

Boosting

1

- Extensions

- autres [fonctions de perte](#)
- fonctions de base [non-binaires](#)
- multiclass, régression

- Applications

- arbres de décisions alternants
- apprentissage des [règles logiques](#)
- traitement d'[images](#)
- traitement de [langue](#)

Boosting

2

- 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

3

- 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

4

- Minimisation d'une [perte arbitraire](#) sur la marge:

$$L((\mathbf{x}, y), f) = L(f(\mathbf{x})y) = L(\rho)$$

- risque empirique:

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

- minimisation [gloutonne](#):

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

- la minimisation analytique n'est pas possible en général

Boosting

5

- Descente de gradient dans l'espace fonctionnel

$$\begin{aligned} h^{(t)}(\cdot) &= \arg \max_h - \frac{\partial \widehat{R}(f^{(t-1)}(\cdot) + \alpha h(\cdot))}{\partial \alpha} \Big|_{\alpha=0} \\ &= \arg \max_h - \frac{1}{n} \sum_{i=1}^n L'(f^{(t-1)}(\mathbf{x}_i) y_i) \cdot h(\mathbf{x}_i) y_i \end{aligned}$$

$$\bullet w_i^{(t)} = \frac{-L'(f^{(t-1)}(\mathbf{x}_i) y_i)}{-\sum_{j=1}^n L'(f^{(t-1)}(\mathbf{x}_j) y_j)} \implies h^{(t)}(\cdot) = \arg \max_h \sum_{i=1}^n w_i^{(t)} h(\mathbf{x}_i) y_i$$

- minimisation scalaire:

$$\alpha^{(t)} = \arg \min_{\alpha} \widehat{R}(f^{(t-1)}(\cdot) + \alpha h^{(t)}(\cdot)) = \arg \min_{\alpha} \sum_{i=1}^n L(f^{(t-1)}(\mathbf{x}_i) y_i + \alpha h^{(t)}(\mathbf{x}_i) y_i)$$

Boosting

6

- Exemples

Algorithme	fonction de perte $L(\rho)$	$\alpha^{(t)}$
ADABOOST	$e^{-\rho}$	line search (analytique)
CONFBOOST	$e^{-\rho}$	line search
ARC-X4	$(1-\rho)^5$	$1/t$
LOGITBOOST	$\ln(1+e^{-\rho})$	Newton-Raphson

• LOGITBOOST

- interprétation statistique de boosting, connection à la régression logistique
- en pratique: plus de tolérance de bruit

Boosting

7

- ADABoost avec abstention

$$\begin{aligned} \bullet h : \mathbb{R}^d &\mapsto \{-1, 0, 1\} \\ \bullet h^{(t)}(\cdot) &= \arg \min_{h: \epsilon_+ > \epsilon_-} (\epsilon_0 + 2\sqrt{\epsilon_+ \epsilon_-}) \\ \bullet \alpha^{(t)} &= \frac{1}{2} \ln \frac{\epsilon_+^{(t)}}{\epsilon_-^{(t)}} \end{aligned}$$

- Problème si $\epsilon_-^{(t)} = 0$

$$\begin{aligned} \bullet \text{solution 1: } \alpha^{(t)} &= \frac{1}{2} \ln \frac{\epsilon_+^{(t)} + \delta}{\epsilon_-^{(t)} + \delta} \\ \bullet \text{solution 2: utilise } \theta &> 0 \text{ dans ADABoost}_{\theta} \end{aligned}$$

