

## Boosting

1

- But:
  - combiner des classifieurs "faibles"  $h^{(1)}, h^{(2)}, \dots \in H$  qui sont à peine meilleurs qu'une décision aléatoire
  - d'une façon itérative
- Exemples de classifieurs de base:
  - petits arbres de décisions
  - arbres d'une feuille: [decision stumps](#)
  - réseaux de neurones avec quelques unités cachées

## Boosting

3

- Démarche
  - entraîner  $h^{(t)}$  sur l'ensemble d'entraînement pondéré par  $\text{BASE}(D_n, \mathbf{w})$
  - calculer  $\alpha^{(t)}$  (le poids de  $h^{(t)}$ )
  - re-pondérer les points (calculer  $w_1, \dots, w_n$ )
- Si  $h^{(t)}$  ne peut pas gérer les points pondérés:
  - re-échantillonner selon la distribution  $w_1, \dots, w_n$

## Boosting

2

- Idée 1: pondération  $w_1, \dots, w_n$  des points d'entraînement
  - attribuer les plus grands poids aux points "difficiles"
  - si  $h^{(t)}(\mathbf{x}_j) = y_j$  alors  $w_j \downarrow$
  - si  $h^{(t)}(\mathbf{x}_j) \neq y_j$  alors  $w_j \uparrow$
- Idée 2: pondération  $\alpha^{(1)}, \dots, \alpha^{(t)}$  des classifieurs de base
  - le poids de  $h^{(t)}$  est monotone décroissant en l'erreur de  $h^{(t)}$

## Boosting

4

```
ADABOOST( $D_n, \text{BASE}(D_n, \mathbf{w}), T$ )
1   $\mathbf{w}^{(1)} \leftarrow (1/n, \dots, 1/n)$     ▷ poids initiaux
2  pour  $t \leftarrow 1$  à  $T$ 
3     $h^{(t)} \leftarrow \text{BASE}(D_n, \mathbf{w})$     ▷ hypothèse de base
4     $\epsilon^{(t)} \leftarrow \sum_{h^{(t)}(\mathbf{x}_i) \neq y_i} w_i$     ▷ erreur pondérée
5    si  $\epsilon^{(t)} \geq 1/2$  alors
6      retourner  $f_{t-1}(\cdot) = \sum_{j=1}^{t-1} \alpha^{(j)} h^{(j)}(\cdot)$ 
7     $\alpha^{(t)} \leftarrow \frac{1}{2} \ln \left( \frac{1-\epsilon^{(t)}}{\epsilon^{(t)}} \right)$     ▷ poids de  $h^{(t)}$ 
8    pour  $i \leftarrow 1$  à  $n$     ▷ re-pondération des points
9      si  $h^{(t)}(\mathbf{x}_i) \neq y_i$  alors
10        $w_i^{(t+1)} \leftarrow \frac{w_i^{(t)}}{2\epsilon^{(t)}}$ 
11     sinon
12        $w_i^{(t+1)} \leftarrow \frac{w_i^{(t)}}{2(1-\epsilon^{(t)})}$ 
13  retourner  $f^{(T)}(\cdot) = \sum_{t=1}^T \alpha^{(t)} h^{(t)}(\cdot)$ 
```

## Boosting

5

- Analyse:

- **normalisation**:  $\sum_{i=1}^n w_i = 1$
- **décorrélation** de  $h^{(t)}$  et l'erreur pondérée:  $\sum_{h^{(t)}(\mathbf{x}_i) \neq y_i} w_i^{(t+1)} = \frac{1}{2}$
- $w_i^{(t+1)} = \frac{w_i^{(t)} e^{-\alpha^{(t)} h^{(t)}(\mathbf{x}_i) y_i}}{2 \sqrt{\varepsilon^{(t)} (1 - \varepsilon^{(t)})}} = \dots = \frac{1}{n \prod_{j=1}^t 2 \sqrt{\varepsilon^{(j)} (1 - \varepsilon^{(j)})}} = \frac{e^{-f(\mathbf{x}_i) y_i}}{\sum_{j=1}^n e^{-f(\mathbf{x}_j) y_j}}$
- **marge** (non-normalisée, fonctionnelle):  $\rho_i = f(\mathbf{x}_i) y_i$
- $\rho_i < 0 \Rightarrow \mathbf{x}_i$  est mal classifié
- $\rho_i$  est **grand**  $\Rightarrow \mathbf{x}_i$  est bien classifié **avec confiance**

