**Resampling for Classifier Design**

* Reusing or selecting data in order to improve classification
* Two most popular
    * Bagging (Breiman, 1994)
    * AdaBoost (Freund and Schapire, 1996)

The idea is to combine the results of multiple "weak" classifiers into a single "strong" classifier.

**The general idea:**

Repeat T times:
    1. Derive rough rule-of-thumb: weak classifier (performs slightly above chance)
    2. Select new sample, derive 2nd rule-of-thumb (weak classifier)
end

**Questions**

1. How to choose samples?
    a. Select multiple random samples?
    b. Concentrate only on the errors?

2. How to combine rules-of-thumb into a single accurate rule?

**More formally:**

Given: training data $(x_1, y_1), \ldots, (x_m, y_m)$, where $x_i \in \mathcal{X}$, $y_i \in \mathcal{Y} = \{-1, +1\}$

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

    1. Train *Weak Learner* on the training set.
       Let $h_t : \mathcal{X} \to \{-1, +1\}$ represent the classifier obtained after training.

    2. Modify the training set somehow

- The final hypothesis $H(x)$ is some combination of all the weak hypotheses:

$$H(x) = f(h(x)) \tag{1}$$

The question is how to modify the training set, and how to combine the weak classifiers.

**Bagging**

The simplest algorithm is called Bagging, used by Breiman 1994

## Algorithm:

Given $m$ training examples, repeat for $t = 1 \ldots T$:

- Select, at random *with replacement*, $m$ training examples.

- Train learning algorithm on selected examples to generate hypothesis $h_t$

Final hypothesis is simple vote: $H(x) = MAJ(h_1(x), \ldots, h_T(x))$.

**Bagging Pros and Cons:**

1. Bagging reduces variance
   a. Helps improve unstable classifiers: i.e., "small" changes in training data lead to significantly different classifiers and "large" changes in accuracy.
   b. no proof for this, however
2. Does not reduce bias

**Boosting:**

Two modifications

1. instead of a random sample of the training data, use a weighted sample to focus learning on most difficult examples.
2. instead of combining classifiers with equal vote, use a weighted vote.

Several previous methods (Schapire, 1990; Freund, 1995) were effective, but had limitations. In the class, we consider the one proposed by Freund and Schapire 1996 called **Adaboost.**

## AdaBoost (Freund and Schapire, 1996)

- Initialize distribution over the training set $D_1(i) = 1/m$

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

  1. Train *Weak Learner* using distribution $D_t$.

  2. Choose a weight (or confidence value) $\alpha_t \in \mathbf{R}$.

  3. Update the distribution over the training set:

  $$D_{t+1}(i) = \frac{D_t(i)e^{-\alpha_t y_i h_t(x_i)}}{Z_t} \tag{2}$$

  Where $Z_t$ is a normalization factor chosen so that $D_{t+1}$ will be a distribution

- Final vote $H(x)$ is a weighted sum:

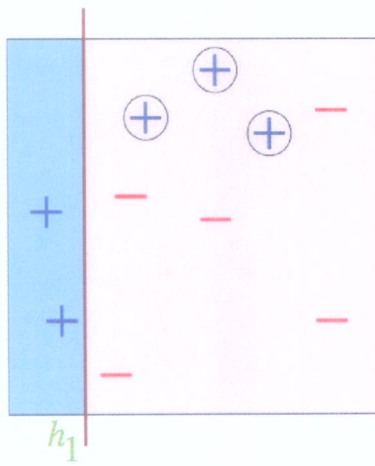$$H(x) = \text{sign}(f(x)) = \text{sign}\left(\sum_{t-1}^{T} \alpha_t h_t(x)\right) \tag{3}$$

