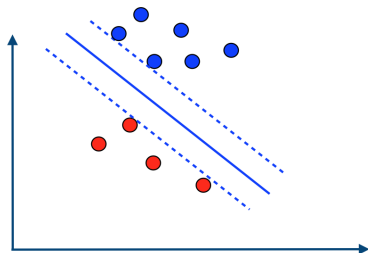
- Equivalent optimization problem:

$$\begin{aligned} & \max_{\alpha} \rho, \quad \text{s.t. : } y_i(\alpha \cdot \mathbf{h}(\mathbf{x}_i)) \geq \rho, \quad \forall i \in [1, m] \\ & \left(\sum_{t=1}^T \alpha_t = 1 \right) \quad \text{and} \quad (\alpha_t \geq 0, \quad \forall t \in [1, T]) \end{aligned}$$

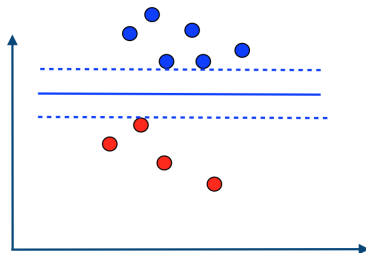
SVM vs. ADABOOST

	SVM	AdaBoost
features or base hypotheses	$\Phi(x) = \begin{bmatrix} \Phi_1(x) \\ \vdots \\ \Phi_N(x) \end{bmatrix}$	$\mathbf{h}(x) = \begin{bmatrix} h_1(x) \\ \vdots \\ h_N(x) \end{bmatrix}$
predictor	$x \mapsto \mathbf{w} \cdot \Phi(x)$	$x \mapsto \alpha \cdot \mathbf{h}(x)$
geom. margin	$\frac{ \mathbf{w} \cdot \Phi(x) }{\ \mathbf{w}\ _2} = d_2(\Phi(x), \text{hyperpl.})$	$\frac{ \alpha \cdot \mathbf{h}(x) }{\ \alpha\ _1} = d_\infty(\mathbf{h}(x), \text{hyperpl.})$
conf. margin	$y(\mathbf{w} \cdot \Phi(x))$	$y(\alpha \cdot \mathbf{h}(x))$
regularization	$\ \mathbf{w}\ _2$	$\ \alpha\ _1$ (L1-AB)

MAXIMUM MARGIN SOLUTIONS



Norm $\| \cdot \|_2$.



Norm $\| \cdot \|_\infty$.

MARGIN DISTRIBUTION

- **Theorem:** for any $\rho > 0$, the following holds:

$$\widehat{\Pr} \left[\frac{yf(\mathbf{x})}{\|\alpha\|_1} \leq \rho \right] \leq 2^T \prod_{t=1}^T \sqrt{\epsilon_t^{1-\rho} (1 - \epsilon_t)^{1+\rho}}$$

- **Proof:** Using the identity $D_{t+1}(i) = \frac{e^{-y_i f(\mathbf{x}_i)}}{m \prod_{t=1}^T Z_t}$

$$\begin{aligned} \frac{1}{m} \sum_{i=1}^m 1_{y_i f(\mathbf{x}_i) - \|\alpha\|_1 \rho \leq 0} &\leq \frac{1}{m} \sum_{i=1}^m \exp(-y_i f(\mathbf{x}_i) + \|\alpha\|_1 \rho) \\ &= \frac{1}{m} \sum_{i=1}^m e^{\|\alpha\|_1 \rho} \left[m \prod_{t=1}^T Z_t \right] D_{T+1}(i) \end{aligned}$$

$$= e^{\|\alpha\|_1 \rho} \prod_{t=1}^T Z_t = 2^T \prod_{t=1}^T \left[\sqrt{\frac{1 - \epsilon_t}{\epsilon_t}} \right]^\rho \sqrt{\epsilon_t (1 - \epsilon_t)}$$

MARGIN DISTRIBUTION: COMMENTS

- if for all $t \in [1, T]$, $\gamma \leq (\frac{1}{2} - \epsilon_t)$, then the upper bound can be bounded by

$$\widehat{\Pr} \left[\frac{yf(\mathbf{x})}{\|\boldsymbol{\alpha}\|_1} \leq \rho \right] \leq [(1 - 2\gamma)^{1-\rho}(1 + 2\gamma)^{1+\rho}]^{T/2}$$

For $\rho < \gamma$, $(1 - 2\gamma)^{1-\rho}(1 + 2\gamma)^{1+\rho} < 1$ and the bound decreases exponentially in T

- For the bound to be convergent: $\rho \gg O(1/\sqrt{m})$, thus $\gamma \gg O(1/\sqrt{m})$ is roughly the condition on the edge value

OUTLIERS

- AdaBoost assigns larger weights to harder examples
- Applications:
 - Detecting mislabeled examples
 - Dealing with noisy data: regularization based on the average weight assigned to a point (soft margin idea for boosting)