## Boosting

7

- Minimisation de la **perte exponentielle** sur la marge:

- étant donné  $h$ :
 
$$\alpha^{(t)} \leftarrow \arg \min_{\alpha} \widehat{R}_e(f^{(t-1)}(\cdot) + \alpha h(\cdot)) = \frac{1}{2} \ln \left( \frac{1 - \varepsilon^{(t)}}{\varepsilon^{(t)}} \right)$$
- si  $\alpha = \frac{1}{2} \ln \left( \frac{1 - \varepsilon^{(t)}}{\varepsilon^{(t)}} \right)$ :
 
$$\widehat{R}_e(f^{(t-1)}(\cdot) + \alpha h(\cdot)) = 2 \sqrt{\varepsilon^{(t)} (1 - \varepsilon^{(t)})} \prod_{j=1}^{t-1} \left[ 2 \sqrt{\varepsilon^{(j)} (1 - \varepsilon^{(j)})} \right]$$
- $h^{(t)}$  doit minimiser  $\varepsilon^{(t)} (1 - \varepsilon^{(t)})$  avec  $\varepsilon^{(t)} < 1/2 \equiv h^{(t)}$  doit minimiser  $\varepsilon^{(t)}$

## Boosting

6

- Minimisation de la **perte exponentielle** sur la marge:

$$L_e((\mathbf{x}, y), f) = e^{-f(\mathbf{x})y} = e^{-\rho}$$

- risque exponentiel empirique:

$$\widehat{R}_e(f) = \frac{1}{n} \sum_{i=1}^n L_e((\mathbf{x}_i, y_i), f) = \frac{1}{n} \sum_{i=1}^n e^{-f(\mathbf{x}_i) y_i}$$

- minimisation **gloutonne**:

$$h^{(t)}(\cdot), \alpha^{(t)} \leftarrow \arg \min_{h(\cdot), \alpha} \widehat{R}_e(f^{(t-1)}(\cdot) + \alpha h(\cdot))$$

## Boosting

8

- Théorème de **convergence**

$$\begin{aligned} \widehat{R}(f) &= \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{f(\mathbf{x}_i) y_i < 0\} \\ &\leq \widehat{R}_e(f) = \frac{1}{n} \sum_{i=1}^n \exp \left( -y_i \sum_{t=1}^T \alpha^{(t)} h^{(t)}(\mathbf{x}_i) \right) \\ &= \prod_{t=1}^T \left[ 2 \sqrt{\varepsilon^{(t)} (1 - \varepsilon^{(t)})} \right] \end{aligned}$$

- soit  $\varepsilon^{(t)} \leq \frac{1}{2} - \delta$ :  $h^{(t)}$  est **un peu** meilleur qu'une décision aléatoire

$$\widehat{R}(f) \leq e^{-2T\delta^2}$$

- $\widehat{R}(f) = 0$  après  $\left\lceil \frac{\ln n}{2\delta^2} \right\rceil + 1$  itérations

## Boosting

9

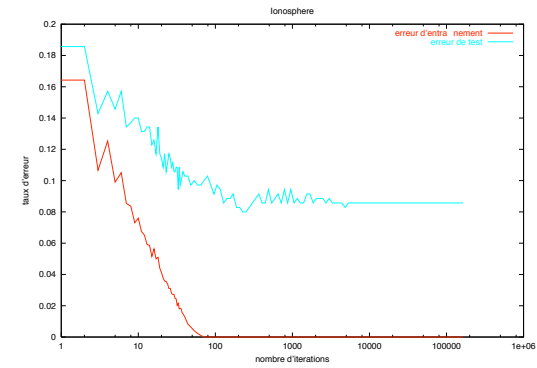
### • Expériences

- **decision stumps**, 4 données de **UCI**, test croisé en **10 blocs**
- l'erreur de test diminue **même après** que l'erreur d'entraînement devienne 0