**How to select alpha?**

To decide how to pick the alphas, we have to understand what the relationship is between the distribution, the alpha_t, and the training error.
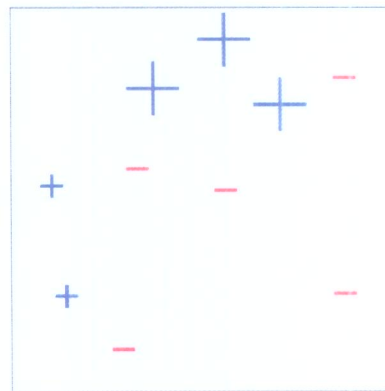
**Toy Example**

$D_1$

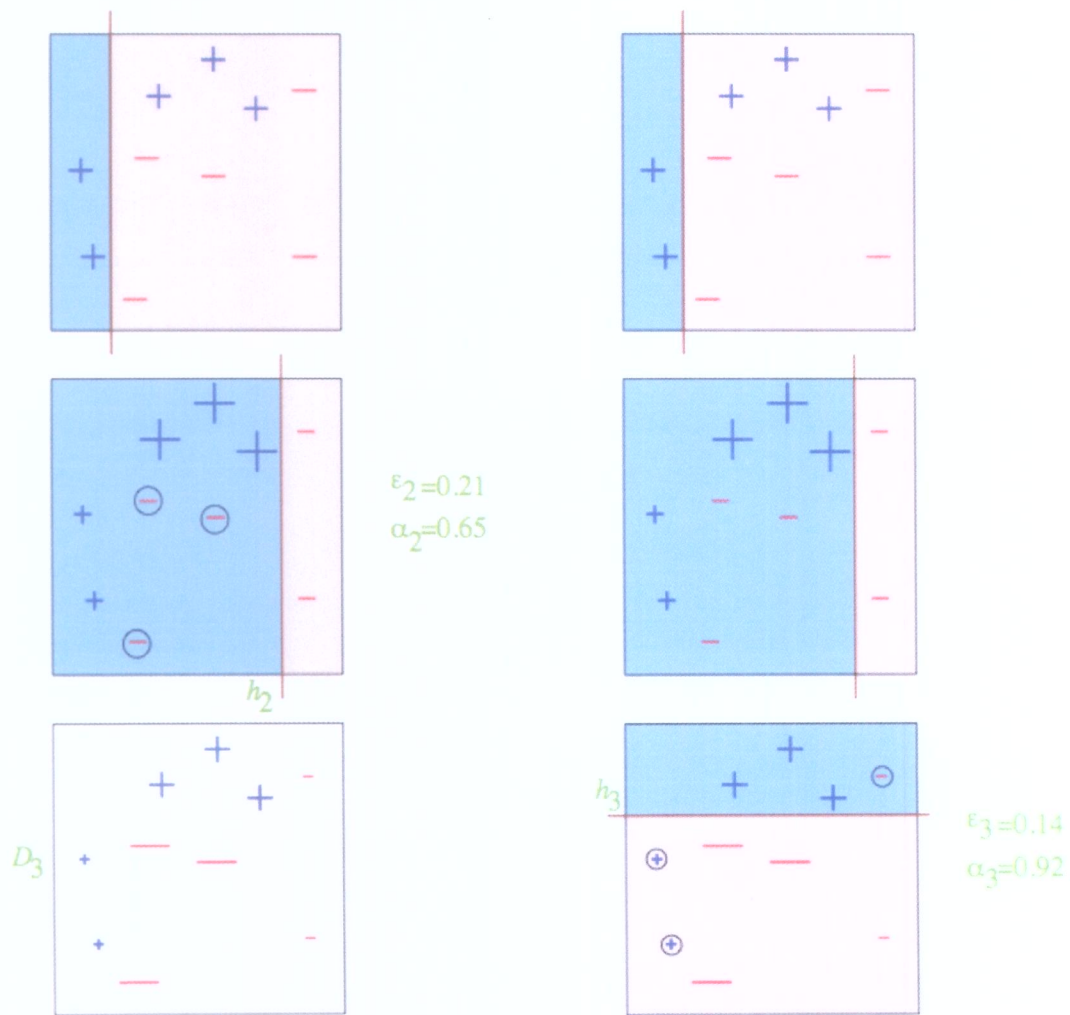**Round 1**

$h_1$

$\varepsilon_1 = 0.30$
$\alpha_1 = 0.42$

$D_2$

**Round 2 and 3, respectively**



$\varepsilon_2 = 0.21$
$\alpha_2 = 0.65$

$h_2$

$D_3$

$h_3$

$\varepsilon_3 = 0.14$
$\alpha_3 = 0.92$

**Final Classification:**

$H_{final}$

$$= \text{sign}\left( 0.42 \quad \boxed{} \quad + 0.65 \quad \boxed{} \quad + 0.92 \quad \boxed{} \right)$$

$$=$$

## Maximising margins in AdaBoost

$$P_{(x,y)\sim S}[yf(x) \leq \theta] \leq 2^T \prod_{t=1}^{T} \sqrt{\epsilon_t^{1-\theta}(1-\epsilon_t)^{1+\theta}} \qquad \text{where } f(x) = \frac{\vec{\alpha} \cdot \vec{h}(x)}{\|\vec{\alpha}\|_1}$$

Choosing $h_t(x)$ with minimal $\epsilon_t$ in each step one minimises the margin

Margin in SVM use the $L_2$ norm instead: $(\vec{\alpha} \cdot \vec{h}(x))/\|\vec{\alpha}\|_2$

## Upper bounds based on margin

With probability $1 - \delta$ over the random choice of the training set $S$