Boosting

8

- CONFBOOST

$$\begin{aligned} \bullet h : \mathbb{R}^d &\mapsto \mathbb{R} \\ \bullet h^{(t)}(\cdot) &= \arg \max_h \sum_{i=1}^n w_i^{(t)} h(\mathbf{x}_i) y_i \\ \bullet \text{pas de solution analytique pour } \alpha^{(t)}: &\text{line search} \end{aligned}$$

```

CONFBOOST( $D_n$ , BASE( $D_n, \mathbf{w}$ ),  $T$ )
1    $\mathbf{w}^{(1)} \leftarrow (1/n, \dots, 1/n)$             $\triangleright$  poids initiaux
2   pour  $t \leftarrow 1$  à  $T$ 
3    $h^{(t)} \leftarrow \text{BASE}(D_n, \mathbf{w})$            $\triangleright \leftarrow \arg \max_h \sum_{i=1}^n w_i^{(t)} h(\mathbf{x}_i) y_i$ 
4    $\alpha^{(t)} \leftarrow \arg \min_{\alpha} \sum_{i=1}^n w_i^{(t)} e^{-\alpha h^{(t)}(\mathbf{x}_i) y_i}$        $\triangleright$  poids de  $h^{(t)}$ 
5   si  $\alpha^{(t)} = \infty$  alors     $\triangleright$  "overfitting":  $h(\mathbf{x}_i) y_i > 0 \forall i = 1, \dots, n$ 
6   retourner  $h^{(t)}(\cdot)$ 
7   si  $\alpha^{(t)} < 0$  alors     $\triangleright$  "underfitting":  $\gamma^{(t)} = \sum_{i=1}^n w_i^{(t)} h^{(t)}(\mathbf{x}_i) y_i < 0$ 
8   retourner  $f^{(t-1)}(\cdot) = \sum_{j=1}^{t-1} \alpha^{(j)} h^{(j)}(\cdot)$ 
9   pour  $i \leftarrow 1$  à  $n$             $\triangleright$  re-péndération des points
10   $w_i^{(t+1)} \leftarrow w_i^{(t)} \frac{e^{-\alpha^{(t)} h^{(t)}(\mathbf{x}_i) y_i}}{\sum_{j=1}^n w_j^{(t)} e^{-\alpha^{(t)} h^{(t)}(\mathbf{x}_j) y_j}}$ 
11  retourner  $f^{(T)}(\cdot) = \sum_{t=1}^T \alpha^{(t)} h^{(t)}(\cdot)$ 

```

9

Boosting

10

- Notations

- avantages (marges) de base: $\gamma_i = h(\mathbf{x}_i) y_i$

- avantage: $\gamma = \sum_{i=1}^n w_i \gamma_i$

- paramètre de relaxation/régularisation: θ

- $E^{(t)}(\alpha) = e^{\theta \alpha} \sum_{i=1}^n w_i^{(t)} e^{-\alpha \gamma_i^{(t)}} = \sum_{i=1}^n w_i^{(t)} e^{-\alpha (\gamma_i^{(t)} - \theta)}$

- Théorème de convergence de l'erreur marginale

$$\widehat{R}^{(\theta)}(f^{(T)}) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\left\{ \tilde{f}^{(T)}(\mathbf{x}_i) y_i < \theta \right\} \leq \widehat{R}_e^{(\theta)}(f^{(T)}) = \prod_{t=1}^T E^{(t)}(\alpha^{(t)})$$

```

CONFBOOSTθ( $D_n$ , BASE( $D_n, \mathbf{w}$ ),  $T, \theta$ )
1    $\mathbf{w}^{(1)} \leftarrow (1/n, \dots, 1/n)$             $\triangleright$  poids initiaux
2   pour  $t \leftarrow 1$  à  $T$ 
3    $h^{(t)} \leftarrow \text{BASE}(D_n, \mathbf{w})$            $\triangleright \leftarrow \arg \max_h \sum_{i=1}^n w_i^{(t)} h(\mathbf{x}_i) y_i$ 
4   pour  $i \leftarrow 1$  à  $n$             $\triangleright$  avantages de base
5    $\gamma_i^{(t)} \leftarrow h^{(t)}(\mathbf{x}_i) y_i$ 
6    $\alpha^{(t)} \leftarrow \arg \min_{\alpha} e^{\theta \alpha} \sum_{i=1}^n w_i^{(t)} e^{-\alpha \gamma_i^{(t)}} = \arg \min_{\alpha} E^{(t)}(\alpha)$        $\triangleright$  poids de  $h^{(t)}$ 
7   si  $\alpha^{(t)} = \infty$  alors     $\triangleright$  "overfitting":  $\gamma_i^{(t)} > \theta \forall i = 1, \dots, n$ 
8   retourner  $h^{(t)}(\cdot)$ 
9   si  $\alpha^{(t)} < 0$  alors     $\triangleright$  "underfitting":  $\gamma^{(t)} < \theta$ 
10  retourner  $f^{(t-1)}(\cdot) = \sum_{j=1}^{t-1} \alpha^{(j)} h^{(j)}(\cdot)$ 
11  pour  $i \leftarrow 1$  à  $n$             $\triangleright$  re-péndération des points
12   $w_i^{(t+1)} \leftarrow w_i^{(t)} \frac{e^{-\alpha^{(t)} \gamma_i^{(t)}}}{\sum_{j=1}^n w_j^{(t)} e^{-\alpha^{(t)} \gamma_j^{(t)}}}$ 
13  retourner  $f^{(T)}(\cdot) = \sum_{t=1}^T \alpha^{(t)} h^{(t)}(\cdot)$ 