## Boosting

10

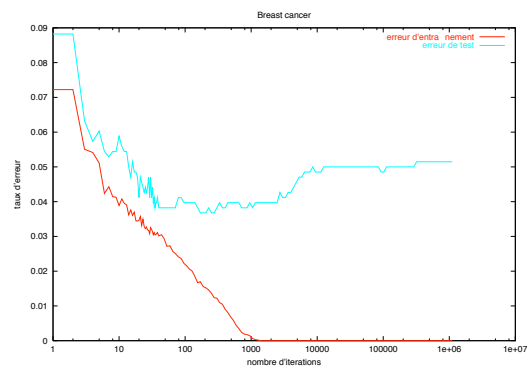
### • Pas beaucoup de sur-apprentissage



## Boosting

11

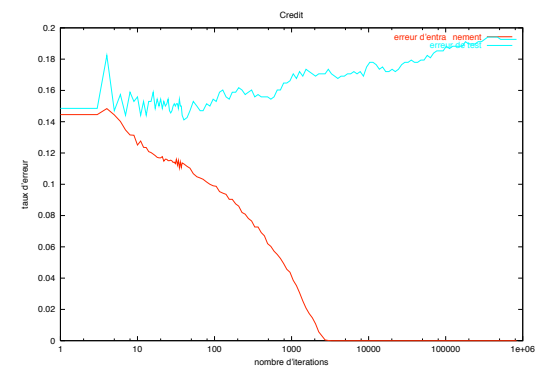
### • Plus de sur-apprentissage



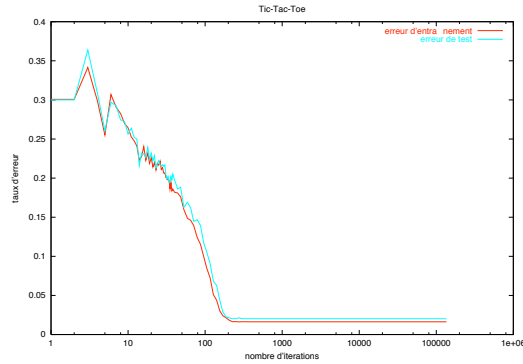
## Boosting

12

### • Beaucoup de sur-apprentissage



- L'erreur d'entraînement asymptotique n'est pas 0



- **Avantage** (edge):  $\gamma(\mathbf{w}) = 1 - 2\varepsilon = \sum_{i=1}^n w_i h(\mathbf{x}_i) y_i$ 
  - juste une transformation linéaire de l'erreur pondérée
  - but de  $h^{(t)}$ : maximiser l'avantage
  - avantage maximal:  $\gamma^*(\mathbf{w}^{(t)}) = \max_{h \in H} \gamma(\mathbf{w}^{(t)})$
  - minimum possible avantage maximal:  $\gamma^* = \min_{\mathbf{w}} \gamma^*(\mathbf{w}) = \min_{\mathbf{w}} \max_{h \in H} \gamma(\mathbf{w})$
  - $\gamma^*(\mathbf{w}^{(t)}) > 2\delta \Rightarrow \widehat{R}(f^{(t)}) = 0$  après  $t = \left\lceil \frac{\ln n}{2\delta^2} \right\rceil + 1$  itérations

- Marge normalisée

- coefficients normalisés:  $\tilde{\alpha}^{(j)} = \frac{\alpha^{(j)}}{\sum_{l=1}^j \alpha^{(l)}}$
- vote normalisé:  $\tilde{f}(\cdot) = \sum_{j=1}^t \tilde{\alpha}^{(j)} h^{(j)}(\cdot) = \frac{\sum_{j=1}^t \alpha^{(j)} h^{(j)}(\cdot)}{\sum_{j=1}^t \alpha^{(j)}} = \frac{f(\cdot)}{\|\alpha\|_1}$
- marge normalisée (géométrique en  $L_1$ ):  $\tilde{\rho}_i(\alpha) = y_i \tilde{f}(\mathbf{x}_i) = \frac{y_i f(\mathbf{x}_i)}{\|\alpha\|_1}$
- marge minimale:  $\tilde{\rho}^*(\alpha) = \min_{i=1, \dots, n} \tilde{\rho}_i(\alpha)$
- maximum possible marge minimale:  $\tilde{\rho}^* = \max_{\alpha} \tilde{\rho}^*(\alpha) = \max_{\alpha} \min_{i=1, \dots, n} \tilde{\rho}_i(\alpha)$