$$P_{(x,y)\sim\mathcal{D}}[yf(x) \leq 0] \leq P_{(x,y)\sim S}[yf(x) \leq \theta] + \mathcal{O}\left(\frac{1}{\sqrt{m}}\left(\frac{d\log^2(m/d)}{\theta^2} + \log(1/\delta)\right)^{1/2}\right)$$

where $\mathcal{D}$ is a distribution over $\mathcal{X} \times \{+1, -1\}$, and $d$ is pseudodimension of $\mathcal{H}$.

**Problem:** The upper bound is very loose. In practice AdaBoost works much better.

---

## Freund & Schapire 1995

Discrete ($h : \mathcal{X} \to \{0, 1\}$)

Multiclass AdaBoost.M1 ($h : \mathcal{X} \to \{0, 1, ..., k\}$)

Multiclass AdaBoost.M2 ($h : \mathcal{X} \to [0, 1]^k$)

Real valued AdaBoost.R ($Y = [0, 1]$, $h : \mathcal{X} \to [0, 1]$)

## Schapire & Singer 1999

Confidence rated prediction ($h : \mathcal{X} \to R$, two-class)

Multilabel AdaBoost.MR, AdaBoost.MH (different formulation of minimised loss)

## Oza 2001

Online AdaBoost

Many other modifications since then: cascaded AB, WaldBoost, probabilistic boosting tree, ...

---

Given: $(x_1, y_1), \ldots, (x_m, y_m); x_i \in \mathcal{X}, y_i \in \{-1, +1\}$
Initialise weights $D_1(i) = 1/m$
For $t = 1, ..., T$:

Find $h_t = \arg\min_{h_j \in \mathcal{H}} \epsilon_j = \sum_{i=1}^{m} D_t(i) [\![y_i \neq h_j(x_i)]\!]$
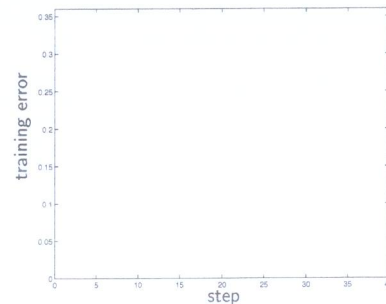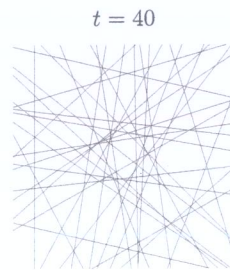
If $\epsilon_t \geq 1/2$ then stop