```

11

Boosting

12

- Multiclasse (K classes)

- étiquettes K -dimensionnelles

- en général, une observation peut appartenir à plusieurs classes (classification "multiétiquettes")

- cas spécial: "uniétiquettes"

$$y_\ell = \begin{cases} 1 & \text{si } \ell \text{ est la classe correcte } \ell(\mathbf{x}), \\ -1 & \text{sinon} \end{cases}$$

Boosting

13

- Distribution des poids $w_{i,\ell}$

- $w_{i,\ell}$ exprime les difficultés de classifier \mathbf{x}_i à sa classe correcte ($\ell = \ell(\mathbf{x}_i)$) ou contre ses classes incorrectes ($\ell \neq \ell(\mathbf{x}_i)$)

- peuvent être initialisés à

$$w_{i,\ell}^{(1)} = \begin{cases} \frac{1}{2n} & \text{si } \ell \text{ est la classe correcte (si } y_{i,\ell} = 1\text{),} \\ \frac{1}{2n(K-1)} & \text{sinon (si } y_{i,\ell} = -1\text{)} \end{cases}$$

Boosting

14

- Classifieurs faibles:

$$\mathbf{h}: X \rightarrow \mathbb{R}^K; h_\ell: X \rightarrow \mathbb{R}$$

- peuvent voter aux plusieurs classes

- Classifieur final:

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

- la meilleure classe (cas “uniétiquettes”):

$$\hat{\ell}(\mathbf{x}) = \arg \max_{\ell} f_\ell^{(T)}(\mathbf{x})$$

Boosting

15

- Avantages de base

- nombre finis de classifieurs faibles $\mathbf{h}_1, \dots, \mathbf{h}_N$

$$\gamma_{j,i,\ell} = h_{j,\ell}(\mathbf{x}_i) y_{i,\ell}$$

- par exemple, si $\mathbf{h}: X \rightarrow \{-1, 0, +1\}^K$:

$$\gamma_{j,i,\ell} = \begin{cases} 1 & \text{si } h_{j,\ell}(\mathbf{x}_i) = y_{i,\ell} \\ -1 & \text{si } h_{j,\ell}(\mathbf{x}_i) \neq y_{i,\ell} \\ 0 & \text{si } h_{j,\ell}(\mathbf{x}_i) = 0 \end{cases}$$

Boosting

16

- Fonction d'énergie:

$$E^{(t)}(\mathbf{h}_j, \alpha) = e^{\theta\alpha} \sum_{i=1}^n \sum_{\ell=1}^K w_{i,\ell}^{(t)} \exp(-\alpha \gamma_{j,i,\ell})$$

- dans chaque itération, sélectionne $\mathbf{h}^{(t)} = \mathbf{h}_{j^{(t)}}$ qui maximise

$$\gamma_j^{(t)} = \sum_{i=1}^n \sum_{\ell=1}^K w_{i,\ell}^{(t)} \gamma_{j,i,\ell} = \theta + \left. \frac{-\partial E^{(t)}(\mathbf{h}_j, \alpha)}{\partial \alpha} \right|_{\alpha=0}$$

- coefficient de $\mathbf{h}^{(t)}$:

$$\alpha^{(t)} = \arg \min_{\alpha} E^{(t)}(\mathbf{h}^{(t)}, \alpha)$$

Boosting

17

- Mise à jour des poids:

$$w_{i,\ell}^{(t+1)} \leftarrow w_{i,\ell}^{(t)} \frac{\exp\left(-\alpha^{(t)} \gamma_{j^{(t)}, i, \ell}\right)}{Z^{(t)}},$$

ou

$$Z^{(t)} = \sum_{\ell=1}^n \sum_{i=1}^K w_{i,\ell}^{(t)} \exp\left(-\alpha^{(t)} \gamma_{j^{(t)}, i, \ell}\right)$$

Boosting

18

- Marges par classe:

$$\rho_{i,\ell}^{(t)} = \sum_{\tau=1}^t \alpha^{(\tau)} \gamma_{j^{(\tau)}, i, \ell} = f_\ell^{(t)}(\mathbf{x}_i) y_{i,\ell}$$