- Théorème min-max de von Neumann [1928]:
 
$$\gamma^* \geq \tilde{\rho}^* \quad (\gamma^* = \tilde{\rho}^* \text{ si } H \text{ est fini})$$
  - théorème principal de la théorie des jeux
  - $\gamma^*(\mathbf{w}^{(t)}) \geq \tilde{\rho}^*$
- Si  $\tilde{\rho}^* > 0$ 
  - $D_n$  est séparable
  - $\tilde{\rho}^*$  est la marge  $L_1$  de séparation
  - $2\delta = \tilde{\rho}^* \Rightarrow \widehat{R}(f^{(t)}) = 0$  après  $t = \left\lceil \frac{2 \ln n}{\tilde{\rho}^{*2}} \right\rceil + 1$  itérations

## Boosting

17

- Comportement **asymptotique**

- séparabilité  $\Rightarrow \|\alpha^{(t)}\|_1 \rightarrow \infty \Rightarrow$  **sur-apprentissage**?

- Théorème de convergence de l'**erreur marginale**

$$\widehat{R}^{(\theta)}(f) = \frac{1}{n} \sum_{i=1}^n \mathbb{I} \left\{ \tilde{f}(\mathbf{x}_i) y_i < \theta \right\} \leq \widehat{R}_\varepsilon^{(\theta)}(f) = 2^T \prod_{t=1}^T \sqrt{\varepsilon^{(t)1-\theta} (1-\varepsilon^{(t)})^{1+\theta}}$$

- $\varepsilon^{(t)} \leq \frac{1}{2} - \theta - \delta \Rightarrow$  convergence en temps fini

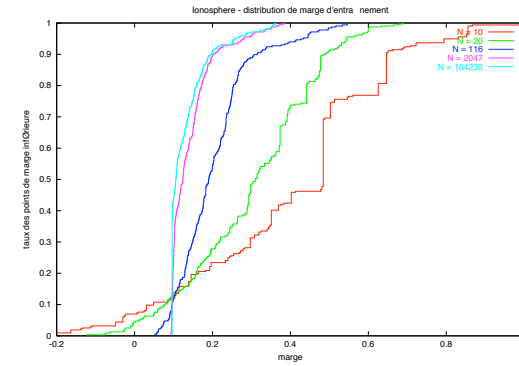
- $\theta < \frac{\tilde{\rho}^*}{2} \Rightarrow$  convergence **asymptotique**

- la marge **minimale asymptotique** est  $\lim_{T \rightarrow \infty} \tilde{\rho}^*(\alpha^{(T)}) \geq \frac{\tilde{\rho}^*}{2}$

## Boosting

18

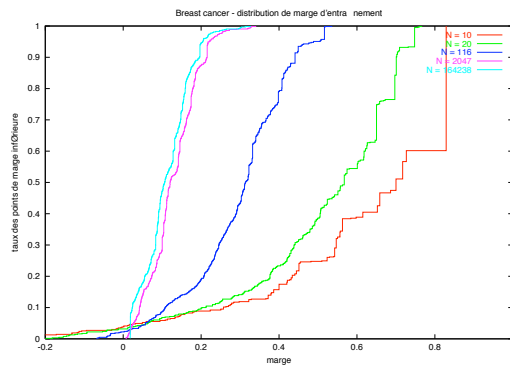
- Distribution de **marge d'entraînement**