Set $\alpha_t = \frac{1}{2}\log(\frac{1-\epsilon_t}{\epsilon_t})$

Update

$$D_{t+1}(i) = \frac{D_t(i)exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

Output the final classifier:

$$H(x) = sign\left(\sum_{t=1}^{T} \alpha_t h_t(x)\right)$$

$t = 40$



---

|  **Offline** | **Online** |
| --- | --- |
| Given: | Given: |
| Set of labeled training samples $\mathcal{X} = \{(x_1, y_1), ..., (x_m, y_m) | y = \pm 1\}$ | *One* labeled training sample $(x, y) | y = \pm 1$ |
| Weight distribution over $\mathcal{X}$ $D_0 = 1/m$ | Strong classifier to update |
|  | Initial importance $\lambda = 1$ |
| For $t = 1, \ldots, T$ | For $t = 1, \ldots, T$ |
| Train a weak classifier using samples and weight distribution | Update the weak classifier using the sample and the importance |
| $h_t(x) = \mathcal{L}(\mathcal{X}, D_{t-1})$ | $h_t(x) = \mathcal{L}(h_t, (x, y), \lambda)$ |
| Calculate error $\epsilon_t$ | Update error estimation $\epsilon_t$ |
| Calculate coeficient $\alpha_t$ from $\epsilon_t$ | Update weight $\alpha_t$ based on $\epsilon_t$ |
| Update weight distribution $D_t$ | Update importance weight $\lambda$ |
| Output: | Output: |
| $F(x) = sign(\sum_{t=1}^{T} \alpha_t h_t(x))$ | $F(x) = sign(\sum_{t=1}^{T} \alpha_t h_t(x))$ |

145

Converges to offline results given the same training set and the number of iterations $N \to \infty$

N. Oza and S. Russell. Online Bagging and Boosting.
Artificial Inteligence and Statistics, 2001.

Y. Freund, R.E. Schapire. **A Decision-theoretic Generalization of On-line Learning and an Application to Boosting**. Journal of Computer and System Sciences. 1997

R.E. Schapire, Y. Freund, P. Bartlett, W.S. Lee. **Boosting the Margin: A New Explanation for the Effectiveness of Voting Methods**. The Annals of Statistics, 1998

R.E. Schapire, Y. Singer. **Improved Boosting Algorithms Using Confidence-rated Predictions**. Machine Learning. 1999

J. Friedman, T. Hastie, R. Tibshirani. **Additive Logistic Regression: a Statistical View of Boosting**. Technical report. 1998

N.C. Oza. **Online Ensemble Learning**. PhD thesis. 2001

`http://www.boosting.org`

**Advantages**

Very simple to implement

General learning scheme - can be used for various learning tasks

Feature selection on very large sets of features

Good generalisation

Seems not to overfit in practice (probably due to margin maximisation)

**Disadvantages**

Suboptimal solution (greedy learning)

**Motivation**

AdaBoost with trees is the best off-the-shelf classifier in the world. (Breiman 1998)

*That's his opinion. Normally, there is a best classif. for each case. So, the real thing is "depends".*

**Outline:**

AdaBoost algorithm

- How it works?
- Why it works?

Online AdaBoost and other variants

Given: $(x_1, y_1), \ldots, (x_m, y_m); x_i \in \mathcal{X}, y_i \in \{-1, +1\}$

Initialise weights $D_1(i) = 1/m$

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

$t = 40$

Find $h_t = \arg \min_{h_j \in \mathcal{H}} \epsilon_j = \sum_{i=1}^{m} D_t(i) [\![ y_i \neq h_j(x_i) ]\!]$