- connection entre les poids et les marges

$$w_{i,\ell}^{(t+1)} = w_{i,\ell}^{(1)} \prod_{\tau=1}^t \frac{\exp\left(-\alpha^{(\tau)} \gamma_{j^{(\tau)}, i, \ell}\right)}{Z^{(\tau)}} = w_{i,\ell}^{(1)} \frac{\exp\left(-\rho_{i,\ell}^{(t)}\right)}{\prod_{\tau=1}^t Z^{(\tau)}}$$

- marges normalisées:

$$\tilde{\rho}_{i,\ell}^{(t)} = \sum_{\tau=1}^t \tilde{\alpha}^{(\tau)} \gamma_{j^{(\tau)}, i, \ell} = \tilde{f}_\ell^{(t)}(\mathbf{x}_i) y_{i,\ell}.$$

Boosting

19

- Erreur marginale:

$$\hat{R}_n^{(\theta)}(\mathbf{f}^{(T)}, \mathbf{w}^{(1)}) = \sum_{i=1}^n \sum_{\ell=1}^K w_{i,\ell}^{(1)} \mathbb{I}\left\{ \tilde{\rho}_{i,\ell}^{(T)} \leq \theta \right\}$$

- si $\forall t : E^{(t)}(\mathbf{h}^{(t)}, \alpha^{(t)}) < 1 - \delta$, alors $\hat{R}_n^{(\theta)}(\mathbf{f}^{(T)}, \mathbf{w}^{(1)}) = 0$ après

$$T^* = \left\lceil \left(\delta \ln \left(\min_{i,\ell : w_{i,\ell}^{(1)} \neq 0} w_{i,\ell}^{(1)} \right) \right)^{-1} \right\rceil + 1$$

- en particulier, si $w_{i,\ell}^{(1)}$ sont initialisés comme ci-dessus:

$$T^* = \left\lceil \frac{\ln(2n(K-1))}{\delta} \right\rceil + 1.$$

Boosting

20

- Erreur de classification:

$$\hat{R}(\mathbf{f}^{(T)}) = \sum_{i=1}^n \mathbb{I}\left\{ \hat{\ell}(\mathbf{x}_i) \neq \ell(\mathbf{x}_i) \right\}$$

- $\hat{R}_n^{(\theta)}(\mathbf{f}^{(T)}, \mathbf{w}^{(1)}) = 0 \implies \hat{R}(\mathbf{f}^{(T)}) = 0$:

$$\forall \ell \neq \ell(\mathbf{x}_i) : \tilde{f}_{i,\ell}^{(T)}(\mathbf{x}_i) = -\tilde{\rho}_{i,\ell}^{(T)} < -\theta \leq 0 \leq \theta < \tilde{\rho}_{i,\ell}^{(T)} = \tilde{f}_{i,\ell}^{(T)}(\mathbf{x}_i).$$

- $\hat{R}(\mathbf{f}^{(T)}) = 0 \not\implies \hat{R}_n^{(\theta)}(\mathbf{f}^{(T)}, \mathbf{w}^{(1)}) = 0$

Boosting

21

- Classifieurs faibles scalaires + vecteurs de votes

$$\mathbf{h} = \varphi \mathbf{v}$$

- classifieurs faibles $\varphi : X \rightarrow \{-1, 1\}$
- vecteurs de votes: $\mathbf{v} : X \rightarrow \{-1, 0, 1\}^K$

- Stratégies

- “one-against-all”: $\mathbf{v} = \{-1, -1, \dots, -1, 1, -1, \dots, -1, -1\}$
- “one-against-one”: $\mathbf{v} = \{0, 0, \dots, 0, -1, 0, \dots, 0, 1, 0, \dots, 0, 0\}$
- ECOC
- optimiser $\mathbf{v}^{(t)}$ (?)

Boosting

22

- Arbres de décisions alternants

- noeuds de prédiction et noeuds de découpage
- chaque noeud de prédiction contient un nombre $\in [-1, 1]$
- la prédiction finale est la signe de la somme des prédictions au long des chemins valides
- utiliser AdaBoost avec abstention

Boosting

23

- Apprentissage des règles logiques: SLIPPER

- classifieurs de base: des conjonctions positives et une règle de défaut négative
- construire une conjonction d'une façon gloutonne en maximisant $\sqrt{\epsilon_+} - \sqrt{\epsilon_-}$ sur un ensemble GROWSET
- élaguer les règles sur un autre ensemble PRUNESSET
- utiliser AdaBoost avec abstention

Boosting

24

- Traitement d'images: détection d'objets

- classifieurs de base: des masques rectangulaires simples
- l'astuce de l'image intégral