## Boosting

19

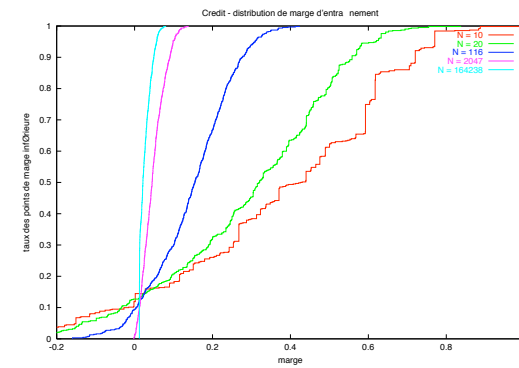
- Distribution de **marge d'entraînement**



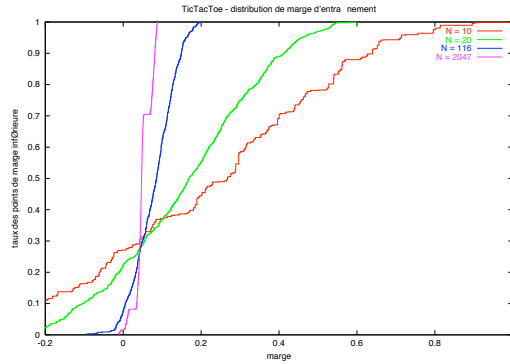
## Boosting

20

- Distribution de **marge d'entraînement**



• Distribution de **marge d'entraînement**



• Maximisation de la marge:  $\lim_{t \rightarrow \infty} \tilde{\rho}^*(\alpha^{(t)}) = \tilde{\rho}^*$ ?

- $h^{(1)}, \dots, h^{(t)}$  donnés: **programmation linéaire**
- boosting **marginal**: minimiser l'**erreur marginale**  $\hat{R}_e^{(\theta)}(f)$
- choisir  $\theta = \theta^{(t)}$  **adaptivement**

```

ADABOOSTθ(Dn, BASE(Dn, w), T, θ)
1  w(1) ← (1/n, ..., 1/n)    ▷ poids initiaux
2  pour t ← 1 à T
3    h(t) ← BASE(Dn, w)    ▷ hypothèse de base
4    ε(t) ← ∑h(t)(xi) ≠ yi wi    ▷ erreur pondérée
5    si ε(t) ≥ 1/2 - θ/2 alors
6      retourner ft-1(·) = ∑j=1t-1 α(j) h(j)(·)
7    α(t) ← 1/2 ln( (1-ε(t)) / ε(t) ) - 1/2 ln( (1+θ) / (1-θ) )    ▷ poids de h(t)
8    pour i ← 1 à n    ▷ re-pondération des points
9      si h(t)(xi) ≠ yi alors
10       wi(t+1) ← wi(t) * (1-θ) / 2ε(t)
11     sinon
12       wi(t+1) ← wi(t) * (1+θ) / 2(1-ε(t))
13  retourner f(T)(·) = ∑t=1T α(t) h(t)(·)
    
```

•  $\theta$  constant

- **weight decay**:  $\alpha^{(t)}$  est plus petit qu'avant
- minimisation de la **perte exponentielle pénalisée** ( $L_1$ ):
 
$$L_e^{(\theta)}((\mathbf{x}, y), f) = e^{-f(\mathbf{x})y + \theta \|\alpha\|_1}$$
- en pratique, la **validation de  $\theta$**  (avec  $T \rightarrow \infty$ ) et la **validation de  $T$**  (avec  $\theta = 0$ ) produisent les résultats **similaires**

## Boosting

25

- $\theta$  adaptatif

- **arc-gv**:  $\theta^{(t)} = \tilde{\rho}^*(\alpha^{(t)})$ : convergence **asymptotique**  $\lim_{t \rightarrow \infty} \tilde{\rho}^*(\alpha^{(t)}) = \tilde{\rho}^*$

- **ADABOOST\***( $\nu$ ):  $\theta^{(t)} = \min_{j=1, \dots, J} \gamma^{(j)} - \nu$ : convergence **en temps fini**:

$$\rho^*(\alpha^{(t)}) > \tilde{\rho}^* - \nu$$

après  $\left\lceil \frac{2 \ln n}{\nu^2} \right\rceil + 1$  itérations

- Observation en pratique

- la maximisation **agressive** de la marge minimale  
**empire la généralisation!!!**
- pas d'explication théorique