If $\epsilon_t \geq 1/2$ then stop

Set $\alpha_t = \frac{1}{2} \log(\frac{1-\epsilon_t}{\epsilon_t})$
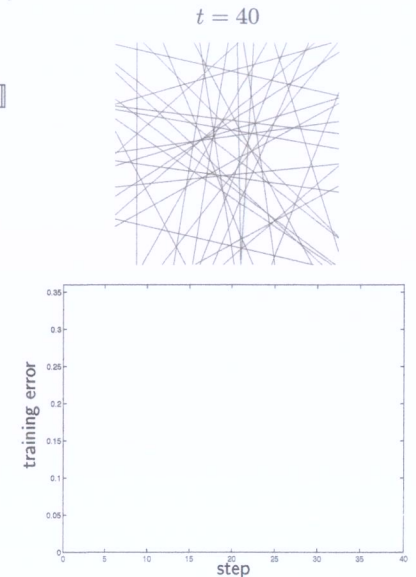
Update

$$D_{t+1}(i) = \frac{D_t(i) exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

where $Z_t$ is normalisation factor

Output the final classifier:

$$H(x) = sign\left(\sum_{t=1}^{T} \alpha_t h_t(x)\right)$$

AdaBoost is an algorithm for constructing a "strong" classifier as linear combination

$$f(x) = \sum_{t=1}^{T} \alpha_t h_t(x)$$

*simple $\neq$ weak*
*complex     strong*

of "simple" "weak" classifiers $h_t(x): \mathcal{X} \to \{-1, +1\}$.

**Terminology**

$h_t(x) \ldots$ "weak" or basis classifier, hypothesis, "feature"

$H(x) = sign(f(x)) \ldots$ "strong" or final classifier/hypothesis

**Interesting properties**

AB is capable reducing both bias (e.g. stumps) and variance (e.g. trees) of the weak classifiers

AB has good generalisation properties (maximises margin)

AB output converges to the logarithm of likelihood ratio

AB can be seen as a feature selector with a principled strategy (minimisation of upper bound on empirical error)

*E.g Viola & Jones face detector*

AB is close to sequential decision making (it produces a sequence of gradually more complex classifiers)

**Effect on the training set**

$$D_{t+1}(i) = \frac{D_t(i) exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

$$exp(-\alpha_t y_i h_t(x_i)) \begin{cases} < 1, & y_i = h_t(x_i) \\ > 1, & y_i \neq h_t(x_i) \end{cases}$$

$\Rightarrow$ Increase (decrease) weight of wrongly (correctly) classified examples

$\Rightarrow$ The weight is the upper bound on the error of a given example

$\Rightarrow$ All information about previously selected "features" is captured in $D_t$

143

**Theorem:** The following upper bound holds on the training error of $H$
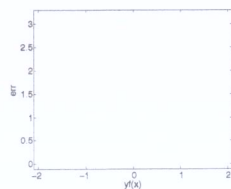
$$\frac{1}{m}|\{i : H(x_i) \neq y_i\}| \leq \prod_{t=1}^{T} Z_t$$

**Proof:** By unravelling the update rule

$$D_{T+1}(i) = \frac{D_t(i)exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

$$= \frac{exp(-\sum_t \alpha_t y_i h_t(x_i))}{m \prod_t Z_t} = \frac{exp(-y_i f(x_i))}{m \prod_t Z_t}$$

If $H(x_i) \neq y_i$ then $y_i f(x_i) \leq 0$ implying that $exp(-y_i f(x_i)) > 1$, thus

$$[\![H(x_i) \neq y_i]\!] \leq exp(-y_i f(x_i))$$

$$\frac{1}{m}\sum_i [\![H(x_i) \neq y_i]\!] \leq \frac{1}{m}\sum_i exp(-y_i f(x_i))$$

$$= \sum_i (\prod_t Z_t) D_{T+1}(i) = \prod_t Z_t$$

Instead of minimising the training error, its upper bound can be minimised

This can be done by minimising $Z_t$ in each training round by:

- Choosing optimal $h_t$, and
- Finding optimal $\alpha_t$

AdaBoost can be proved to maximise margin

AdaBoost iteratively fits an additive logistic regression model

We attempt to minimise $Z_t = \sum_i D_t(i) exp(-\alpha_t y_i h_t(x_i))$:

$$\frac{dZ}{d\alpha} = -\sum_{i=1}^{m} D(i) y_i h(x_i) e^{-y_i \alpha_i h(x_i)} = 0$$

$$-\sum_{i:y_i=h(x_i)} D(i)e^{-\alpha} + \sum_{i:y_i \neq h(x_i)} D(i)e^{\alpha} = 0$$

$$-e^{-\alpha}(1-\epsilon) + e^{\alpha}\epsilon = 0$$

$$\alpha = \frac{1}{2}\log\frac{1-\epsilon}{\epsilon} \quad \left(\text{proof in the verse}\right)$$

⇒ The minimisator of the upper bound is

$$\alpha = \frac{1}{2} \cdot \log\left(\frac{1 - \mathcal{E}_t}{\mathcal{E}_t}\right) \qquad \mathcal{E}_t = \text{Error training}$$
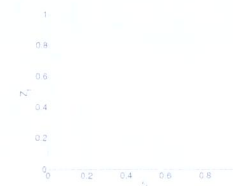
**Weak classifier examples**

> Decision tree (or stump), Perceptron − $\mathcal{H}$ *infinite*
>
> Selecting the best one from given *finite* set $\mathcal{H}$

**Justification of the weighted error minimisation**

Having $\alpha_t = \frac{1}{2}\log\frac{1-\epsilon_t}{\epsilon_t}$

$$Z_t = \sum_{i=1}^{m} D_t(i) e^{-y_i \alpha_i h_t(x_i)}$$

$$= \sum_{i:y_i=h_t(x_i)} D_t(i) e^{-\alpha_t} + \sum_{i:y_i \neq h_t(x_i)} D_t(i) e^{\alpha_t}$$

$$= (1 - \epsilon_t) e^{-\alpha_t} + \epsilon_t e^{\alpha_t}$$

$$= 2\sqrt{\epsilon_t(1 - \epsilon_t)}$$

⇒ $Z_t$ is minimised by selecting $h_t$ with minimal weighted error $\epsilon_t$

144

$$-\sum_{i\,:\,y_i = h_t(x_i)} D(i) \cdot \ell^{-\alpha} + \sum_{i\,:\,y_i \neq h(x_i)} D(i) \cdot \ell^{\alpha} = 0$$

$$-\ell^{-\alpha} \underbrace{\left[\sum_{=} D(i)\right]}_{\text{accuracy}} + \ell^{\alpha} \underbrace{\left[\sum_{\neq} D(i)\right]}_{\text{Error } \mathcal{E}} = 0$$

$$-\frac{1}{\ell^{\alpha}} \cdot (1 - \mathcal{E}) + \ell^{\alpha} \cdot \mathcal{E} = 0$$

$$\frac{1}{\ell^{\alpha}} (1 - \mathcal{E}) = \ell^{\alpha} \mathcal{E} *$$

$$(1 - \mathcal{E}) = \ell^{2\alpha} \mathcal{E}$$

$$\log(1 - \mathcal{E}) = \log\left(\ell^{2\alpha} \mathcal{E}\right)$$

$$\log(1 - \mathcal{E}) = 2\alpha + \log \mathcal{E}$$

$$\alpha = \frac{\log(1-\mathcal{E})}{\log \mathcal{E}} \cdot \frac{1}{2}$$

$$2\alpha = \log(1 - \mathcal{E}) - \log \mathcal{E}$$

$$\alpha = \frac{1}{2}\left(\log(1-\mathcal{E}) - \log \mathcal{E}\right)$$

$$\boxed{\alpha = \frac{1}{2} \cdot \log\left(\frac{1-\mathcal{E}}{\mathcal{E}}\right)}$$